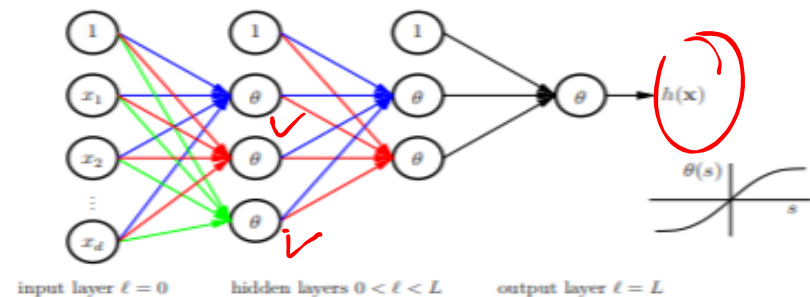# Machine Learning from Data

Lecture 21: Spring 2021
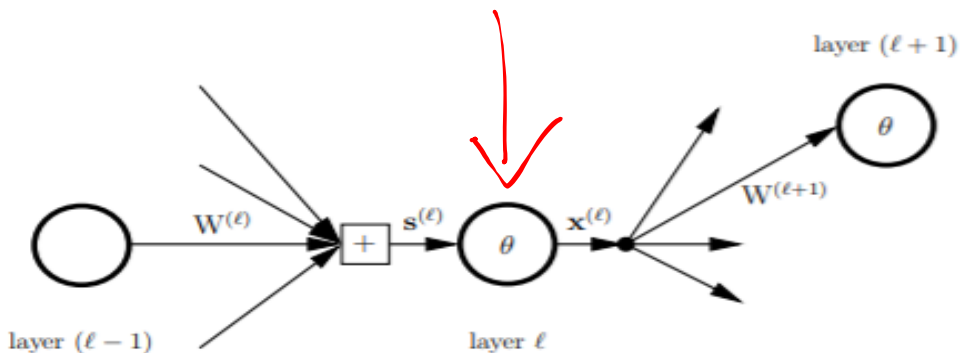
# Today's Lecture

- Neural Networks
    - Forward Propagation ✓
    - Backward Propagation
    - Overfitting

# Recap



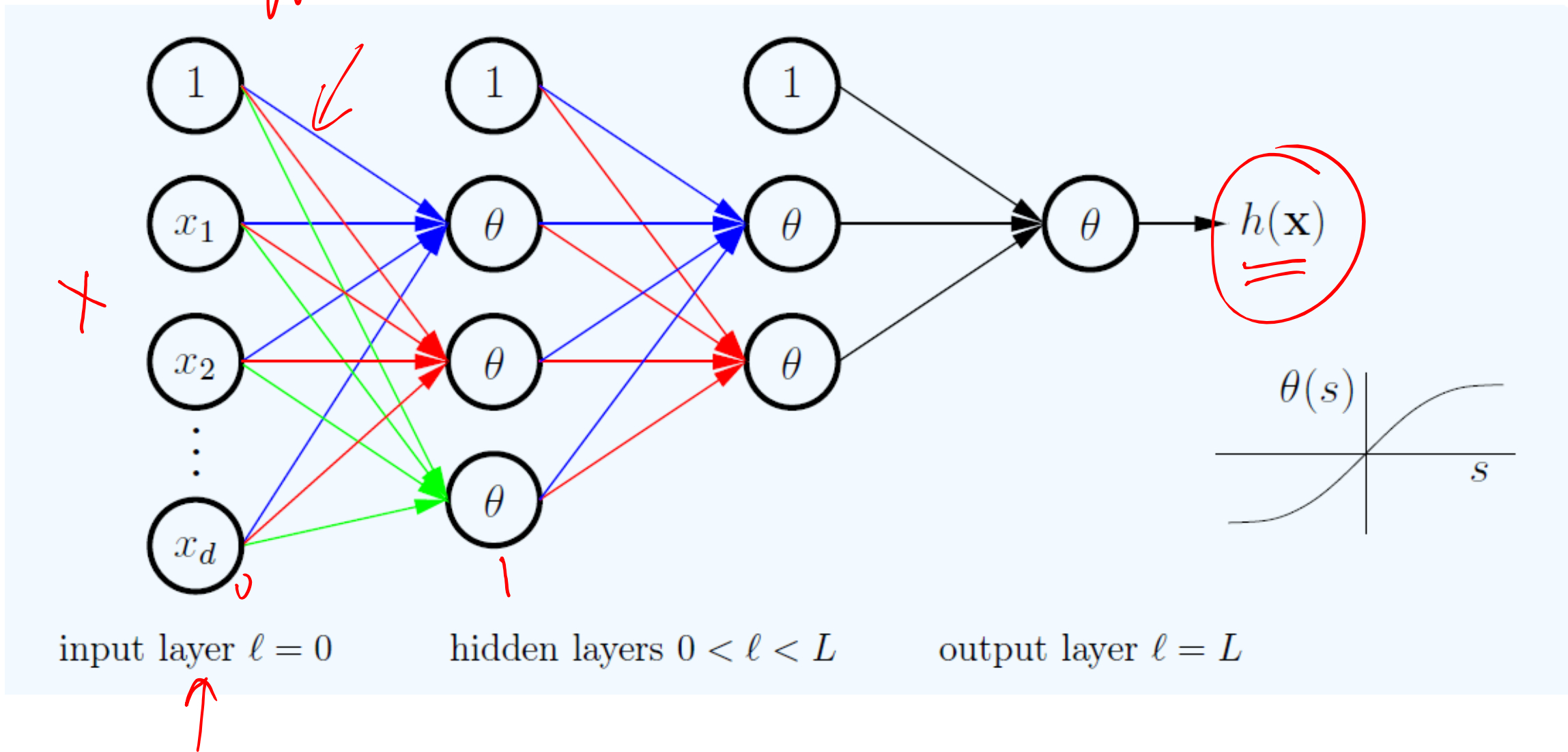input layer $\ell = 0$  hidden layers $0 < \ell < L$  output layer $\ell = L$

NN

layer $(\ell + 1)$

$W^{(\ell+1)}$

$W^{(\ell)}$  $s^{(\ell)}$  $\theta$  $x^{(\ell)}$

layer $(\ell - 1)$  layer $\ell$

## layer $\ell$ parameters

| | | |
|---|---|---|
| signals in | $s^{(\ell)}$ | $d^{(\ell)}$ dimensional input vector |
| outputs | $x^{(\ell)}$ | $d^{(\ell)} + 1$ dimensional output vector |
| weights in | $W^{(\ell)}$ | $(d^{(\ell-1)} + 1) \times d^{(\ell)}$ dimensional matrix |
| weights out | $W^{(\ell+1)}$ | $(d^{(\ell)} + 1) \times d^{(\ell+1)}$ dimensional matrix |

layers $\ell = 0, 1, 2, \ldots, L$
layer $\ell$ has "dimension" $d^{(\ell)} \implies d^{(\ell)} + 1$ nodes

$$W^{(\ell)} = \begin{bmatrix} \mathbf{w}_1^{(\ell)} & \mathbf{w}_2^{(\ell)} & \cdots & \mathbf{w}_{d^{(\ell)}}^{(\ell)} \\ | & | & \vdots & | \end{bmatrix}$$

$W^{(1)}$  $W^{(2)}$  $\cdots$  $W^{(L)}$

$W^{(1)}$  $W^{(2)}$  $W^{(3)}$

$\mathbf{x}$

$1$

$x_1$

$x_2$

$\vdots$

$x_d$

$0$

$1$

$\theta$

$\theta$

$\theta$

$1$

$1$

$\theta$

$\theta$

$\theta$

$h(\mathbf{x})$

$\theta(s)$

$s$

input layer $\ell = 0$     hidden layers $0 < \ell < L$     output layer $\ell = L$

$$\ell-1 \qquad w^{(\ell)} \qquad \ell$$

$$w_i^{(\ell)}$$

$$x_1$$

$$x_2 \qquad s_i^{(\ell)} \qquad x_i^{(\ell)}$$

$$x_d^{(\ell-1)}$$

$$d^{(\ell)}$$

$$h(x)$$

$$S^{(\ell)} = \begin{bmatrix} S_1^{(\ell)} \\ S_2^{(\ell)} \\ \vdots \\ S_{d^{(\ell)}}^{(\ell)} \end{bmatrix}$$

$$x^{(\ell)} = \begin{bmatrix} 1 \\ x_1^{(\ell)} \\ x_2^{(\ell)} \\ \vdots \\ x_d^{(\ell)} \end{bmatrix} d^{\ell}+1$$

$$W^{(l)} = \left[ W_1^{(l)} \quad W_2^{(l)} \quad \cdots \quad W_d^{(l)} \right] \Big\} \, d^{(l-1)} + 1$$

perception

$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{d\text{-dimensional}}$

# Forward Propogation Algorithm

$$W = \left\{ W^{(1)}, \; W^{(2)}, \; W^{(3)} \right\}$$

$\underset{0,1}{\nwarrow} \qquad \underset{1,2}{\nwarrow} \qquad \underset{2,3.}{\nwarrow}$

$$\underline{\Theta}_i^l \longrightarrow \quad S_i^{(\ell)} = \left(W_i^{(\ell)}\right)^T x^{(\ell-1)}$$

weight of perception

$$\underbrace{\left(W_i^{(\ell)}\right)^T x^{(\ell-1)}}_{\text{linear signal}}$$

Signal

$$S^l = \left(W^{(\ell)}\right)^T x^{(\ell-1)}$$

$$S^l = \begin{bmatrix} S_1^{(\ell)} \\ S_2^{(\ell)} \\ \vdots \\ S_d^{(\ell)} \end{bmatrix} = \begin{bmatrix} - (W_1^{(\ell)})^T - \\ - (W_2^{(\ell)})^T - \\ \\ - (W_d^{(\ell)})^T - \end{bmatrix} \quad x^{(\ell-1)}$$
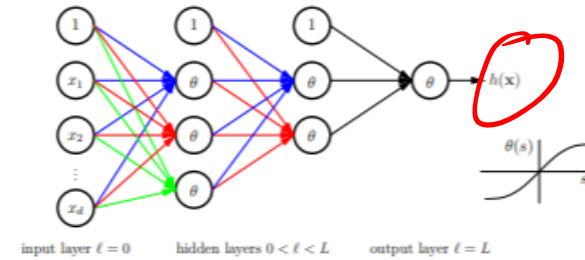
$$X^l = \begin{bmatrix} 1 \\ \Theta(S^l) \end{bmatrix}$$

# The Linear Signal

Input $\mathbf{s}^{(\ell)}$ is a linear combination (using weights) of the outputs of the previous layer $\mathbf{x}^{(\ell-1)}$.

$$\mathbf{s}^{(\ell)} = (\mathbf{W}^{(\ell)})^{\mathsf{T}} \mathbf{x}^{(\ell-1)}$$



input layer $\ell = 0$     hidden layers $0 < \ell < L$     output layer $\ell = L$

$$
\begin{bmatrix} s_1^{(\ell)} \\ s_2^{(\ell)} \\ \vdots \\ s_j^{(\ell)} \\ \vdots \\ s_{d^{(\ell)}}^{(\ell)} \end{bmatrix}
=
\begin{bmatrix} (\mathbf{w}_1^{(\ell)})^{\mathsf{T}} \rule{2cm}{0.4pt} \\ (\mathbf{w}_2^{(\ell)})^{\mathsf{T}} \rule{2cm}{0.4pt} \\ \vdots \\ (\mathbf{w}_j^{(\ell)})^{\mathsf{T}} \rule{2cm}{0.4pt} \\ \vdots \\ (\mathbf{w}_{d^{(\ell)}}^{(\ell)})^{\mathsf{T}} \rule{2cm}{0.4pt} \end{bmatrix}
\mathbf{x}^{(\ell-1)}
$$

$$s_j^{(\ell)} = (\mathbf{w}_j^{(\ell)})^{\mathsf{T}} \mathbf{x}^{(\ell-1)}$$

(recall the linear signal $s = \mathbf{w}^{\mathsf{T}}\mathbf{x}$)

$$\mathbf{s}^{(\ell)} \xrightarrow{\ \theta\ } \mathbf{x}^{(\ell)}$$

## Forward prop. algo.

$W = \{ W^{(1)}, W^{(2)} \cdots W^{(L)} \}$

1) Initialize $x^{(0)} = x$, $l = 0$  $h(x) \to$ goal

$$x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

2) Propogate, set $l = l+1$

$$\underline{\underline{S^l}} = \left( W^{(l)} \right)^T x^{(l-1)}$$

$$x^l = \begin{bmatrix} 1 \\ \theta(S^l) \end{bmatrix}$$  Repeat

3)  $x^{(L)} \longrightarrow$  $h(x) = \underline{\underline{x_1^{(L)}}}$

# Forward Propagation: Computing $h(\mathbf{x})$

$$\mathbf{x} = \mathbf{x}^{(0)} \xrightarrow{W^{(1)}} \mathbf{s}^{(1)} \xrightarrow{\theta} \mathbf{x}^{(1)} \xrightarrow{W^{(2)}} \mathbf{s}^{(2)} \xrightarrow{\theta} \mathbf{x}^{(2)} \cdots \xrightarrow{W^{(L)}} \mathbf{s}^{(L)} \xrightarrow{\theta} \mathbf{x}^{(L)} = h(\mathbf{x}).$$

---

**Forward propagation to compute $h(\mathbf{x})$:**

1: $\mathbf{x}^{(0)} \leftarrow \mathbf{x}$          [Initialization]

2: **for** $\ell = 1$ to $L$ **do**      [Forward Propagation]

3:     $\mathbf{s}^{(\ell)} \leftarrow (W^{(\ell)})^{\mathsf{T}} \mathbf{x}^{(\ell-1)}$

4:     $\mathbf{x}^{(\ell)} \leftarrow \begin{bmatrix} 1 \\ \theta(\mathbf{s}^{(\ell)}) \end{bmatrix}$

5: **end for**

6: $h(\mathbf{x}) = \mathbf{x}^{(L)}$          [Output]

## Data

$x_1 \quad y_1 \xrightarrow{\text{NN}} h(x_1)$

$x_2 \quad y_2 \xrightarrow{\text{NN}} h(x_2)$

$x_n \quad y_n \xrightarrow{\text{NN}} h(x_n)$

$W^{(1)} \quad W^{(2)} \quad W^{(3)} \dots W^{(L)}$

$$E_{in}(h) = \frac{1}{N} \sum_{i=1}^{N} \left( h(x_n) - y_n \right)^2$$

$$E_{in}(h) = \frac{1}{N} \sum_{i=1}^{N} e\left( h(x_n), y_n \right)$$

Know $\longrightarrow$ W's $\longrightarrow$ find weights that minimize $E_{in}$

$$W^{(l)} \longleftarrow W^{(l)} - \alpha \frac{\partial E_{in}}{\partial W^{(l)}}$$

learning rate

$$\frac{\partial E_{in}}{\partial W^{(l)}} = \frac{1}{N} \sum_{n=1}^{N} \frac{\partial e(h(x_n), y_n)}{\partial W^{(l)}}$$

Summary $\longrightarrow$ Compute Output X
Getting weights $\iff$ fitting data.

## In practice

$$|W| = W^1 + W^2 + \dots W^L$$

$$|V| = \text{No. of nodes} \;(\theta) \longrightarrow \tanh$$

$E_{in} \longrightarrow$ How many computations?

$$N|W| + N|V|$$

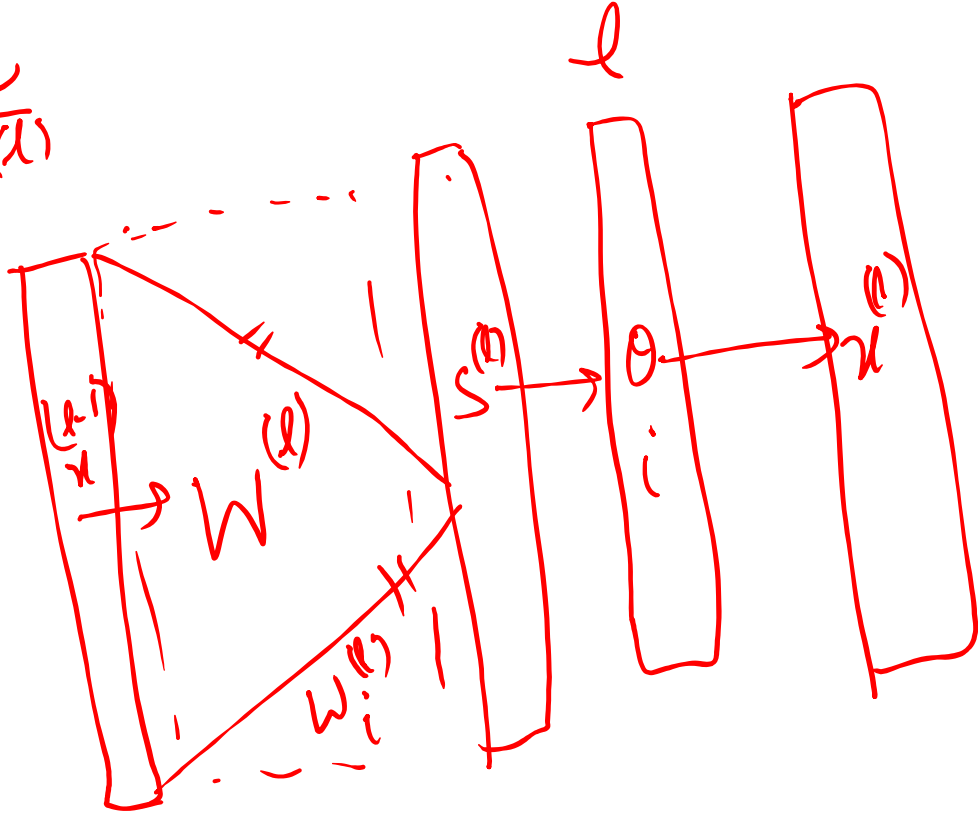$$w_{ij}^{(l)} \longrightarrow \frac{\partial E_{in}}{\partial w_{ij}^{(l)}} \simeq \frac{E_{in}\left(w_{ij}^{(l)} + \triangle\right) - E_{in}\left(w_{ij}^{(l)} - \triangle\right)}{2\triangle}$$

$\triangle \longrightarrow$ tolerance

$$O\left(\left(N|W| + N|V|\right)|W|\right)$$

# BACKPROPOGATION ALGORITHM

$$\frac{\partial e}{\partial W^{(l)}} \Rightarrow$$

$\ell$



$x^{(l-1)}$ $\xrightarrow{W^{(l)}}$ $s^{(l)}$ $\rightarrow$ $\theta$ $\rightarrow$ $x^{(l)}$

$w_i^{(l)}$

$$\frac{\partial e}{\partial W_i^{(l)}}$$

Chain Rule

$$\frac{\partial e}{\partial W_i^{(l)}} = \frac{\partial S_i^{(l)}}{\partial W_i^{(l)}} \times \frac{\partial e}{\partial S_i^{(l)}}$$

$$S_i^{(l)} = \left(W_i^{(l)}\right)^T x^{(l-1)}$$

$$\frac{\partial S_i^{(l)}}{\partial W_i^{(l)}} = x^{(l-1)}$$

$$\therefore \frac{\partial e}{\partial W_i^{(l)}} = x^{(l-1)} \times \frac{\partial e}{\partial S_i^{(l)}}$$

$$\rightarrow \frac{\partial e}{\partial W_i^{(l)}}$$

$$\frac{\partial e}{\partial W^{(l)}} = \left[ x^{(l-1)} \frac{\partial e}{\partial S_1^{(l)}} \quad x^{(l-1)} \frac{\partial e}{\partial S_2^{(l)}} \quad \cdots \cdots \quad x^{(l-1)} \frac{\partial e}{\partial S_{d^{(l)}}^{(l)}} \right]$$

$$\frac{\partial e}{\partial W^{(l)}} = x^{(l-1)} \underset{\substack{\nearrow \\ \text{previous} \\ \text{layer}}}{} \left( \frac{\partial e}{\partial S^{(l)}} \right)^T \longrightarrow \left(\delta^{(l)}\right)^T \longrightarrow \text{sensitivity}$$

$$\underset{\text{vector of derivatives}}{\longleftarrow}$$

$$\boxed{\frac{\partial e}{\partial S^{(l)}}} \qquad \therefore \quad \frac{\partial e}{\partial S_i^{(l)}} = \frac{\partial e}{\partial x_i^{(l)}} \times \frac{\partial x_i^{(l)}}{\partial S_i^{(l)}} \Bigg] \quad x_i^{(l)} = \theta\left(S_i^{(l)}\right)$$

$$\therefore \quad \frac{\partial e}{\partial S_i^{(l)}} = \theta'\left(S_i^{(l)}\right) \cdot \frac{\partial e}{\partial x_i^{(l)}}$$
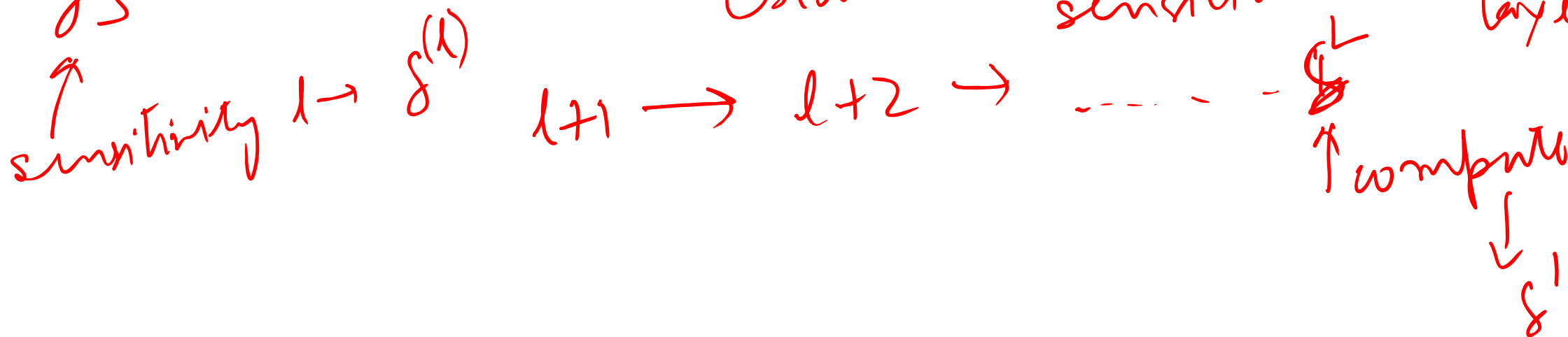
$$\frac{\partial e}{\partial S^{(l)}} = \theta'\left(S^{(l)}\right) \otimes \left[\frac{\partial e}{\partial x^{(l)}}\right]_1^{d^{(l)}} \leftarrow$$

$$\frac{\partial e}{\partial x^{(l)}} = \sum_{k=1}^{d^{(l+1)}} \frac{\partial e}{\partial S_k^{(l+1)}} \cdot \frac{\partial S_k^{(l+1)}}{\partial x^{(l)}}$$

$$\therefore \quad S_k^{(l+1)} = \left(W_k^{l+1}\right)^T x^{(l)}$$

$$\frac{\partial e}{\partial x^{(l)}} = \sum_{k=1}^{d^{(l+1)}} W_k^{(l+1)} \frac{\partial e}{\partial S^{(l+1)}}$$

$$\frac{\partial e}{\partial x^{(\ell)}} = W^{(\ell+1)} \underline{\delta^{(\ell+1)}} \longleftarrow \text{sensitivity}$$

$$\frac{\partial e}{\partial S^{(\ell)}} = \theta'(S^{(\ell)}) \otimes \left[ W^{(\ell+1)} \, \delta^{(\ell+1)} \right]^{d^{(\ell)}}$$

Established a relation between sensitivities across layers.

$\uparrow$ sensitivity $\ell \longrightarrow \delta^{(\ell)}$      $\ell+1 \longrightarrow \ell+2 \longrightarrow \cdots - \overset{L}{\underset{\downarrow}{\delta}}$

$\uparrow$ compute

$\downarrow$
$\delta^{\prime}$

1) $\dfrac{\partial e}{\partial W^{(\ell)}} = x^{(\ell-1)} (\delta^{\ell})^{T}$ ✓

$(1+d^{(\ell-1)}) \times d^{\ell}$

2) $\delta^{(\ell)} = \theta'(S^{(\ell)}) \otimes [W^{(\ell+1)} \delta^{(\ell+1)}]_{1}^{d^{(\ell)}}$ ✓

$$\underline{\text{STOCHASTIC GRADIENT DESCENT}}$$

1) Pick $x_n$

2) $x_n = x^{(0)} \xrightarrow{W^{(1)}} S^{(1)} \xrightarrow{\theta} x^{(1)} \longrightarrow S^{(2)} \xrightarrow{\theta} x^{(2)} \cdots \cdots$

$\cdots \longrightarrow x^{(L)} = h(x_n)$

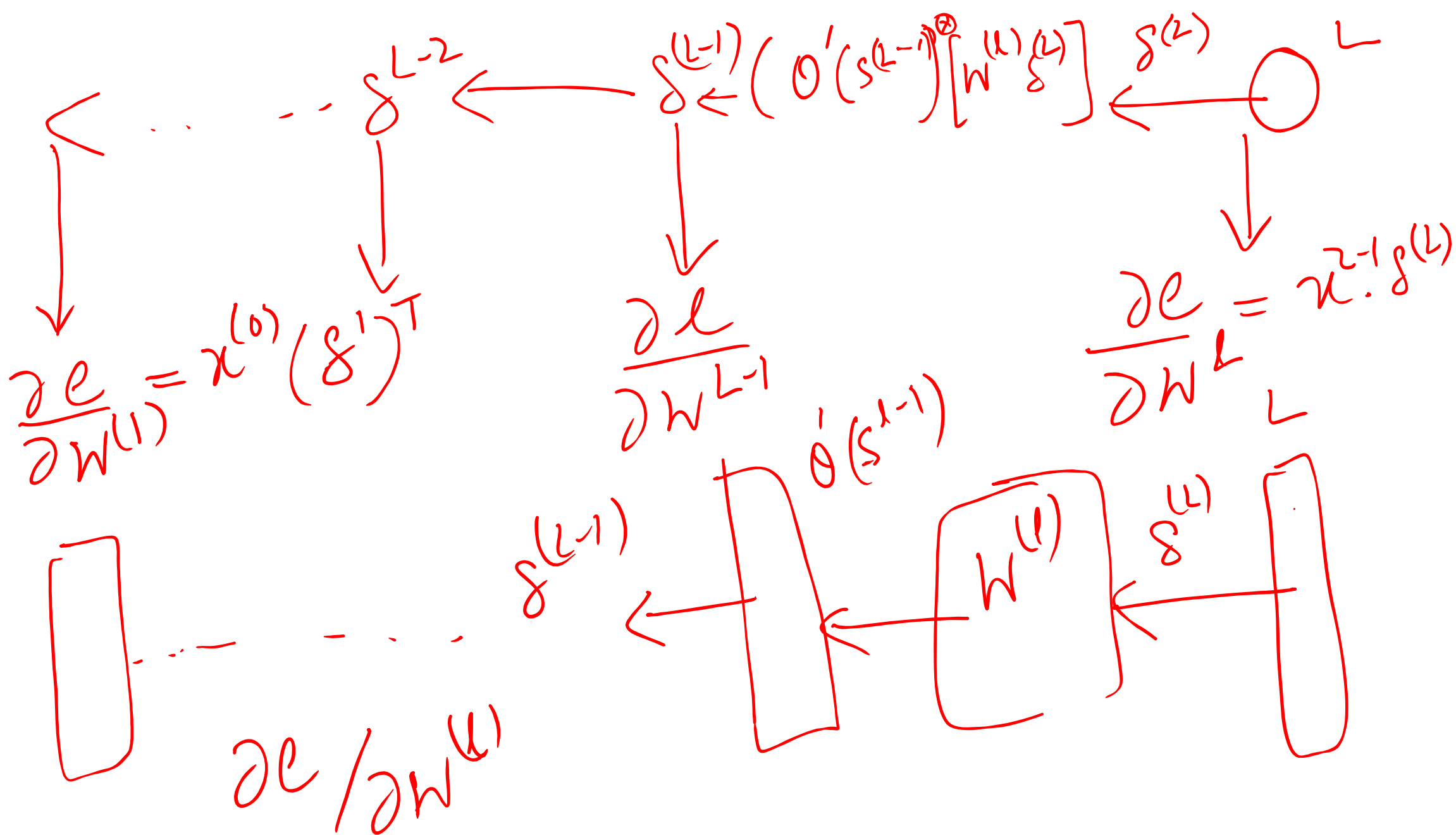$e(h(x_n), y_n)$

3) Run BACKPROPOGATION

$$\delta^{(L)} = \frac{\partial e}{\partial S^{(L)}} \quad , \quad e = (x^{L} - y_n)^2 \qquad x_n y_n$$

$$\delta^{(L)} = \frac{\partial (x^{(L)} - y_n)^2}{\partial S^{(L)}} = 2(x^{(L)} - y_n) \cdot \frac{\partial x^{(L)}}{\partial S^{(L)}}$$

$$= 2(x^{(L)} - y_n) \cdot \theta'(S^{(L)})$$

$$\delta^{(L)} = 2(x^{(L)} - y_n)\left(1 - \tanh^2(S^{(L)})\right)$$

$$= 2(x^{(L)} - y_n)\left(1 - (x^{(L)})^2\right)$$

$$\frac{\partial e}{\partial (W^{(l)})} = x^{(L-1)} \cdot \delta^{(L)}$$

$$\delta^{L-2} \longleftarrow \delta^{(L-1)} \left( \theta'(S^{(L-1)})^{\otimes} [W^{(L)} \delta^{(L)}] \right) \overset{\delta^{(L)}}{\longleftarrow} \bigcirc \; L$$

$$\frac{\partial e}{\partial W^{(1)}} = x^{(0)} (\delta^1)^T$$

$$\frac{\partial \ell}{\partial W^{L-1}}$$

$$\frac{\partial e}{\partial W^L} = x^{L-1} \delta^{(L)}$$

$$\delta^{(L-1)} \longleftarrow \theta'(S^{L-1}) \longleftarrow W^{(l)} \overset{\delta^{(L)}}{\longleftarrow}$$

$$\partial \ell / \partial W^{(l)}$$

## Update Step

$$W^{(l)} \leftarrow W^{(l)} - \eta \frac{\partial e}{\partial W^{(l)}} \checkmark$$

SGD , Batch grad. descent.

# Computing $\boldsymbol{\delta}^{(\ell)}$ Using the Chain Rule

$$\boldsymbol{\delta}^{(1)} \longleftarrow \boldsymbol{\delta}^{(2)} \cdots \longleftarrow \boldsymbol{\delta}^{(L-1)} \longleftarrow \boldsymbol{\delta}^{(L)}$$

Multiple applications of the chain rule:

$$\Delta\mathbf{s}^{(\ell)} \xrightarrow{\theta} \Delta\mathbf{x}^{(\ell)} \xrightarrow{\mathrm{W}^{(\ell+1)}} \Delta\mathbf{s}^{(\ell+1)} \cdots \longrightarrow \Delta\mathbf{e}(\mathbf{x})$$



don't use $0^{\text{th}}$ component (bias)
$\downarrow$

$$\boldsymbol{\delta}^{(\ell)} = \theta'(\mathbf{s}^{(\ell)}) \otimes [\mathrm{W}^{(\ell+1)}\boldsymbol{\delta}^{(\ell+1)}]_1^{d^{(\ell)}}$$

$\uparrow$
componentwise multiplication

# The Backpropagation Algorithm

$$\delta^{(1)} \longleftarrow \delta^{(2)} \cdots \longleftarrow \delta^{(L-1)} \longleftarrow \delta^{(L)}$$

**Backpropagation to compute sensitivities $\delta^{(\ell)}$:**
*(Assume $\mathbf{s}^{(\ell)}$ and $\mathbf{x}^{(\ell)}$ have been computed for all $\ell$)*

1: $\delta^{(L)} \leftarrow 2(x^{(L)} - y) \cdot \theta'(s^{(L)})$     [Initialization] ✓

2: **for** $\ell = L - 1$ to 1 **do**     [Back-Propagation]

3:    Compute (for tanh hidden node):

$$\theta'(\mathbf{s}^{(\ell)}) = \left[1 - \mathbf{x}^{(\ell)} \otimes \mathbf{x}^{(\ell)}\right]_1^{d^{(\ell)}}$$

4:    $\delta^{(\ell)} \leftarrow \theta'(\mathbf{s}^{(\ell)}) \otimes \left[W^{(\ell+1)}\delta^{(\ell+1)}\right]_1^{d^{(\ell)}}$     $\leftarrow$ componentwise multiplication

5: **end for**

# Algorithm for Gradient Descent on $E_{\text{in}}$

---

**Algorithm to Compute $E_{\text{in}}(\mathbf{w})$ and $\mathbf{g} = \nabla E_{\text{in}}(\mathbf{w})$:**
**Input:** weights $\mathbf{w} = \{\mathrm{W}^{(1)}, \ldots, \mathrm{W}^{(L)}\}$; data $\mathcal{D}$.
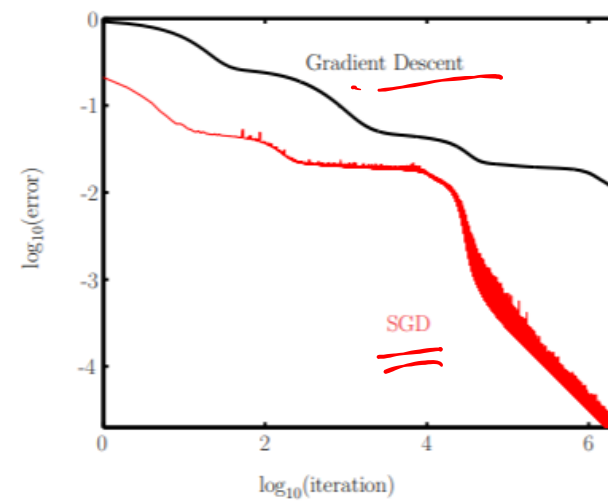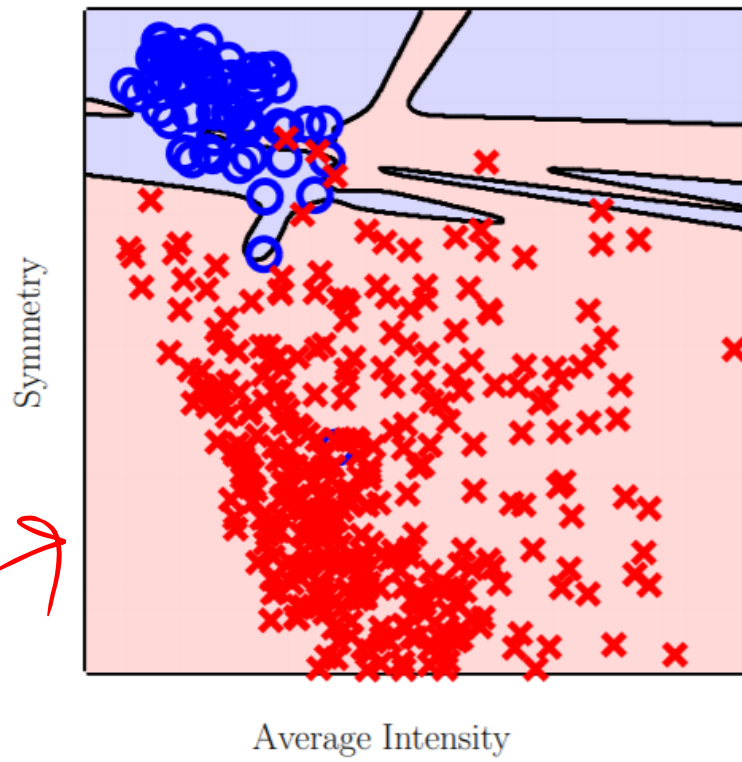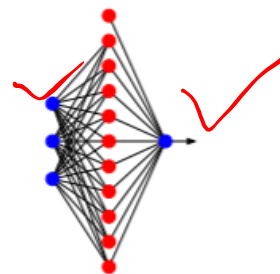**Output:** error $E_{\text{in}}(\mathbf{w})$ and gradient $\mathbf{g} = \{\mathrm{G}^{(1)}, \ldots, \mathrm{G}^{(L)}\}$.

1: Initialize: $E_{\text{in}} = 0$; for $\ell = 1, \ldots, L$, $\mathrm{G}^{(\ell)} = 0 \cdot \mathrm{W}^{(\ell)}$ .

2: **for** Each data point $\mathbf{x}_n$ $(n = 1, \ldots, N)$ **do**

3:     Compute $\mathbf{x}^{(\ell)}$ for $\ell = 0, \ldots, L$.    [forward propagation]

4:     Compute $\boldsymbol{\delta}^{(\ell)}$ for $\ell = 1, \ldots, L$.    [backpropagation]

5:     $E_{\text{in}} \leftarrow E_{\text{in}} + \frac{1}{N}(\mathbf{x}_1^{(L)} - y_n)^2$.

6:     **for** $\ell = 1, \ldots, L$ **do**

7:         $\mathrm{G}^{(\ell)}(\mathbf{x}_n) = [\mathbf{x}^{(\ell-1)}(\boldsymbol{\delta}^{(\ell)})^{\mathsf{T}}]$

8:         $\mathrm{G}^{(\ell)} \leftarrow \mathrm{G}^{(\ell)} + \frac{1}{N}\mathrm{G}^{(\ell)}(\mathbf{x}_n)$.

9:     **end for**

10: **end for**

---

Can do batch version or sequential version (SGD).

# Digits Data



*Overfitting* (handwritten annotation)

Symmetry (vertical axis label)

Average Intensity (horizontal axis label)

Gradient Descent

SGD

$\log_{10}(\text{error})$

$\log_{10}(\text{iteration})$

# Thanks!