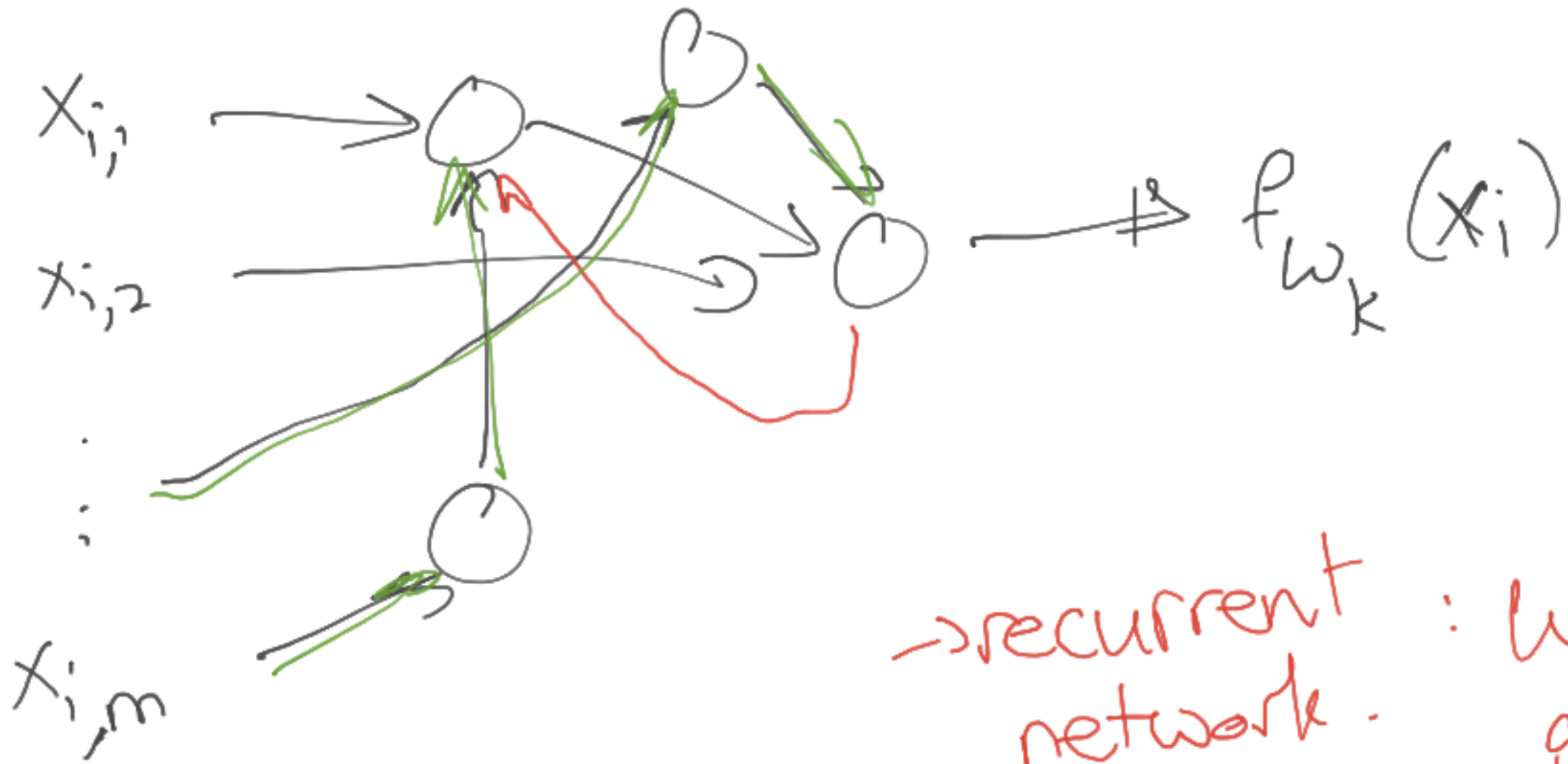


Artificial Neural Networks (ANN)

Lecture 8
-ANNs!



"Network architecture"
↳ how perceptrons
are wired together.
(and how many).

→ recurrent
network.

: when the
graph is
cyclic.

① → ~~perceptron~~
→ unit
→ "node"
→ vertex.

- We consider
non-recurrent
networks.

Perceptron

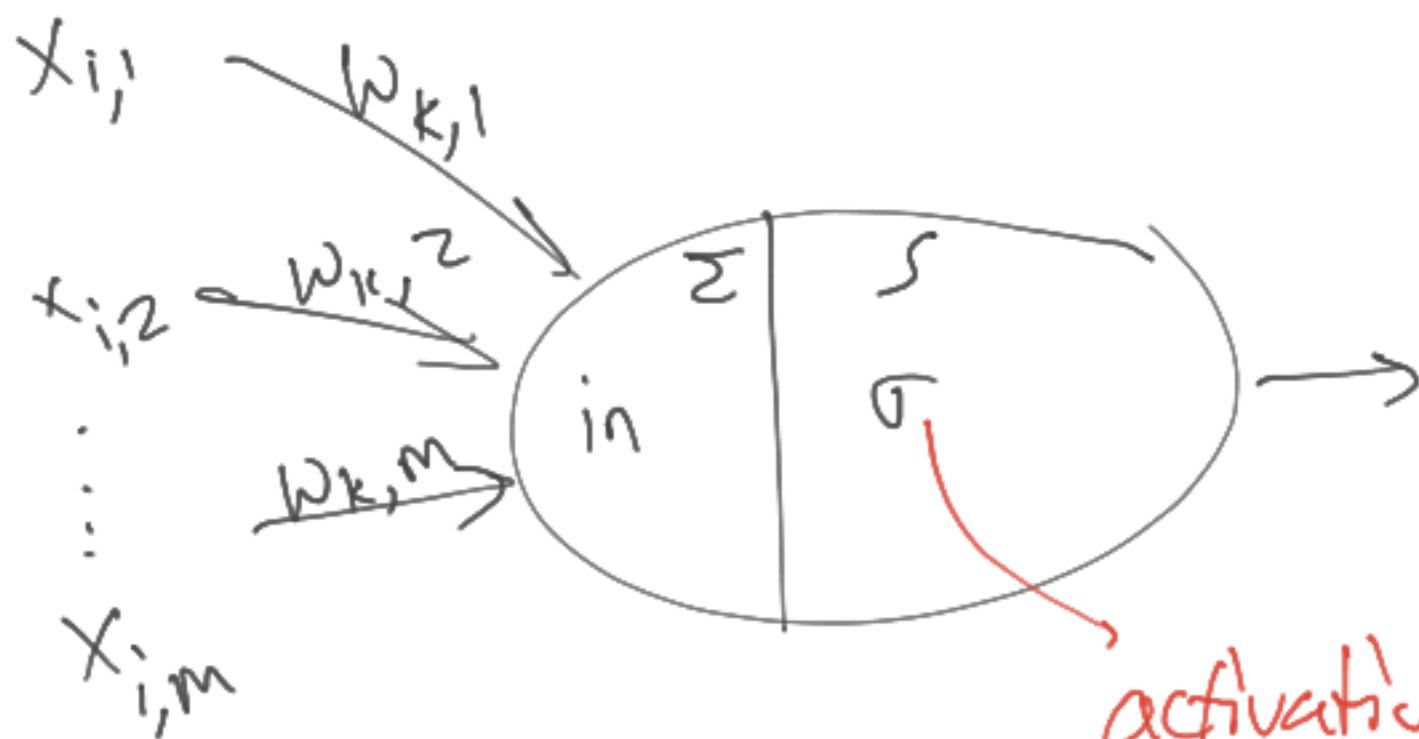
$$(x_i, y_i)_{i=1}^n$$

$$l(w_k) = \sum_{i=1}^n (y_i - f_{w_k}(x_i))^2$$

$$\theta_j : w_{k+1,j} = w_{k,j} - \alpha \frac{\partial l(w_k)}{\partial w_{k,j}}$$

$$\frac{\partial l(w_k)}{\partial w_{k,j}} = \sum_i (y_i - f_{w_k}(x_i)) \frac{\partial f_{w_k}(x_i)}{\partial w_{k,j}}$$

new notation:
 $\frac{\partial f_w(x)}{\partial w_{i,j}}$

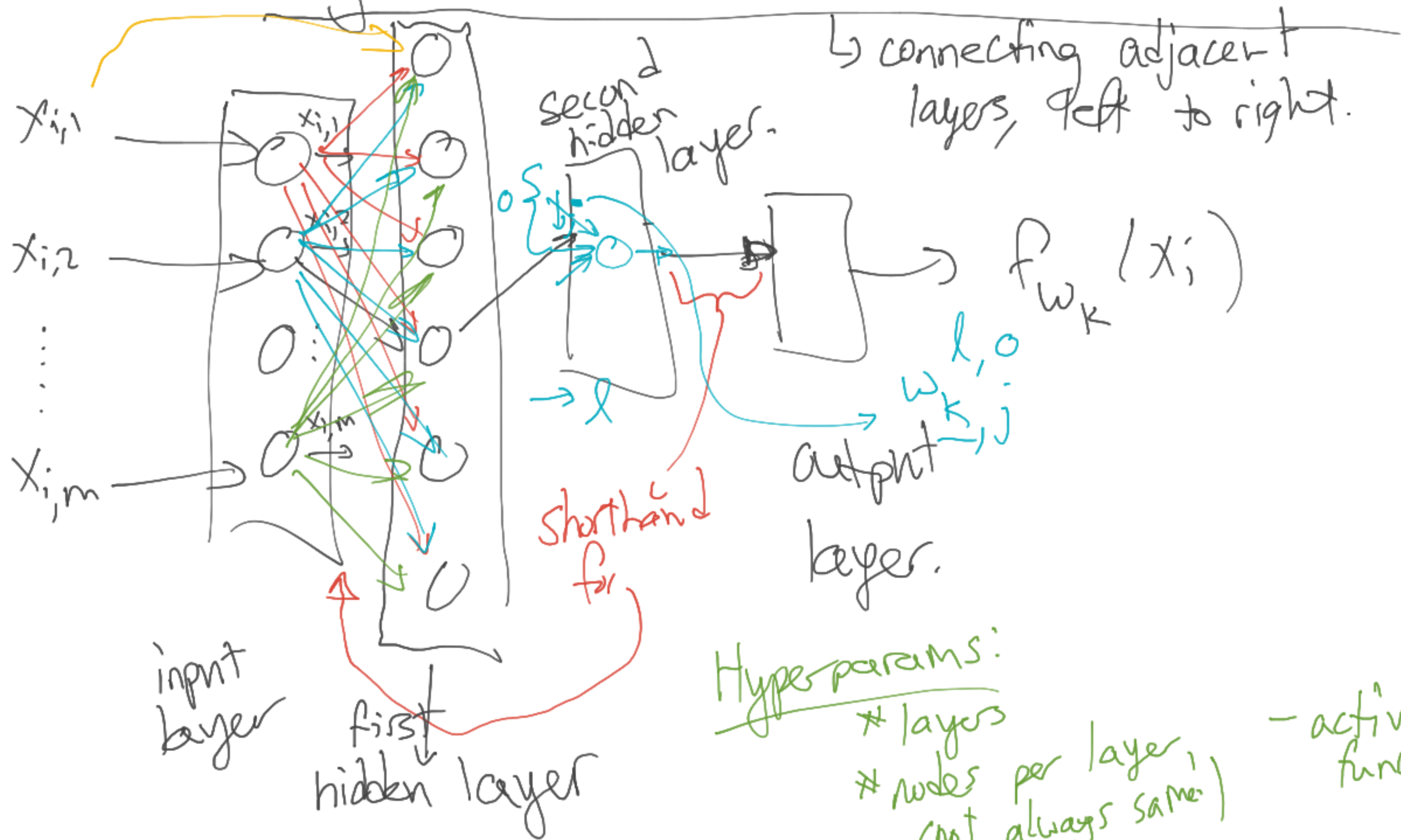


$$in = \sum_{j=1}^m x_{i,j} w_{k,j}$$

$$f_{w_k}(x_i) = \sigma(in)$$

activation
 function.
 "nonlinearity"

Fully Connected Feedforward Network.



Hyperparams:
layers
nodes per layer,
(not always same)

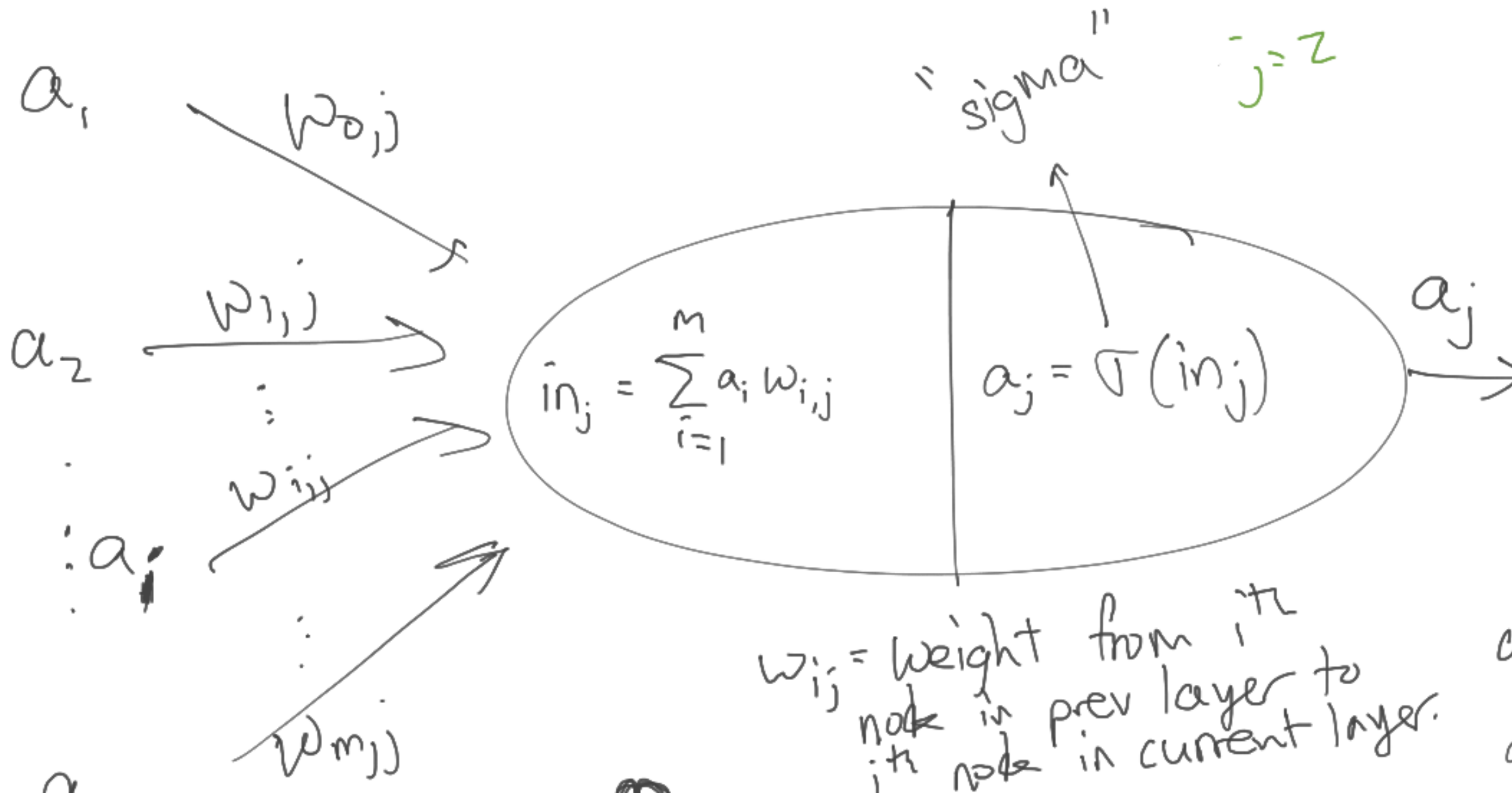
- activation function.

New Notation

Trick: Assume

"current layer"
 $\rightarrow j$
 "previous layer"
 $\rightarrow i$

"next layer"
 $\rightarrow k$
 i j k
 prev ↓ current ↘ next.



$j = z$

$w_{i,j}$ = weight from i th node in prev layer to j th node in current layer.

a_i : output of i th node in prev layer.
 a_j : output of j th node in current layer.

$x \in \mathbb{R}^m$
 \hookrightarrow input to network. } not x_i any more.

Forward Pass

// load input layer

for each input x_j in X → input to network.

$a_j \leftarrow x_j$ // a_i from input layer.

// Run hidden layers.

for each hidden layer

for each node j in current hidden layer do:

$$in_j = \sum_i a_i w_{ij}$$

$$a_j = \sigma(in_j)$$

// Return: output of output layer.

weights 4-D.

$$w[l][j][i]$$

layer node i th input.

$f(x)$

$$x_{k+1} = x_k - \alpha \nabla f(x_k)$$

$$x \leftarrow x - \alpha \nabla f(x)$$

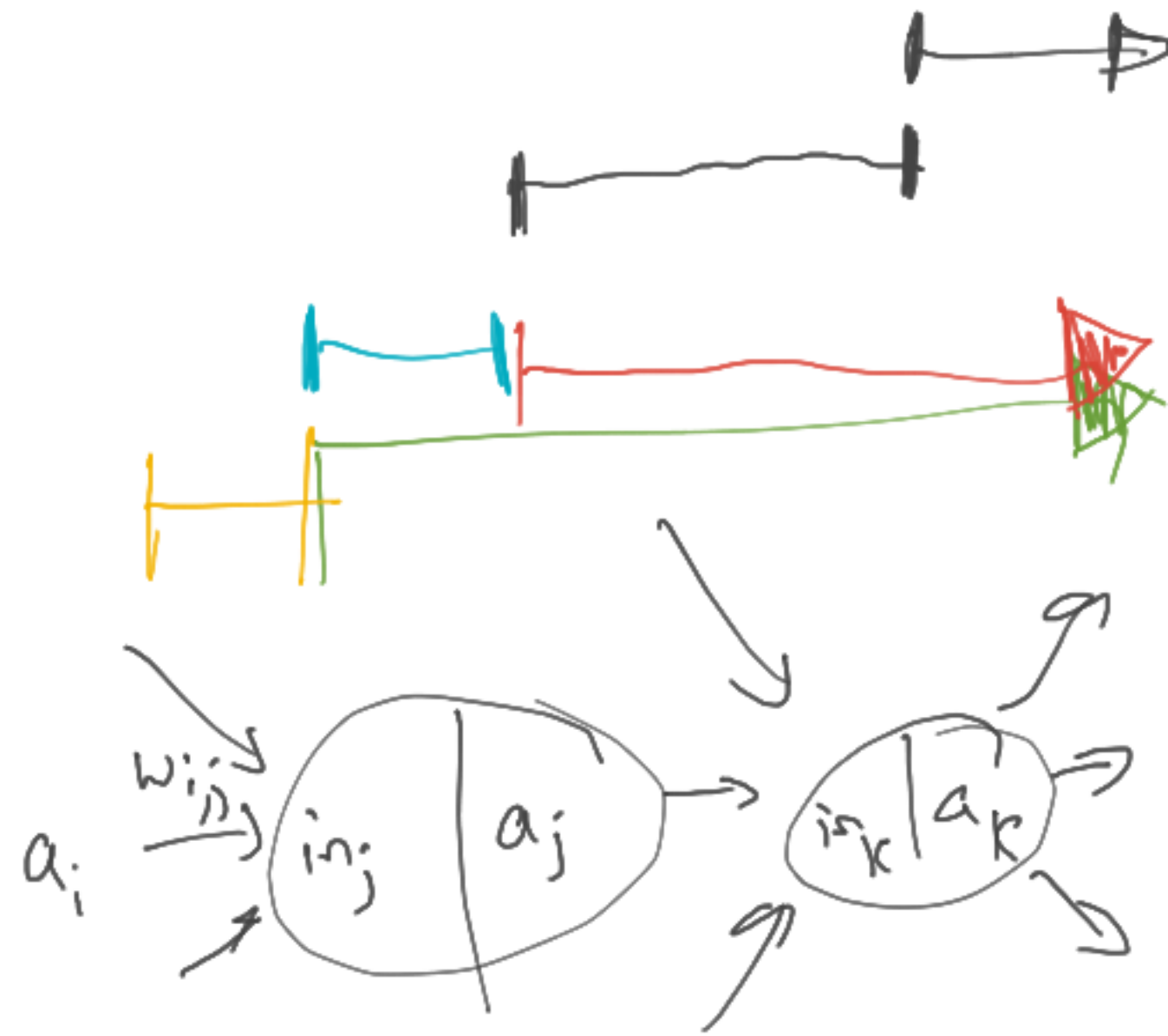
~~$x \leftarrow x - \alpha \nabla f(x)$~~

$$\frac{\partial f_w(x)}{\partial w_{ij}} = \underbrace{\frac{\partial f_w(x)}{\partial in_j}}_{\Delta_j} \frac{\partial in_j}{\partial w_{ij}}$$

$$\Delta_j = \frac{\partial f_w(x)}{\partial in_j} \quad (\text{pink})$$

$$= \underbrace{\frac{\partial f_w(x)}{\partial a_j}}_{\Delta_j} \frac{\partial a_j}{\partial in_j} \frac{\partial in_j}{\partial w_{ij}}$$

lecture 9
- Back propagation



$$\frac{\partial in_j}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \sum_{\beta} a_{\beta} w_{\beta j}$$

$$= \underline{a_i}$$

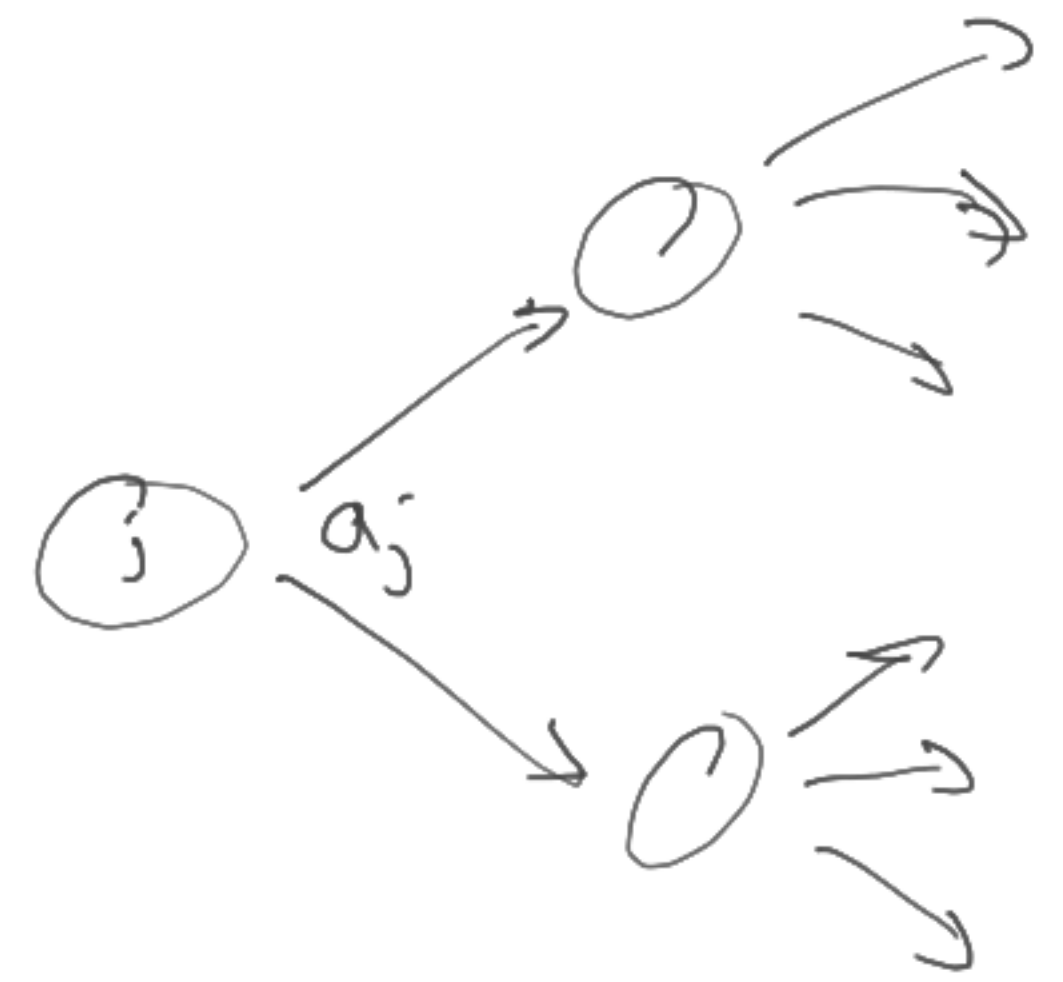
$$\frac{\partial a_j}{\partial in_j} = \frac{\partial \sigma(in_j)}{\partial in_j}$$

$$= \underline{\sigma(in_j) (1 - \sigma(in_j))}$$

$$\frac{\partial \sigma(x)}{\partial x} = \sigma(x) (1 - \sigma(x))$$

$$\frac{\partial f_w(x)}{\partial a_j} = \underbrace{\frac{\partial f_w(x)}{\partial \text{in}_k}}_{\Delta_k} \underbrace{\frac{\partial \text{in}_k}{\partial a_j}}_{= w_{jk}}$$

$$\frac{\partial f_w(x)}{\partial a_j} = \sum_k \frac{\partial f_w(x)}{\partial \text{in}_k} \frac{\partial \text{in}_k}{\partial a_j}$$



$$\frac{\partial f_w(x)}{\partial w_{ij}} = \underbrace{\sum_k \Delta_k w_{j,k}}_{\Delta_j} \sigma(\text{in}_j) (1 - \sigma(\text{in}_j)) a_j$$

Δ_j for output layer?

$$a_j = \sigma(\text{in}_j)$$

$$\Delta_j = \frac{\partial f_w(x)}{\partial \text{in}_j} = \frac{\partial f_w(x)}{\partial a_j} \frac{\partial a_j}{\partial \text{in}_j} = \underline{1} \sigma(\text{in}_j) (1 - \sigma(\text{in}_j))$$

$$w_{ij} \leftarrow w_{ij} - \alpha \sum_{\substack{(x,y) \\ \text{in} \\ \text{data set}}} (y - f_w(x)) \frac{\partial f_w(x)}{\partial w_{ij}}$$

prev slide

Incremental gradient descent. (convergence similar to grad descent)

✓ For each point (x, y) in data set:

for each weight w_{ij}

$$w_{ij} \leftarrow w_{ij} - \alpha (y - f_w(x)) \frac{\partial f_w(x)}{\partial w_{ij}}$$

"estimate of gradient from one point"

Backpropagation

Input: Data set consisting of many input-label pairs, (x, y)

Output: Weights w that approximately minimize the least squares loss, for a fully connected feedforward ANN.

Hyperparams: # of hidden layers, number of nodes per layer, initial weight distribution p
step size α , stopping criterion.

// initialize weights w

for each weight w_{ij} in w

$$w_{ij} \sim p$$

// Run gradient descent

while stopping criterion not satisfied

for each point (x, y) in data set

$\hat{y} = \text{Forwards Pass}(x, w)$

// compute Δ_j for output layer.

for each node j in output layer:

$$\Delta_j = \sigma(\text{in}_j) (1 - \sigma(\text{in}_j))$$

1

2



1) compute Δ_j values for hidden layers.

For each hidden layer FROM OUTPUT TO INPUT

For each node j in hidden layer

$$\Delta_j = \sum_k \Delta_k w_{jk} \sigma'(in_j) (1 - \sigma(in_j))$$

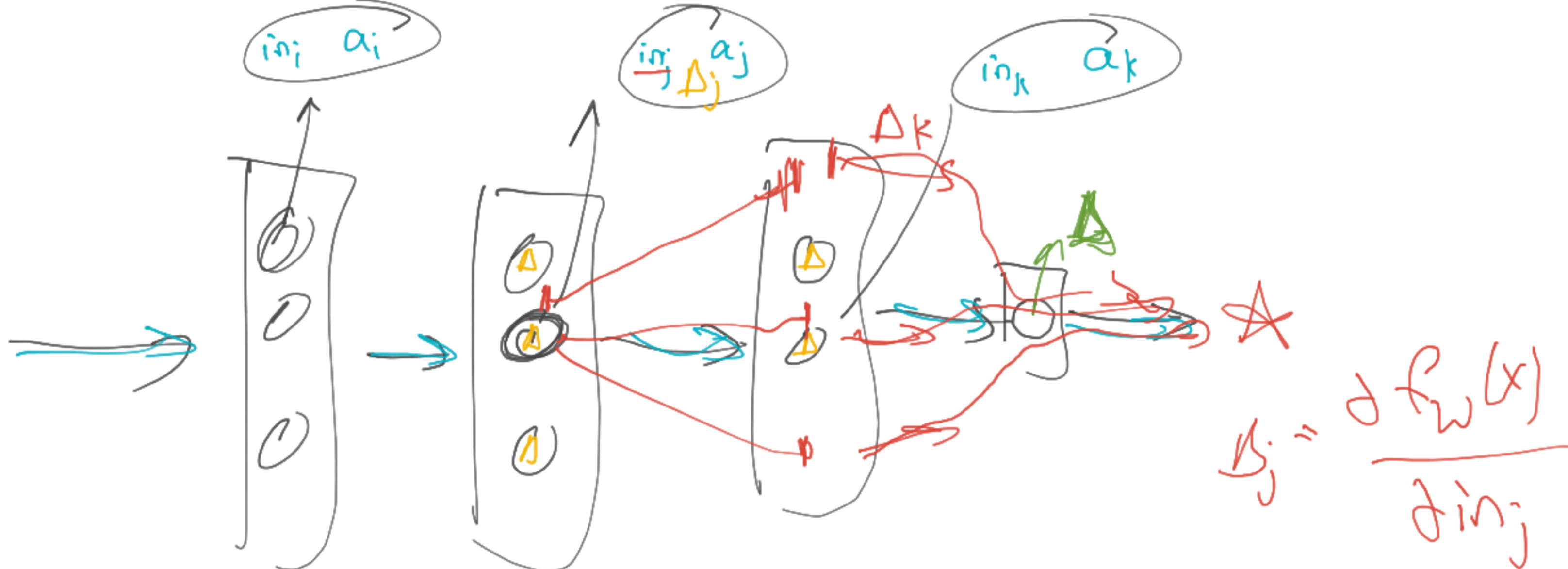
3

1) perform incremental gradient step.

For each weight w_{ij} in w (all weights in all layers)

$$w_{ij} \leftarrow w_{ij} - \alpha \sum (y - \hat{y}) \frac{\partial f}{\partial w_{ij}} \Delta_j a_i$$

return w



(X, W)

- ① Forward pass
- ② Δ_j output special case.
- ③ Backwards pass compute Δ

