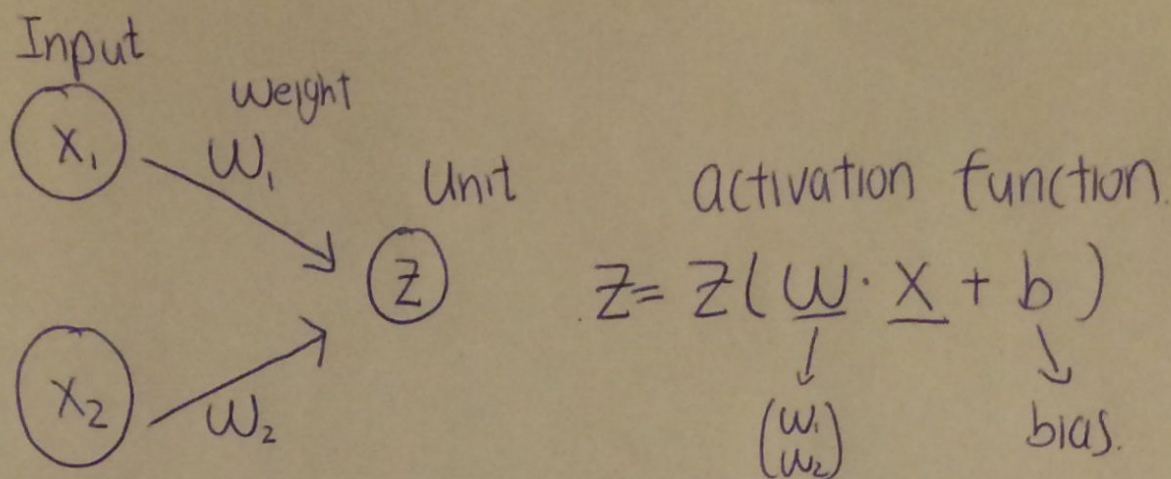


①

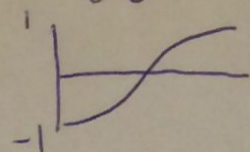


Possible act. functions.

$z(t) = \text{rectified linear}$ $\max(0, t)$

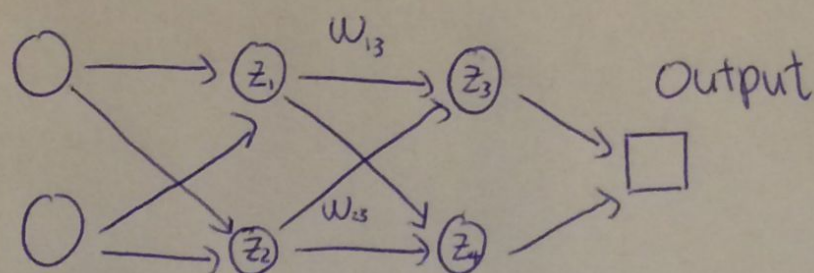
Sigmoid $\frac{1}{1+e^{-t}}$

\tanh $\frac{e^{2x} - 1}{e^{2x} + 1}$



$$\tanh(x) = 2 \text{Sig}(2x) - 1$$

②



Recursively define $z_3 = Z(\begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \cdot \begin{pmatrix} w_{13} \\ w_{23} \end{pmatrix} + b_3)$

Let \underline{z}^L denote the last layer of nodes
and \underline{w}^L denote the last layer of weights

Output: \square

◦ linear regression

$$y = \underline{w}^L \cdot \underline{z}^L + b^L$$

◦ Classification

$$y_i = \exp[w^{L,i} \cdot z^L] / \sum_j \exp[w^{L,j} \cdot z^L]$$

⑤ Non-linear function approximation.

$$y = f^*(x) + \text{noise}, \quad f^* \text{ is the "true" mapping}$$

$$y = f(x; \underbrace{w, b}_{\text{param. } \theta}).$$

$$\text{Cost function } C = \frac{1}{n} \sum_x [y(x) - f(x; \theta)]^2$$

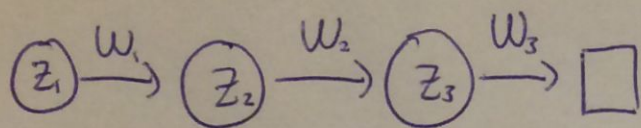
Learning θ

• Init θ_0 .

• Iterate $\theta_t = \theta_{t-1} - \alpha_t \underbrace{\nabla_{\theta} C}_{\text{gradient}}$

$$-\frac{2}{n} \sum_x [y(x) - f(x; \theta)] \times \begin{pmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \\ \vdots \end{pmatrix}$$

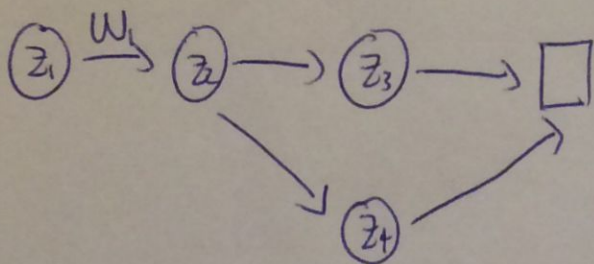
Example.



$$\text{out} = w_3 \cdot z_3$$

$$f = w_3 \cdot z_3(w_2 \cdot z_2(w_1 \cdot z_1))$$

$$\frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial z_3} \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_1}$$



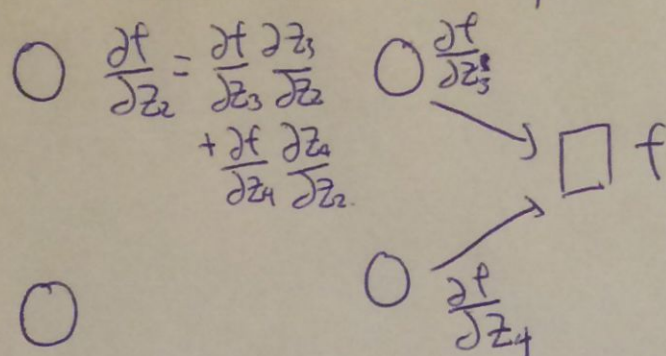
$$\begin{aligned} \frac{\partial f}{\partial w_1} &= \frac{\partial f}{\partial z_3} \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_1} \\ &+ \frac{\partial f}{\partial z_4} \frac{\partial z_4}{\partial z_2} \frac{\partial z_2}{\partial w_1} \end{aligned}$$

2 path of influence
from $\textcircled{w_1}$ to output \square .

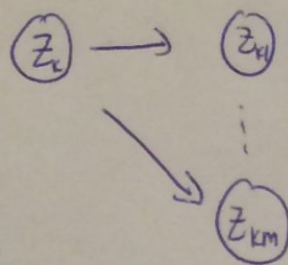
④ Backprop = chain rule + dynamic programming
Reuse computation.

• If $(z_i) \rightarrow (z_j)$, compute $\frac{\partial z_j}{\partial z_i}$

• Start from the last layer and go backward.



In general, If



$$\frac{\partial f}{\partial z_k} = \sum_{j=1}^m \frac{\partial f}{\partial z_{kj}} \frac{\partial z_{kj}}{\partial z_k}$$

• Now we have all $\frac{\partial f}{\partial z}$

If $w \rightarrow (z) \quad \frac{\partial f}{\partial w} = \frac{\partial f}{\partial z} \frac{\partial z}{\partial w}$

Cost of computing $\nabla_{\theta} \mathcal{L}$ is $O(\# \text{ of edges})$.