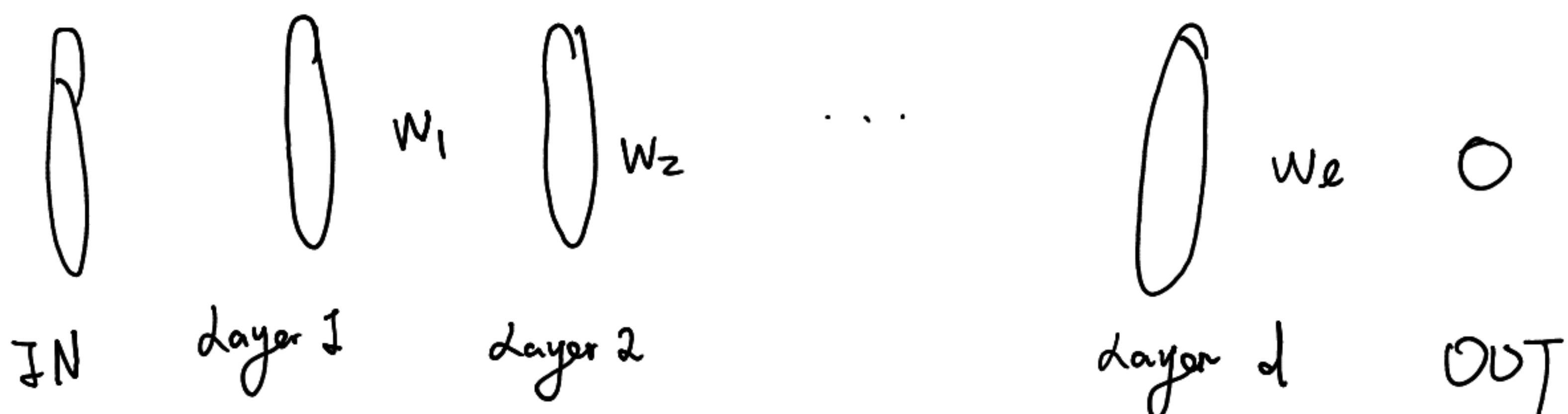


## Lecture Plan :

- \* Background (Recap from last lecture)
- \* Warmup : Multi-layer linear feedforward neural network
- \* The general case:
  - Forward pass
  - Backward pass (with a first attempt)
- \* Extensions :
  - Acyclic computation graphs
  - Weight sharing (e.g. ConvNets)
  - Cyclic computation graphs (RNN)
  - Hessian - vector product

Background. A natural algorithm for training a neural network is gradient descent. How do we run the gradient descent algorithm efficiently?

Warmup: linear neural network.



$$\begin{aligned} \text{loss}(x, y) &= \left( w_d w_{d-1} \dots w_1 x - y \right)^2 \\ &= \ell(w_1, w_2, \dots, w_d) \end{aligned}$$

Gradients.

$$\begin{aligned} \frac{\partial \ell(w_1, \dots, w_d)}{\partial w_i} &:= \ell(w_1, \dots, w_{i-1}, w_i + \delta, w_{i+1}, \dots, w_d) \\ &\quad - \ell(w_1, \dots, w_{i-1}, w_i, w_{i+1}, \dots, w_d) \\ &= (w_d \dots (w_i + \delta) \dots x - y)^2 - (w_d \dots (w_i) \dots x - y)^2 \\ &= 2(w_d \dots w_{i+1} \cdot \delta \cdot w_{i-1} \dots w_1 x - y) \cdot (w_d \dots w_i x - y) \\ &\quad + O(\delta^2) \end{aligned}$$

↑  
(1)

Matrix calculus • (Taylor expansion of matrix functionals) (Expand a bit more)

$$d(\dots, w_i + \delta, \dots) - d(\dots, w_i, \dots)$$

$$= \underbrace{\frac{\partial d(w_1, \dots, w_d)}{\partial w_i}}_{\delta}, \delta \quad (2)$$

Remark: inner product  $\langle a, b \rangle = \sum_{i=1}^n a_i b_i$ , for two vectors  $a, b \in \mathbb{R}^n$ , and

$$\langle X, Y \rangle = \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} X_{i,j} Y_{i,j}, \text{ for two matrices } X, Y \in \mathbb{R}^{n_1 \times n_2},$$

So, our goal is to simplify Eq(1) to a form like Eq(2).

$$W_d \dots W_{i+1} \delta W_i \dots W_1 X$$

$$= \text{Tr} [W_d \dots W_{i+1} \delta W_i \dots W_1 X]$$

$$= \text{Tr} [W_i \dots W_1 X W_d \dots W_{i+1} \delta] \quad (\text{Tr}[ABC] = \text{Tr}[CAB])$$

$$= \langle (W_i \dots W_1 X W_d \dots W_{i+1})^T, \delta \rangle$$

$$= \langle (W_d \dots W_{i+1})^T (W_i \dots W_1 X)^T, \delta \rangle$$

$$\Rightarrow \text{Eq}(2) = 2(W_d \dots W_1 X - y) \cdot (\langle (W_d \dots W_{i+1})^T (W_i \dots W_1 X)^T, \delta \rangle - g)$$

$$\text{Hence } \frac{\partial L}{\partial w_i} = 2(w_d \cdots w_1 x - y) \cdot (w_d \cdots w_{i+1})^T \cdot (w_i \cdots w_1 x)^T. \quad (3)$$

How do we implement the above gradient computation efficiently?

### Forward Pass.

IN:  $x$

Output of layer 1:  $w_1 \cdot x \leftarrow a_1$

⋮

Output of layer  $i$ :  $w_i \cdots w_1 x \leftarrow a_i$

Notice that  $a_i$  appears in Eq (3)!

### Backward Pass.

OUT:  $y$

Gradient wrt to quadratic loss function:

$$g_d = 2(w_d \cdots w_1 x - y)$$

Gradient wrt to the input to the  $d$ -th layer:

$$\begin{aligned} \hat{g}_{d-1} &= \frac{\partial (w_d \cdot z_d - y)^2}{\partial z_d} = 2w_d^T (w_d \cdot z_d - y) \\ &= w_d^T \cdot g_d \end{aligned}$$

Gradient wrt to the input to the  $(d-1)$ -th layer:

$$\begin{aligned}
 g_{d-2} &= \frac{\partial (W_d \cdot W_{d-1} \cdot z_{d-1} - y)^2}{\partial z_{d-1}} \\
 &= 2 W_{d-1}^T W_d^T (W_d \cdot W_{d-1} \cdot z_{d-1} - y) \\
 &= W_{d-1}^T \cdot g_{d-1} \\
 &\vdots \\
 g_i &= W_{i+1}^T \cdot g_{i+1}
 \end{aligned}$$

Put together.

Observe that  $\frac{\partial L}{\partial w_i} = g_i \cdot a_i^T$ !

Therefore, our full algorithm is :

---

### Backpropagation for linear neural net

IN :  $(x_1, y_1), \dots, (x_m, y_m)$

Req :  $(w_1^{(t)}, \dots, w_d^{(t)})$ ,  $\eta$

OUT :  $(w_1^{(t+1)}, \dots, w_d^{(t+1)})$

---

Forward pass :

compute  $a_1, a_2, \dots, a_d$

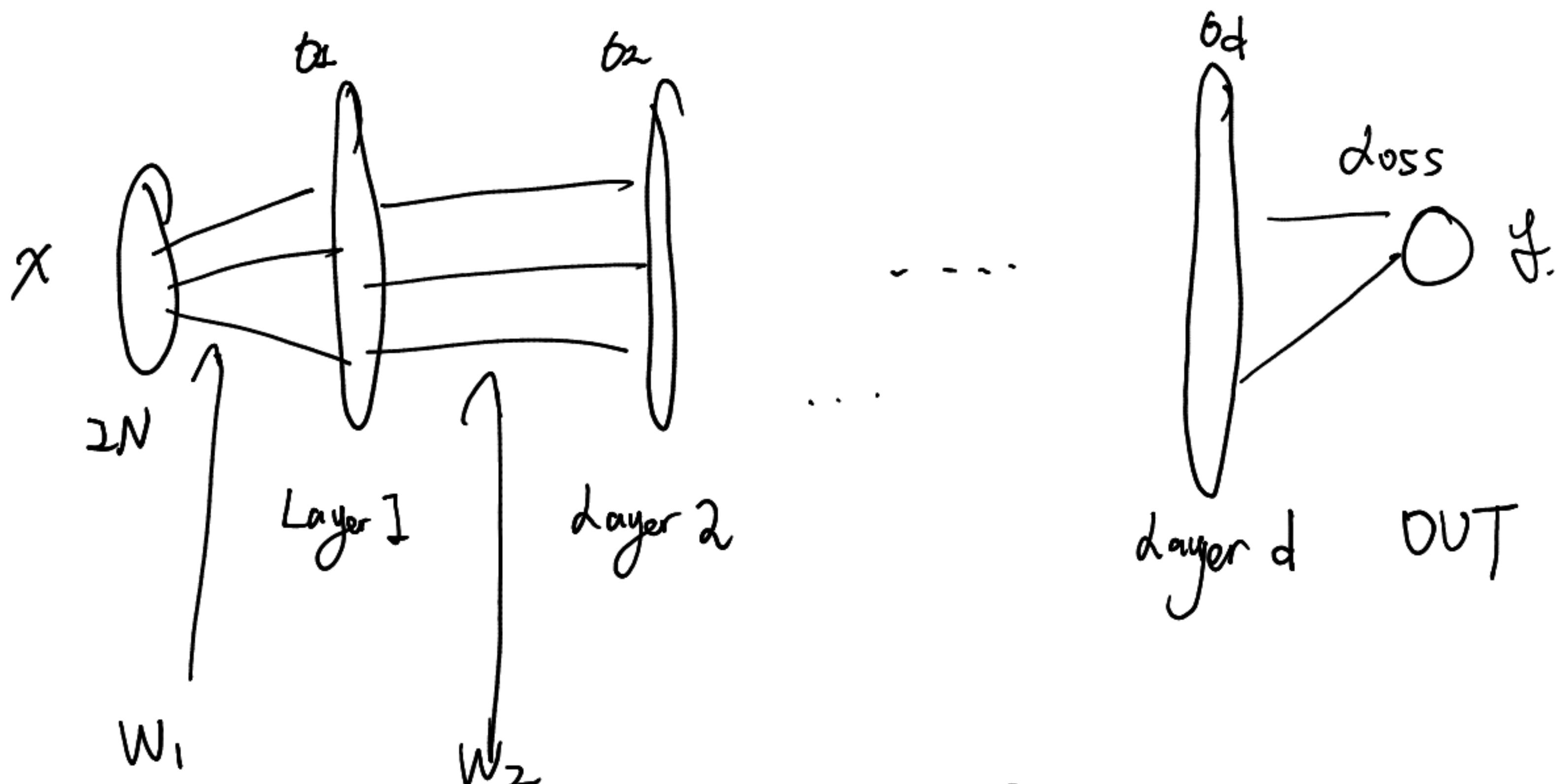
Backward pass :

compute  $g_d, g_{d-1}, \dots, g_1$

update  $w_i^{(t+1)} \leftarrow w_i^{(t)} - \eta \cdot g_i \cdot a_i^T$

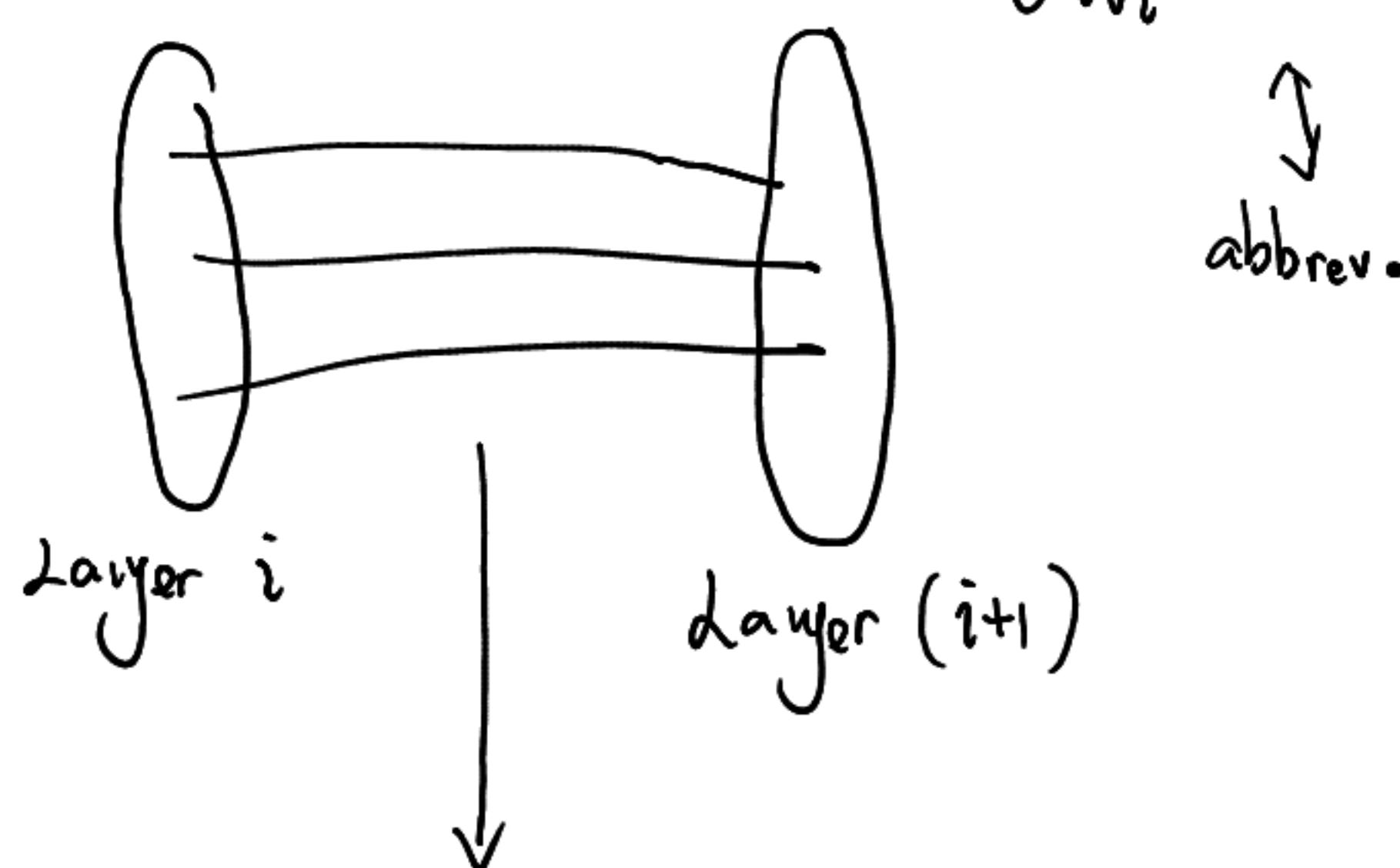
The general case. Non-Linear activation, arbitrary loss function, bias in activation function.

Model.



Chain rule. We would like to calculate

$$\frac{\partial f(w_1, w_2, \dots, w_d)}{\partial w_i} ?$$



Weight matrix:  $W_i \in \mathbb{R}^{r_i \times r_{i+1}}$

Output of layer  $i$ :  $a_i \in \mathbb{R}^{r_i}$

Input to layer  $i+1$ :  $z_i \in \mathbb{R}^{r_{i+1}}$

Claim 1.

$$\frac{\partial f}{\partial w_i} = \frac{\partial f}{\partial z_i} \cdot a_i^T$$

Proof. Chain rule :  $\frac{\partial f}{\partial w_i} = \frac{\partial f}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_i} = \frac{\partial f}{\partial z_i} \cdot a_i^T$ .

By Def.  $z_i = w_i \cdot a_i = \langle a_i^T, w_i \rangle$ , hence  $\frac{\partial z_i}{\partial w_i} = a_i^T$ .

Claim 2.  $\frac{\partial f}{\partial z_i} = (w_{i+1}^T \cdot \frac{\partial f}{\partial z_{i+1}}) \odot b_{i+1}'(z_i)$ .

Notation : for two vectors  $x, y \in \mathbb{R}^n$ , their Kronecker product

$$x \odot y = (x_1 y_1, x_2 y_2, \dots, x_n y_n) \in \mathbb{R}^n.$$

- for a function  $t: \mathbb{R} \rightarrow \mathbb{R}$ , we use the vector notation  $t(x) = (t(x_1), t(x_2), \dots, t(x_n))$ , and  $t'(x) = (t'(x_1), t'(x_2), \dots, t'(x_n))$ .

Proof. By def.  $\nabla f = \langle \frac{\partial f}{\partial z_i}, \nabla z_i \rangle$ , for some  $\nabla z_i$  that goes to zero.

So, suppose we perturb  $z_i$  by  $\delta$ , we would like to know how much is  $f$  changed. Note that

$$\nabla f = \langle \frac{\partial f}{\partial z_{i+1}}, \nabla z_{i+1} \rangle.$$

So we are going to track how much  $z_{i+1}$  changes.

$$z_{i+1} = w_{i+1} \cdot a_{i+1}$$

$$= w_{i+1} \cdot b_{i+1}'(z_i + \delta)$$

$$= w_{i+1} \cdot (b_{i+1}'(z_i) + b_{i+1}''(z_i) \odot \delta)$$

$$= W_{\bar{i}+1} \cdot (b'_{\bar{i}+1}(z_i) + b'_{\bar{i}+1}(z_i) \odot \delta)$$

$$= W_{\bar{i}+1} \cdot b_{\bar{i}+1}(z_i) + W_{\bar{i}+1} \cdot (b'_{\bar{i}+1}(z_i) \odot \delta)$$

$$\Rightarrow \nabla z_{\bar{i}+1} = W_{\bar{i}+1} \cdot (b'_{\bar{i}+1}(z_i) \odot \delta). \quad \text{Therefore}$$

$$\nabla f = \left\langle \frac{\partial f}{\partial z_{\bar{i}+1}}, W_{\bar{i}+1} \cdot (b'_{\bar{i}+1}(z_i) \odot \delta) \right\rangle$$

$$\left\langle \frac{\partial f}{\partial z_i}, \delta \right\rangle = \text{Tr} \left[ \left( \frac{\partial f}{\partial z_{\bar{i}+1}} \right)^T \cdot W_{\bar{i}+1} \cdot (b'_{\bar{i}+1}(z_i) \odot \delta) \right]$$

$$= \left\langle W_{\bar{i}+1}^T \cdot \frac{\partial f}{\partial z_{\bar{i}+1}}, b'_{\bar{i}+1}(z_i) \odot \delta \right\rangle.$$

$$\Rightarrow \frac{\partial f}{\partial z_i} = \left( W_{\bar{i}+1}^T \cdot \frac{\partial f}{\partial z_{\bar{i}+1}} \right) \odot b'_{\bar{i}+1}(z_i).$$

**Remark.** Suppose that  $b_{\bar{i}+1}(z_i) = z_i$  (i.e. linear activation),

then  $b'_{\bar{i}+1}(z_i) = 1$ , therefore

$$\frac{\partial f}{\partial z_i} = W_{\bar{i}+1}^T \cdot \frac{\partial f}{\partial z_{\bar{i}+1}}. \quad \text{Combining Claim 1 and 2,}$$

we have shown that

$$\frac{\partial f}{\partial w_i} = W_{\bar{i}+1}^T \cdot \frac{\partial f}{\partial z_{\bar{i}+1}} \cdot a_i^T, \quad \text{which recovers our earlier result!}$$

Put together, we have that

$$\frac{\partial f}{\partial w_i} = \frac{\partial f}{\partial z_i} \cdot a_i^T,$$

$$\frac{\partial f}{\partial z_i} = \left( w_{i+1}^T \cdot \frac{\partial f}{\partial z_{i+1}} \right) \cdot b'_{i+1}(z_i),$$

Goal: Compute  $\{a_i\}$ ,  $\left\{ \frac{\partial f}{\partial z_i} \right\}$ ,  $\left\{ b'_{i+1}(z_i) \right\}$ .

Recall that  $z_i$  is the input to the  $(i+1)$ -th layer,  
and  $a_i$  is the output of the  $i$ -th layer.

Therefore,

$$z_i = w_i \cdot a_i,$$

Forward Pass  $\{a_i\}$ .  $a_0 = x$

$$\begin{aligned} a_{i+1} &= b(z_{i+1}) \\ &= b(w_i \cdot a_i). \end{aligned}$$

Backward Pass.  $\left\{ \frac{\partial f}{\partial z_i} \right\}$ .  $\frac{\partial f}{\partial z_d} = \frac{\partial \text{loss}(b_{d+1}(z_d), y^{\text{true}})}{\partial z_d}$ .

$$\frac{\partial f}{\partial z_i} = \left( w_{i+1}^T \cdot \frac{\partial f}{\partial z_{i+1}} \right) \cdot b'_{i+1}(z_i).$$

Put together, we have derived the backpropagation algorithm.

---

## Backpropagation for feedforward neural networks

---

Forward Pass:

Compute  $\{a_i\}$ ,  $\{z_i\}$ .

Backward Pass:

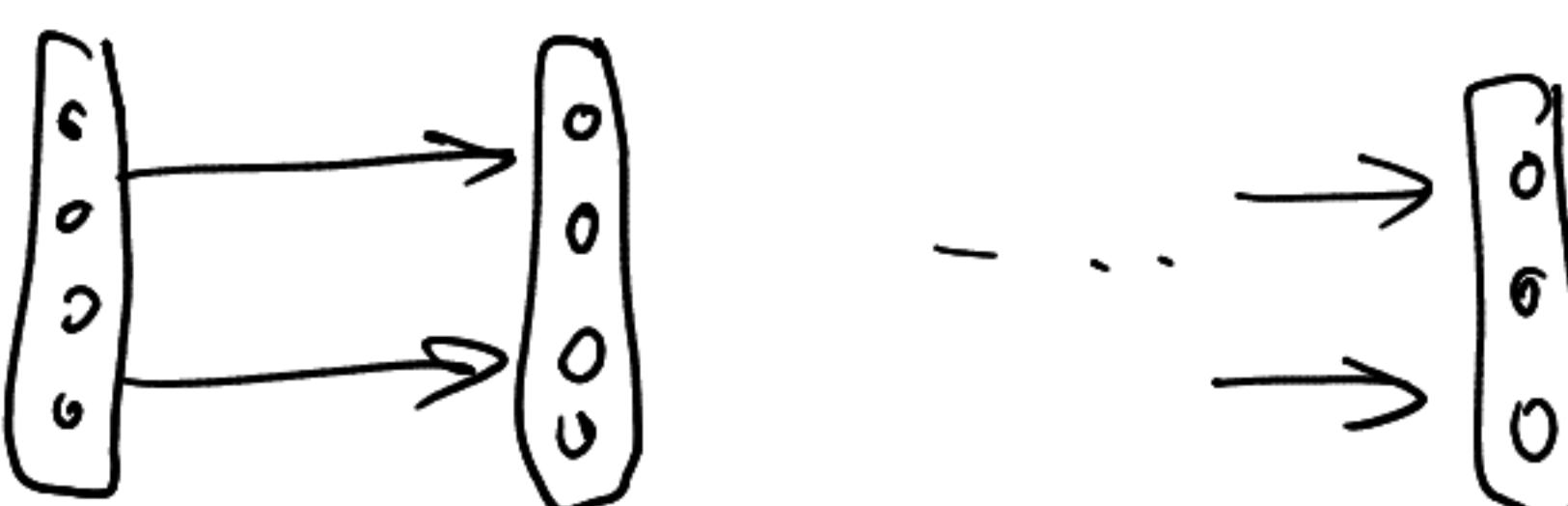
Compute  $\left\{ \frac{\partial f}{\partial z_{ij}} \right\}$ , hence  $\left\{ \frac{\partial f}{\partial w_i} \right\}$ .

Apply gradient descent over  $\{w_i\}$ .

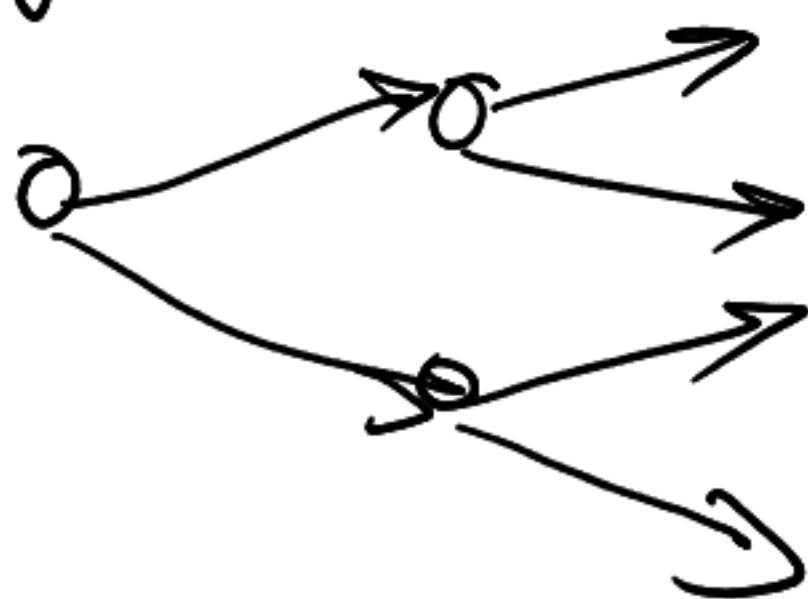
---

### Extensions.

- Acyclic neural network

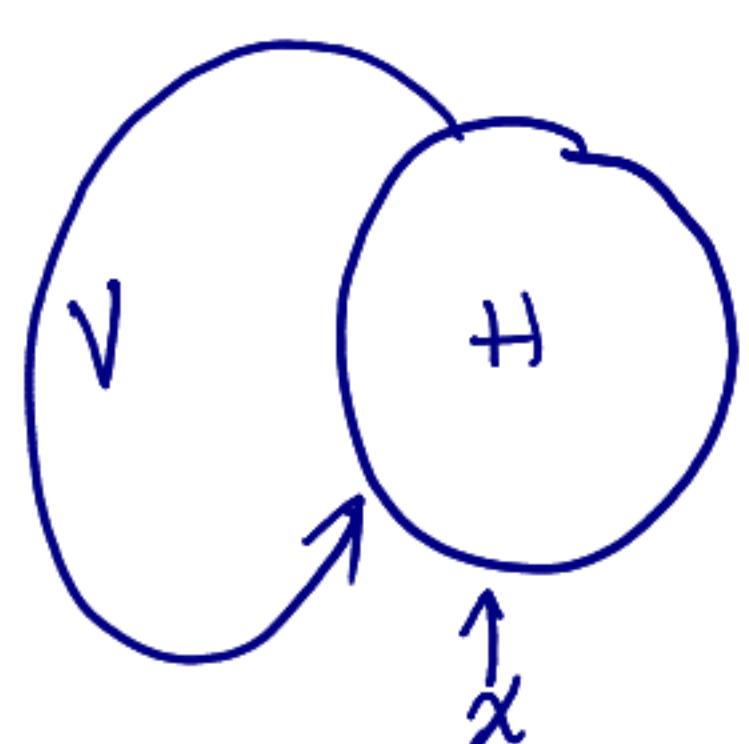


In general, the algorithm works as long as there are no cycles in the network!



What if there are cycles? e.g. RNN? yes but we may run into issues.

Exploding or vanishing gradients.

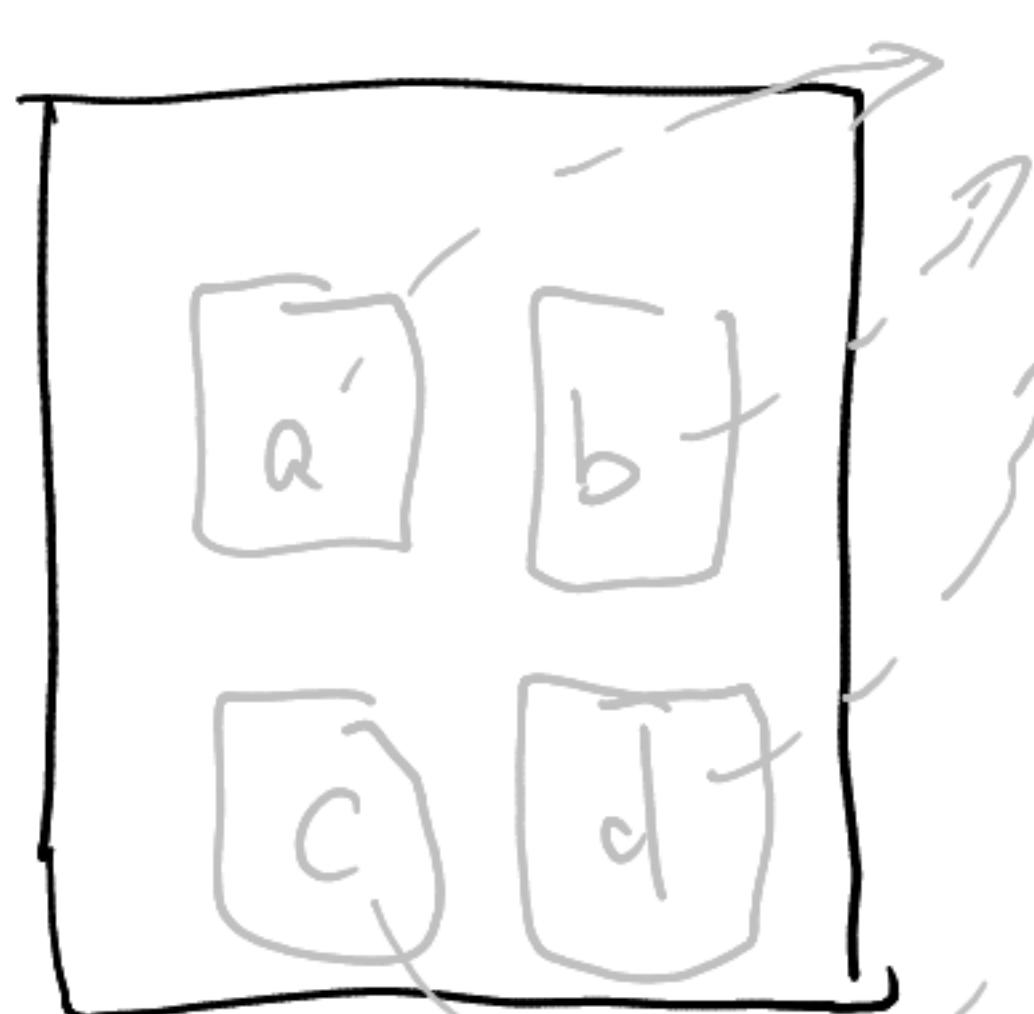


- i) if  $\|V\|_2 < 1$ , then  $\|V^i \cdot x\| < \sigma_{\max}(V) \cdot \|x\|$
- ii) if  $\|V\|_2 > 1$ , then  $\|V^i \cdot x\| > \sigma_{\min}(V) \cdot \|x\|$

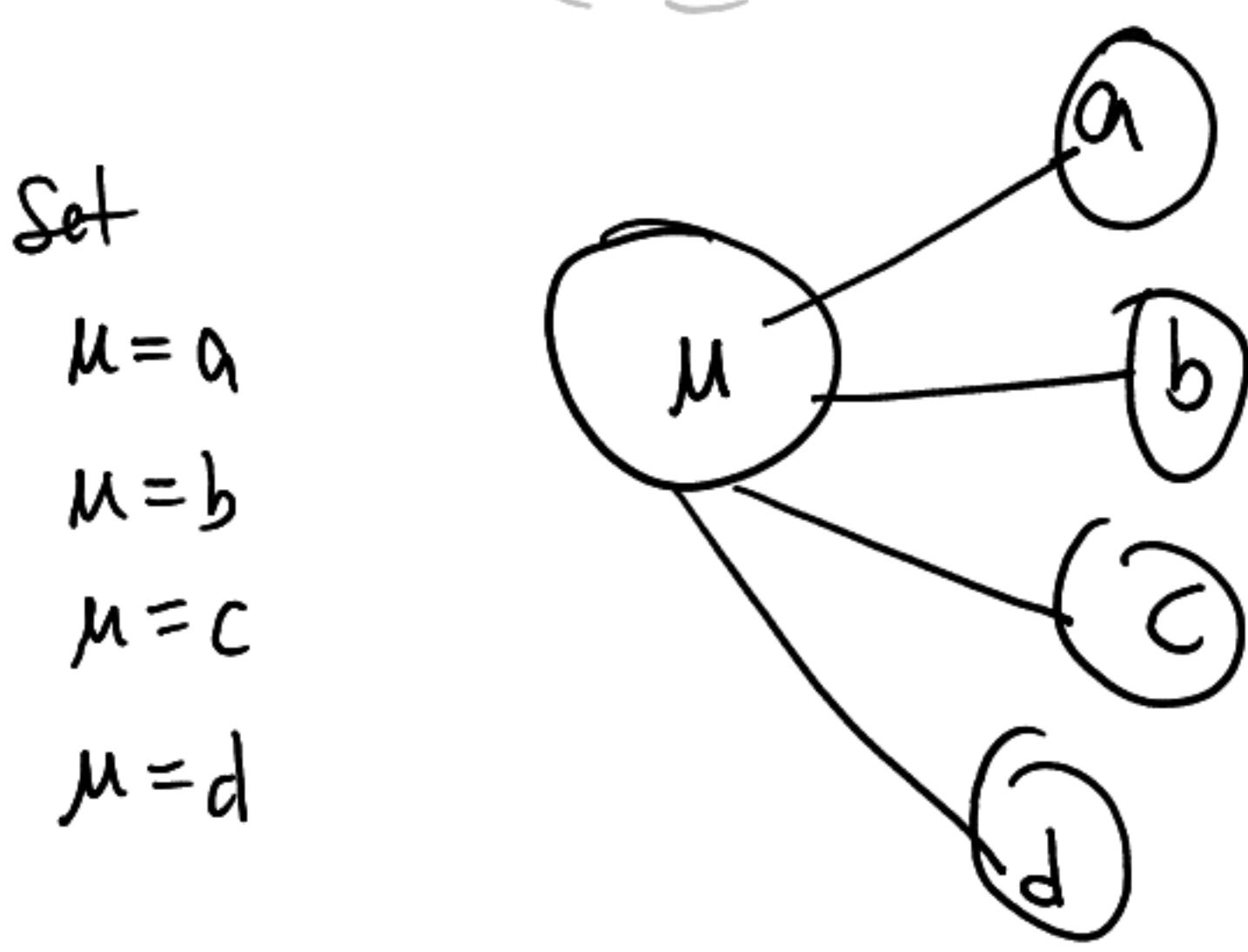
Fix. Gradient clipping on long short-term memory.

Fix. Backpropagation in time. Expand over the prev. k steps.

- Weight sharing.



multiple patches share the same weight matrix W.



Set

$$\mu = a$$

$$\mu = b$$

$$\mu = c$$

$$\mu = d$$

$$\begin{aligned}\frac{\partial f}{\partial \mu} &= \frac{\partial f}{\partial a} \cdot \frac{\partial a}{\partial \mu} + \dots \\ &= \frac{\partial f}{\partial a} + \frac{\partial f}{\partial b}.\end{aligned}$$

## Second-order methods:

- Hessian :  $\frac{\partial^2 f}{\partial w_i \partial w_j}$
- Traditional use of Hessian
- Hessian in neural nets

## Conclusions:

- We derived the Backpropagation algorithm, which involves a forward pass and a backward pass.
- Backpropagation provides a linear-time implementation of gradient descent.
  - & Key idea: Chain rule in calculus.
- Extensions: Acyclic networks, Weight sharing,
- Discussion: second-order methods.