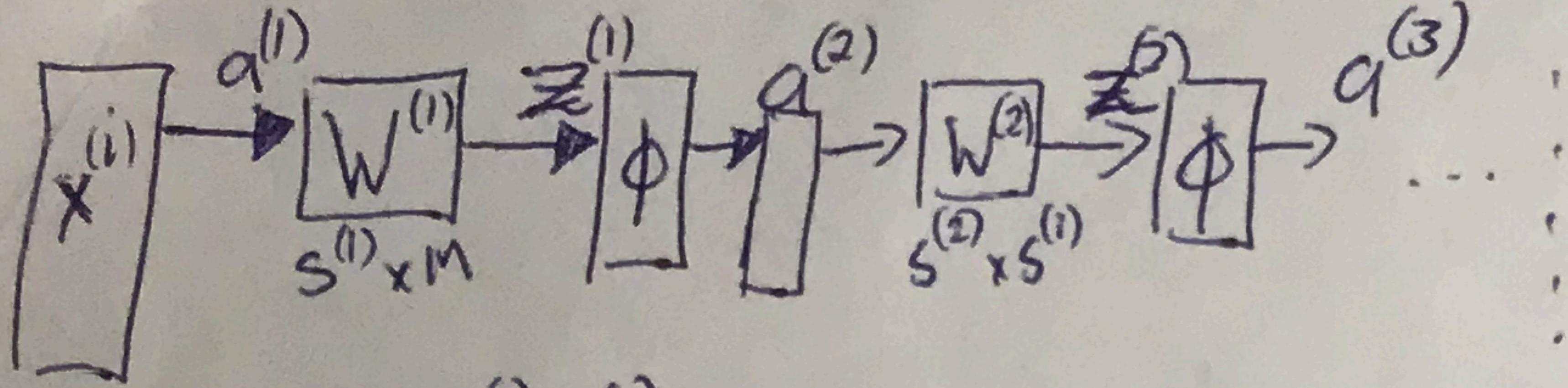


① Recall how a Neural Network Works



$$Z^{(1)} = W^{(1)} a^{(1)}$$

$$Z^{(1)} = \begin{bmatrix} -N^{(1)} \\ -N^{(2)} \\ \vdots \\ -N^{(M)} \end{bmatrix} \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ \vdots \\ a_s^{(1)} \end{bmatrix}$$

$$= \begin{bmatrix} N^{(1)} a^{(1)} \\ 2N^{(2)} a^{(1)} \\ \vdots \\ sN^{(M)} a^{(1)} \end{bmatrix}$$

$$A^{(2)} = \phi(Z^{(1)})$$

$$= \begin{matrix} | & | \\ \phi(Z^{(1)}) & \phi(Z^{(1)}) \\ | & | \\ (1) & (2) \end{matrix} \dots = A^{(2)}$$

$$Z^{(2)} = W^{(2)} \cdot A^{(2)}$$

...

$$Z^{(3)} = W^{(3)} \cdot A^{(3)}$$

ETC.

- TO UPDATE $W^{(k)}$, WE NEED GRADIENT OF THE OBJECTIVE FUNCTION, $J(w)$

$$w_{ij}^{(k)} \leftarrow w_{ij}^{(k)} + \frac{\eta \partial J(w)}{\partial w_{ij}^{(k)}}$$

① REMEMBER WHY THIS DIFFICULT

② DERIVE $\frac{\partial J}{\partial w}$ FOR ONE INSTANCE $x^{(i)}$

③ DEFINE BACK PROPAGATION

④ DERIVE FAST LINEAR ALGEBRA NOTATION FOR ALL INSTANCES

$$X = \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(M)} \end{bmatrix}$$

$$X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(M)} \end{bmatrix}$$

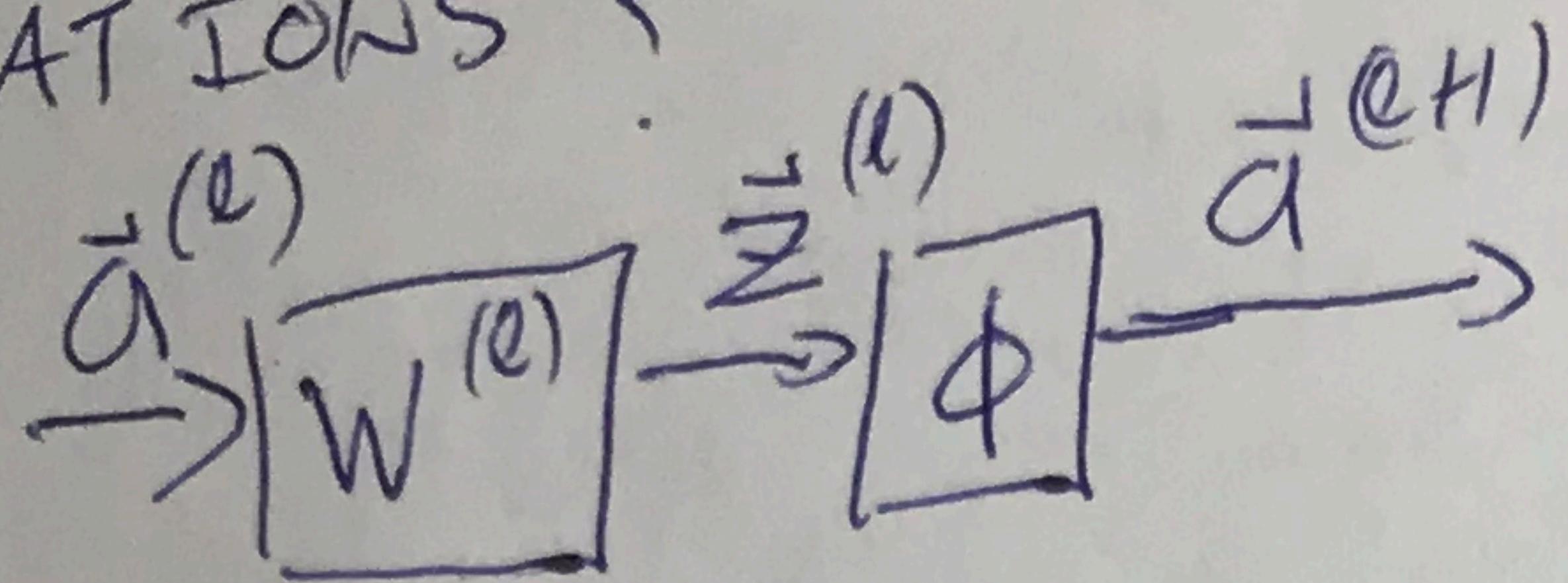
$$A^{(1)} = X^T$$

$$Z^{(1)} = \begin{bmatrix} -N^{(1)} \\ -N^{(2)} \\ \vdots \\ -N^{(M)} \end{bmatrix} \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ \vdots \\ a_s^{(1)} \end{bmatrix}$$

$$Z^{(1)} = W \cdot A^{(1)}$$

$$= \begin{bmatrix} | & | & | \\ Z^{(1)}_1 & Z^{(1)}_2 & Z^{(1)}_M \\ | & | & | \\ (1) & (2) & (M) \end{bmatrix}$$

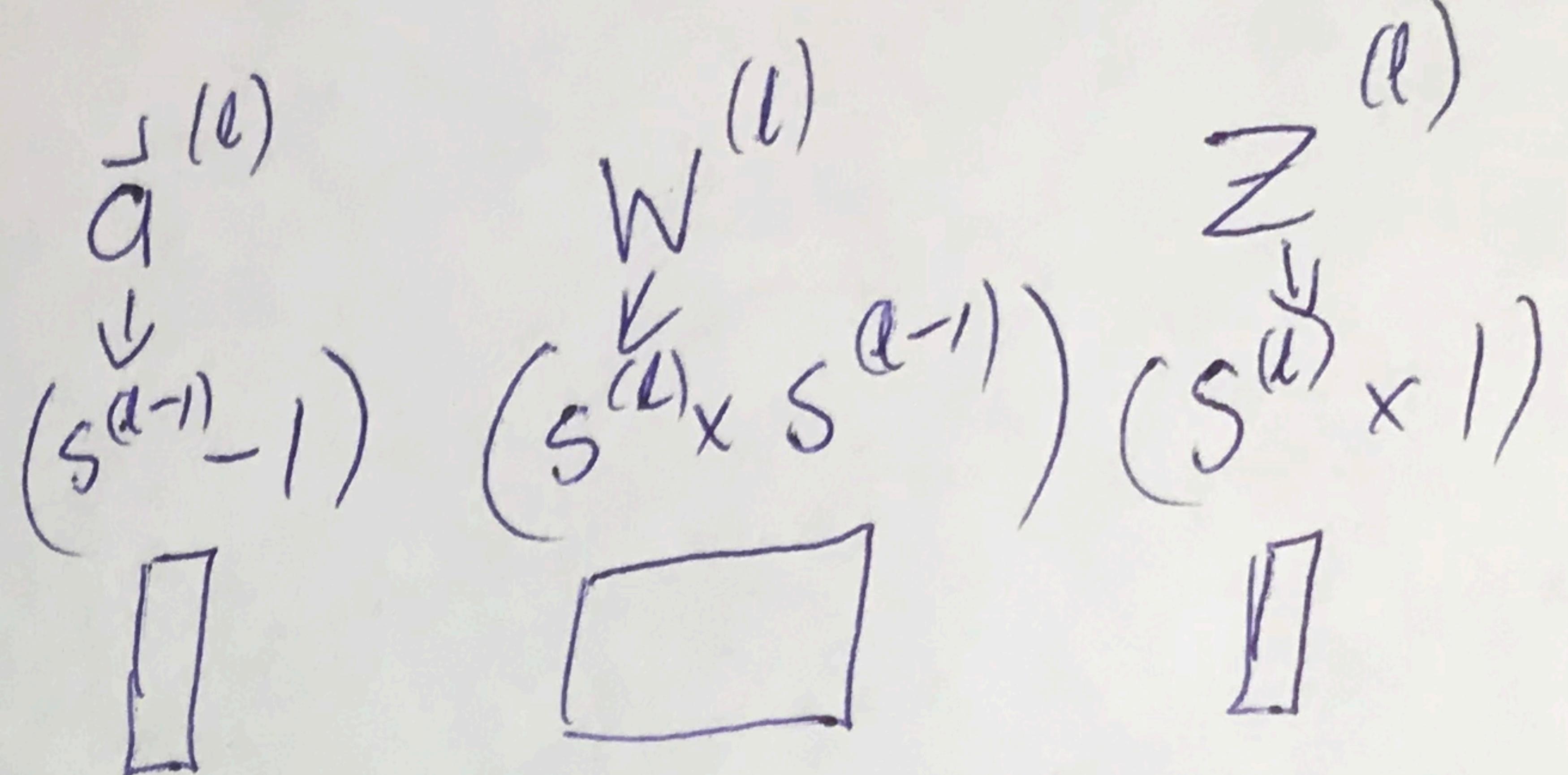
CAN WE WRITE DOWN THE FUNCTIONS OF THE NETWORK AS A SERIES OF COMBINED EQUATIONS?



$$z^{(l)} = W^{(l)} a^{(l)}$$

$$d^{(l+1)} = \phi(W^{(l)} a^{(l)})$$

$$\begin{aligned} a^{(l+1)} &= \phi(W^{(l)} \phi(W^{(l-1)} a^{(l-1)})) \\ &= \phi(W^{(l)} \phi(W^{(l-1)} \phi(W^{(l-2)} a^{(l-2)}))) \end{aligned}$$



IF WE WANT TO KNOW DERIVATIVES,
WE NEED TO USE CHAIN RULE!!

- WHY DO WE NEED DERIVATIVE?

WE NEED GRADIENT OF OBJECTIVE FUNCTION

WITH RESPECT TO EACH $W^{(l)}$

$$J(w) = \sum_{k=1}^M |y^{(k)} - \underbrace{[a^{(l)}]^{(k)}}_{}|^2$$

TO GET OUR UPDATE EQUATIONS, WE NEED

$$W^{(l)} \leftarrow W^{(l)} + \eta \frac{\partial J(w)}{\partial W^{(l)}}$$

OR

$$w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} + \eta \frac{\partial J(w)}{\partial w_{ij}^{(l)}} \quad \text{GETTING THIS IS AN EXERCISE IN USING CHAIN RULE..}$$

SO LET'S REVIEW THE CHAIN RULE

$$\frac{\partial}{\partial x} f(g(x)) = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x}$$

$$f(x) = \ln(\underbrace{1+x^2}_g) \quad \frac{\partial}{\partial g} \ln(g) = \frac{1}{g}$$

$$\frac{\partial f}{\partial x} = \frac{1}{1+x^2} \frac{\partial}{\partial x} (1+x^2)$$

$$= \frac{1}{1+x^2} 2x = \frac{2x}{1+x^2}$$

AND IT WORKS CHAINED AS
MANY TIMES AS WE WANT!

$$f = (1 + \ln(x^3))^2$$

$$\frac{\partial f}{\partial x} = 2(1 + \ln(x^3)) \frac{\partial}{\partial x} (1 + \ln(x^3))$$

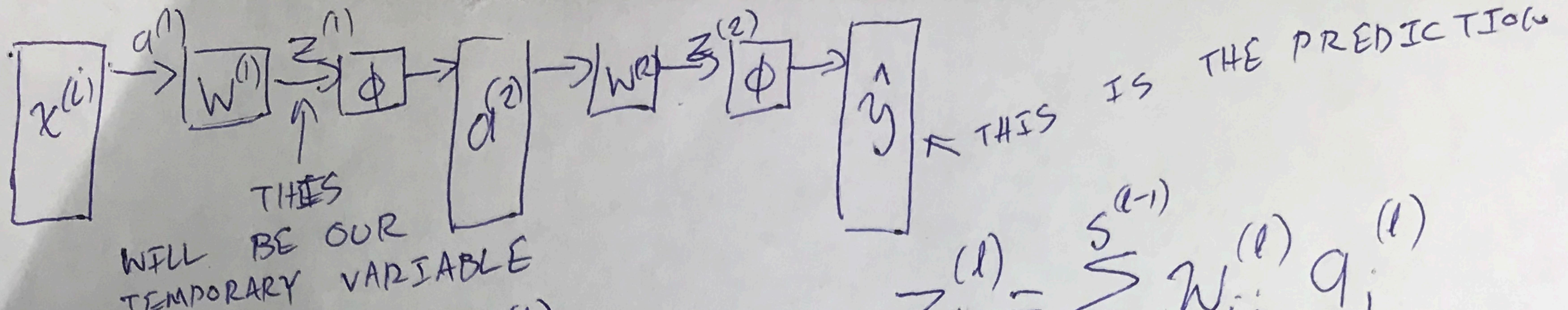
$$= 2(1 + \ln(x^3)) \frac{1}{x^3} \frac{\partial}{\partial x} x^3$$

$$= 2(1 + \ln(x^3)) \frac{1}{x^3} 3x^2$$

$$= \frac{6(1 + \ln(x^3))}{x}$$

SO WE CAN

✓ USE CHAIN RULE
ON OUR NEURAL
NETWORK. WE WILL
USE "Z" AS OUR
INTERMEDIATE
VARIABLE!!



WE NEED

$$\frac{\partial J}{\partial w_{ij}^{(l)}} = \frac{\partial J}{\partial z_i^{(l)}} \frac{\partial z_i^{(l)}}{\partial w_{ij}^{(l)}}$$

$$= \frac{\partial J}{\partial z_i^{(l)}} \frac{\partial}{\partial w_{ij}^{(l)}} \left(\sum_{q=1}^{S^{(l-1)}} w_{iq}^{(l)} a_q^{(l)} \right)$$

↑ DUMMY VARIABLE

$$= \frac{\partial J}{\partial z_i^{(l)}} a_j^{(l)}$$

DEFINE AS

$$= v_i^{(l)} a_j^{(l)}$$

$$z_i^{(l)} = \sum_{j=1}^{S^{(l-1)}} w_{ij}^{(l)} a_j^{(l)}$$

IF $q \neq j$ THEN
DERIVATIVE IS \emptyset

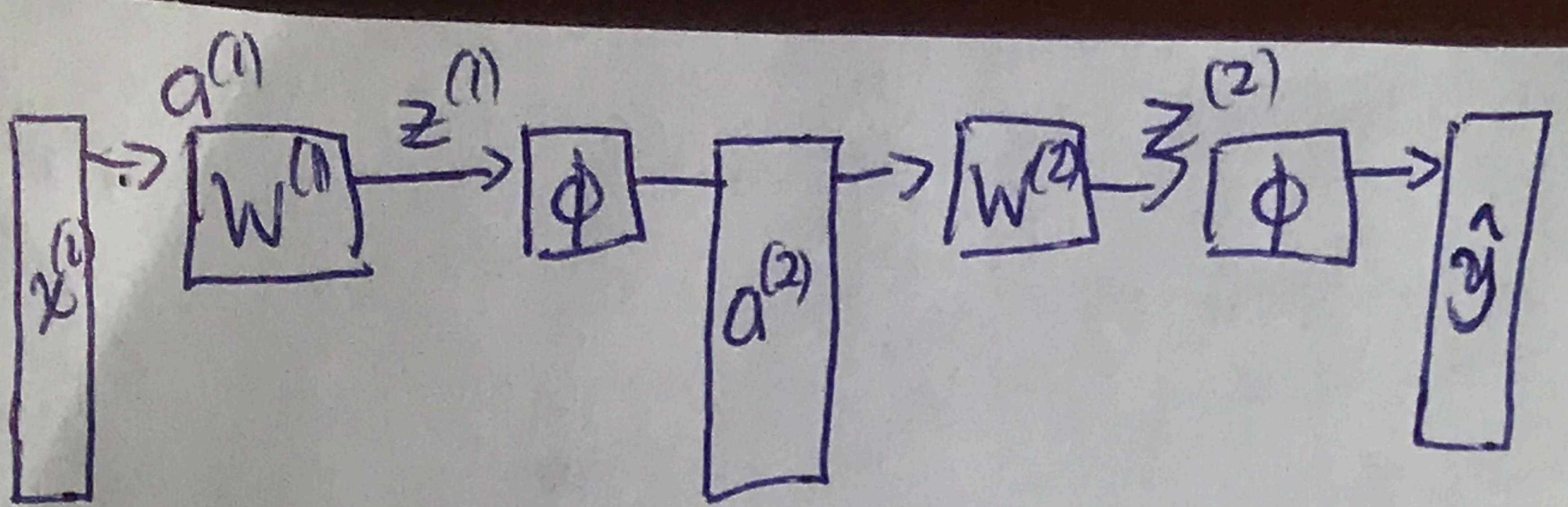
THE "SENSITIVITY", $v_i^{(l)}$

$$W^{(l)} \leftarrow W^{(l)} + \eta \tilde{V}^{(l)} [\vec{a}^{(l)}]^T$$

$$\begin{bmatrix} v_1^{(l)} \\ v_2^{(l)} \\ \vdots \\ v_{S^l}^{(l)} \end{bmatrix} \begin{bmatrix} a_1^{(l)} & a_2^{(l)} & \dots & a_{S^{(l-1)}}^{(l)} \end{bmatrix} = \begin{bmatrix} a_1 v_1 & a_2 v_1 & \dots & a_{S^{(l-1)}} v_1 \\ a_1 v_2 & \vdots & & \vdots \\ \vdots & & & \vdots \\ a_1 v_{S^l} & a_2 v_{S^l} & \dots & a_{S^{(l-1)}} v_{S^l} \end{bmatrix}$$

BUT WE NEED TO ADD UP THE
RESPONSE FOR ALL ACTIVATIONS FROM ALL
INSTANCES

$$W^{(l)} \leftarrow W^{(l)} + \eta \sum_{k=1}^M [\tilde{V}^{(l)}]^{(k)} [\vec{a}^{(l)T}]^{(k)}$$



$$z_i^{(l)} = \sum_{j=1}^{s^{(l+1)}} w_{ij}^{(l)} a_j^{(l)}$$

$$\begin{aligned}\frac{\partial J}{\partial w_{ij}^{(l)}} &= \frac{\partial J}{\partial z_i^{(l)}} \frac{\partial z_i^{(l)}}{\partial w_{ij}^{(l)}} \\ &= \frac{\partial J}{\partial z_i^{(l)}} \frac{\partial}{\partial w_{ij}^{(l)}} \left(\sum_{q=1}^{s^{(l+1)}} w_{iq}^{(l)} a_q^{(l)} \right) \\ &= \underbrace{\frac{\partial J}{\partial z_i^{(l)}}}_{i^{\text{th}} \text{ SENSITIVITY}, v_i^{(l)}} a_j^{(l)} \\ &= v_i^{(l)} a_j^{(l)}\end{aligned}$$

$$W^{(l)} \leftarrow W^{(l)} + \eta \tilde{V}^{(l)} [\bar{a}^{(l)}]^T$$

$$\begin{bmatrix} v_1^{(l)} \\ v_2^{(l)} \\ \vdots \\ v_{s^l}^{(l)} \end{bmatrix} [a_1^{(l)} \dots a_{s^{l+1}}^{(l)}] = \begin{bmatrix} a_1 v_1 & a_2 v_1 & \dots & a_{s^{l+1}} v_1 \\ a_1 v_2 & a_2 v_2 & & \\ \vdots & & & \\ a_1 v_{s^l} & \dots & \dots & a_{s^{l+1}} v_{s^l} \end{bmatrix}$$

$s^l \times s^{l+1}$

BUT WE NEED TO DO THIS FOR ALL ACTIVATIONS

$$W^{(l)} \leftarrow W^{(l)} + \eta \sum_{k=1}^M [\tilde{V}^{(l)}]^{(k)} [\bar{a}^{(l)}]^T$$

~~$$W^{(l)} \leftarrow W^{(l)} + \eta \tilde{V}^{(l)} [\bar{A}^{(l)}]^T$$~~

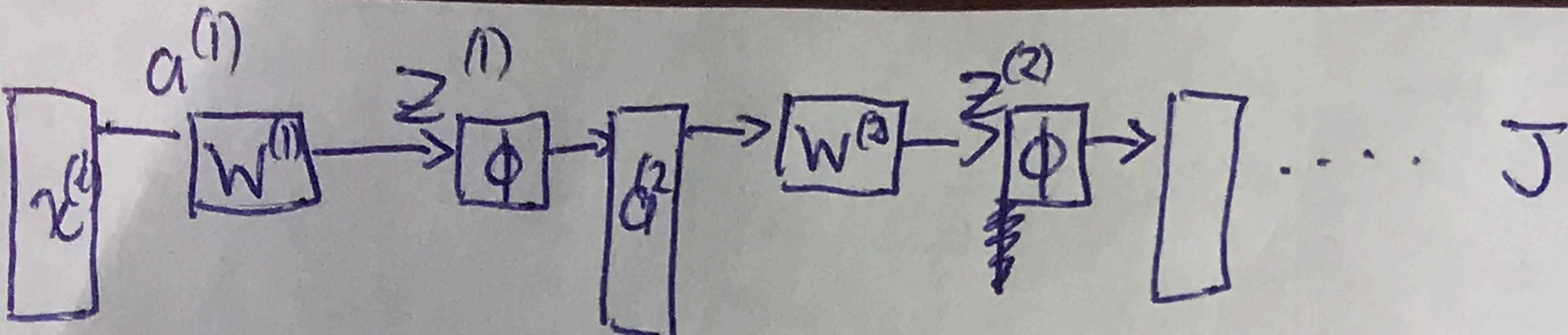
$(s^l \times M) (M \times s^{l+1})$

$$W^{(l)} \leftarrow W^{(l)} + \eta V^{(l)} [A^{(l)}]^T$$

$$V^{(l)} = \begin{bmatrix} 1 & & & \\ \vdots & \ddots & \dots & 1 \\ 1 & & & \end{bmatrix}$$

$$A^{(l)T} = \begin{bmatrix} -[a^{(l)}_1] \\ \vdots \\ -[a^{(l)}_{s^{l+1}}] \end{bmatrix}$$

HOW DO WE CALCULATE THE SENSITIVITIES?



$$\bar{V}^{(l)} = \frac{\partial J}{\partial \bar{z}^{(l)}} = \frac{\partial J}{\partial \bar{z}^{(l+1)}} \quad \frac{\partial \bar{z}^{(l+1)}}{\partial \bar{z}^{(l)}} = \underbrace{\begin{bmatrix} \frac{\partial \bar{z}^{(l+1)}}{\partial z^{(l)}} \\ \vdots \\ \frac{\partial \bar{z}^{(l+1)}}{\partial z^{(l)}} \end{bmatrix}}_{\text{JACOBIAN } D}^T \frac{\partial J}{\partial \bar{z}^{(l+1)}}$$

CHAIN RULE FOR MATRICES
OF EQUATIONS IS DEFINED
BY JACOBIAN

$V^{(l)} = D^T V^{(l+1)}$

THIS IS A RECURRENCE RELATION!
WE CAN BACK PROPAGATE THE SENSITIVITY!

$$D_{ij} = \frac{\partial z_i^{(l+1)}}{\partial z_j^{(l)}} = \begin{bmatrix} \frac{\partial z_1^{(l+1)}}{\partial z_1^{(l)}} & \frac{\partial z_1^{(l+1)}}{\partial z_2^{(l)}} & \dots & \frac{\partial z_1^{(l+1)}}{\partial z_{S^{(l+1)}}^{(l)}} \\ \frac{\partial z_2^{(l+1)}}{\partial z_1^{(l)}} & & & \\ \vdots & & & \\ \frac{\partial z_{S^{(l+1)}}^{(l+1)}}{\partial z_1^{(l)}} & & & \frac{\partial z_{S^{(l+1)}}^{(l+1)}}{\partial z_{S^{(l+1)}}^{(l)}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial z_i^{(l+1)}}{\partial z_j^{(l)}} &= \frac{\partial}{\partial z_j^{(l)}} \left(\sum_{q=1}^{S^{(l+1)}} w_{iq}^{(l+1)} q_q^{(l+1)} \right) \\ &= \frac{\partial}{\partial z_j^{(l)}} \sum_{q=1}^{S^{(l+1)}} w_{iq}^{(l+1)} \phi(z_q^{(l)}) \quad q=j \text{ for nonzero!} \\ &= w_{ij}^{(l+1)} \phi'_j \quad \text{if } \phi \text{ is sigmoid } \phi'_j = a_j(1-a_j) \end{aligned}$$

$$D_{ij} \stackrel{\text{FOR SIG.}}{=} w_{ij}^{(l+1)} a_j^{(l+1)} (1 - a_j^{(l+1)})$$

$$D = W^{(l+1)} \Phi'^{(l+1)}$$

$$V^{(l)} = D^T V^{(l+1)}$$

$$V^{(l)} = [W^{(l+1)} \dot{\Phi}^{(l+1)}]^T V^{(l+1)} = \dot{\Phi}^{(l+1)} [W^{(l+1)}]^T V^{(l+1)}$$

$$= \dot{\Phi}^{(l+1)} \otimes [W^{(l+1)}]^T V^{(l+1)} = \begin{bmatrix} \dot{\Phi} \\ \vdots \\ \dot{\Phi} \end{bmatrix} = \begin{bmatrix} W V \\ \vdots \\ W V \end{bmatrix} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$\dot{\Phi}' = \text{DIAG}(\dot{\Phi}') = \begin{bmatrix} \dot{\Phi}_1^{(l+1)} & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \dot{\Phi}_{S^{(l+1)}}^{(l+1)} \end{bmatrix}$$

$$\begin{aligned}
 V^{(L)} &= \begin{bmatrix} | & | & | \\ [V^{(1)}]^{(1)} & [V^{(2)}]^{(2)} & \dots & [V^{(M)}]^{(M)} \\ | & | & | \end{bmatrix} \\
 &= \left[\begin{bmatrix} \vec{j}^{(L+1)} \end{bmatrix}^{(1)} \otimes \begin{bmatrix} W^{(L+1)} \end{bmatrix}^T \begin{bmatrix} V^{(L+1)} \end{bmatrix}^{(1)} \right. \dots \left. \begin{bmatrix} \vec{\phi}^{(L+1)} \end{bmatrix}^{(M)} \otimes \begin{bmatrix} W^{(L+1)} \end{bmatrix}^T \begin{bmatrix} V^{(L+1)} \end{bmatrix}^{(M)} \right] \\
 &= \vec{\phi}^{(L+1)} \otimes \begin{bmatrix} W^{(L+1)} \end{bmatrix}^T V^{(L+1)} \\
 &\quad (S^L \times M) \quad (\underbrace{S^L \times S^{(L+1)}}_{(S^L \times S^{L+1})} \times \underbrace{(S^{L+1} \times M)}_{(S^{L+1} \times M)})
 \end{aligned}$$

$$\begin{aligned}
 \vec{j} &= \vec{a} \otimes (1 - \vec{a}) \\
 \vec{\phi}^{(L+1)} &= A^{(L+1)} (1 - A^{(L+1)})
 \end{aligned}$$

$$V^{(L)} = A^{(L+1)} (1 - A^{(L+1)}) \otimes \begin{bmatrix} W^{(L+1)} \end{bmatrix}^T V^{(L+1)}$$

FAST CALCULATION OF RECURRENCE RELATION!

NEED $V^{(L)}$ at THE LAST LAYER

$$\begin{aligned} J(w) &= \sum_{k=1}^M |\tilde{y}^{(k)} - [a^{(L)}]^{(k)}|^2 \\ &= \sum_{k=1}^M |y^{(k)} - \phi([\tilde{z}^{(L-1)}]^{(k)})|^2 \end{aligned}$$

$$[\nabla_i^{(L-1)}]^{(k)} = \left[\frac{\partial J}{\partial z_i^{(L-1)}} \right]^{(k)}$$

~~$$= \frac{\partial}{\partial z_i^{(L-1)}} (\tilde{y}^{(k)} - \phi([\tilde{z}^{(L-1)}]^{(k)}))^2$$~~

$$= \frac{\partial}{\partial z_i^{(L-1)}} (y_i^{(k)} - \phi([\tilde{z}_i^{(L-1)}]^{(k)}))^2$$

$$= \frac{\partial}{\partial z_i^{(L-1)}} (y_i^{(k)} - \phi([z_i^{(L-1)}]^{(k)}))^2$$

$$= (y_i^{(k)} - \phi(\cdot)) 2 \left(-\frac{\partial}{\partial z_i^{(L-1)}} \phi(z_i^{(L-1)}) \right)$$

$$= -2(y_i^{(k)} - \phi(\cdot)^{(k)}) [a_i^{(L)}]^{(k)} (1 - [a_i^{(L)}]^2)$$

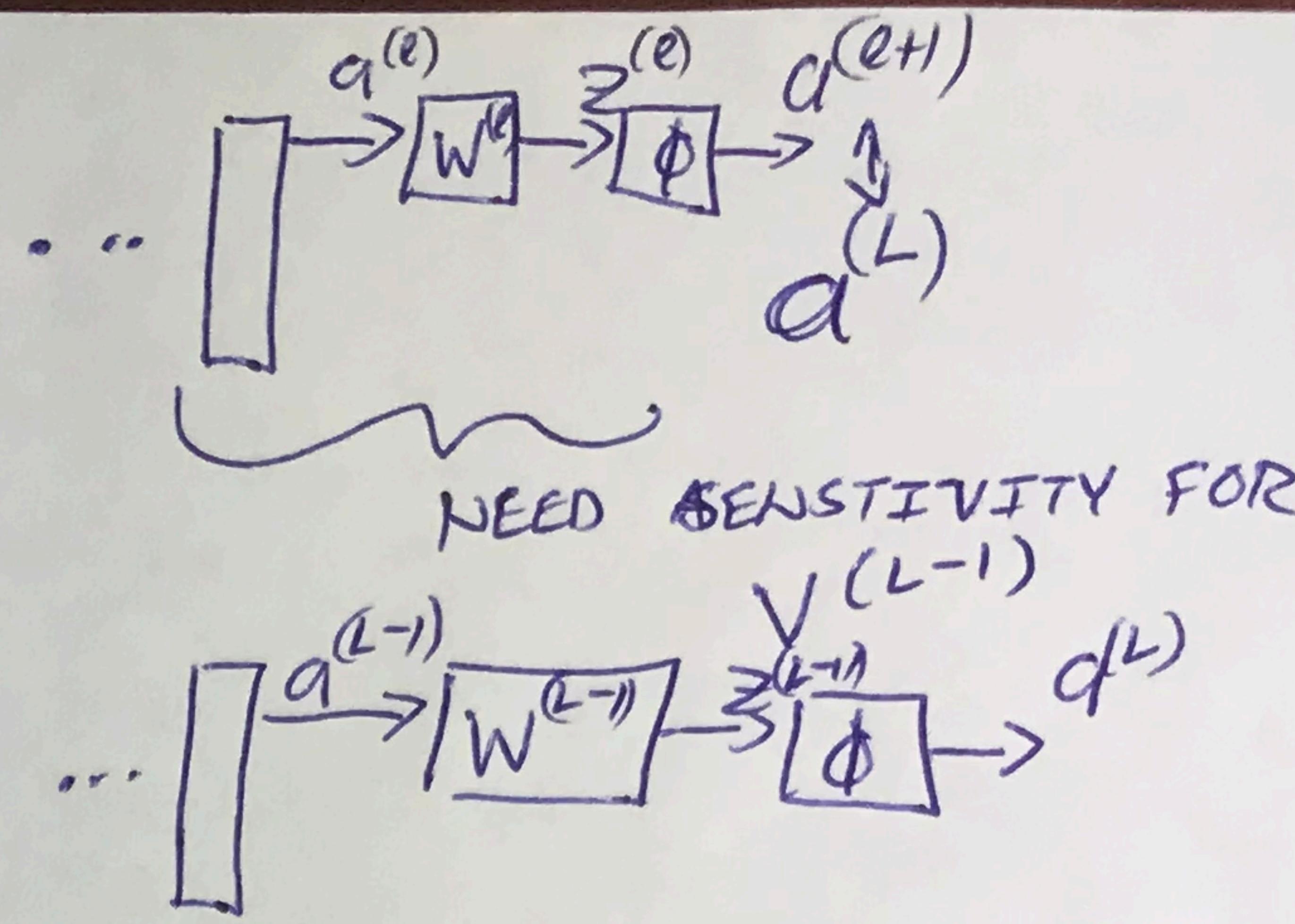
~~$$\tilde{y}^{(k)} = \tilde{z}^{(L-1)} \otimes \phi$$~~

~~$$\tilde{y}^{(k)} = f_2 \times \phi$$~~

$$= -2(y_i^{(k)} - [a_i^{(L)}]) [a_i^{(L)}]^{(k)} (1 - [a_i^{(L)}]^2)$$

$$[\nabla^{(L-1)}]^{(k)} = -2(\tilde{y}^{(k)} - [a^{(L)}]^{(k)}) \otimes [\tilde{a}^{(L)}]^{(k)} \otimes (1 - [a^{(L)}]^2)$$

$$\nabla^{(L-1)} = -2(Y - A^{(L)}) * A^{(L)} * (I - A^{(L)})$$



$$\nabla^{(L-1)} = \begin{bmatrix} 1 & & & & 1 \\ [\nabla^{(L-1)}]^{(1)} & \cdots & [\nabla^{(L-1)}]^{(M)} \\ 1 & & & & 1 \end{bmatrix}$$

\uparrow
SOLVE I ST

SUMMARY

- ① PERFORM FEED FORWARD

$$x^T = A^{(1)} \quad z^{(1)} = W^{(0)} A^{(1)}$$

$$z^{(l)} = W^{(l)} A^{(l)}$$

$$\phi(z^{(l)}) = A^{(l+1)}$$

- ② GET FINAL LAYER SENSITIVITY

$$V^{(L-1)} = -2(Y - A^{(L)}) * A^{(L)} * (1 - A^{(L)})$$

- ③ BACK PROP. ALL SENS.

$$V^{(l)} = A^{(l+1)} \odot (1 - A^{(l+1)}) * [W^{(l+1)}]^T V^{(l+1)}$$

- ④ UPDATE ALL WEIGHT MATRICES

$$W^{(l)} \leftarrow W^{(l)} + \eta V^{(l)} [A^{(l)}]^T$$

- Self Test: We assumed that no bias term was in each layer. Can we update our notation to include bias vectors for each layer?

