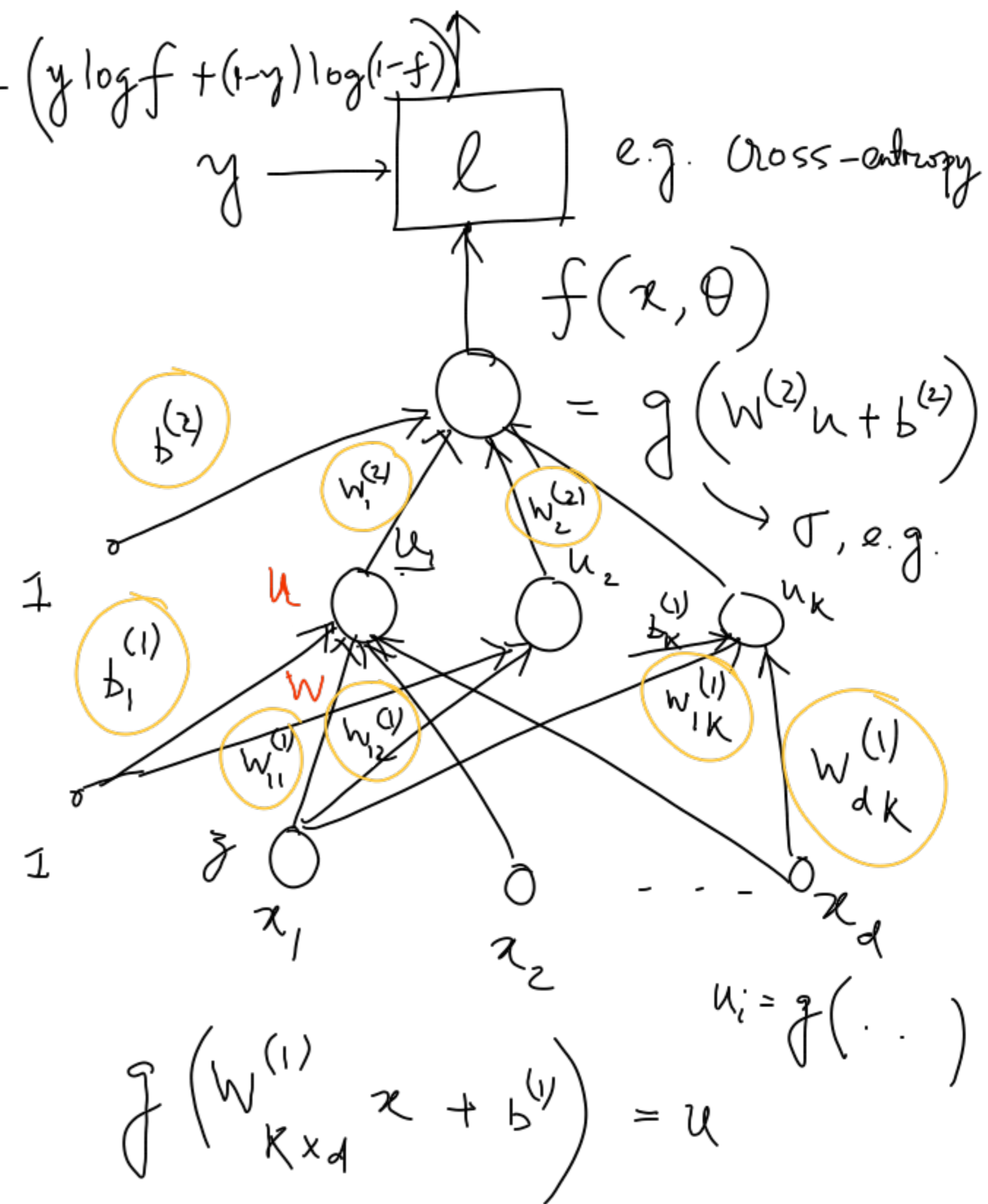


Lec 12: Neural Network Training Algorithm (Feed forward)

Input: $NN(x; \theta)$, training data $(x_i, y_i)_{i=1}^n - (y \log f + (1-y) \log(1-f))$
 loss function l e.g. cross-entropy

- Randomly initialize θ
 - Repeat until stopping criteria is met:
 - Pick randomly a batch of examples $(x_i, y_i)_{i \in B}$ - $B \subset N, |B| = T < n$
 - Compute the gradient of the loss function $\nabla_{\theta} l$ at θ
 - update $\theta \leftarrow \theta - \eta \nabla_{\theta} l$
 - Return θ
- ↑ hyperparameter

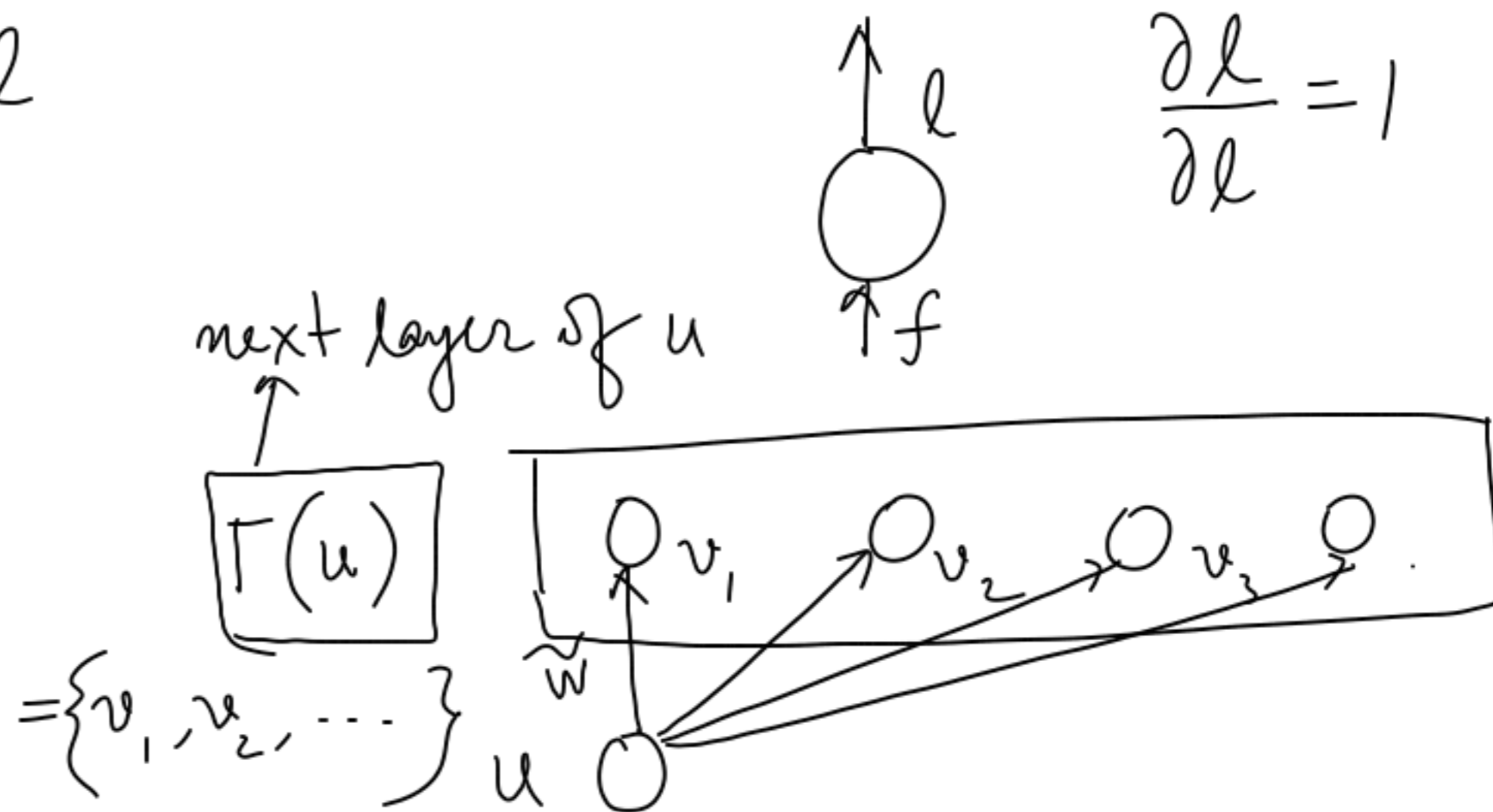


Goal: learn weights w, b that minimize \mathcal{L}

Gradient descent:

$$W_{ij}^{(k)} \leftarrow W_{ij}^{(k)} - \eta \frac{\partial \mathcal{L}}{\partial W_{ij}^{(k)}}$$

$\frac{\partial \mathcal{L}}{\partial W_{ij}^{(k)}}$ ← layer
 ij ← output node
 ← input node of that layer



Q: How to efficiently compute the gradients?

① Compute $\frac{\partial \mathcal{L}}{\partial u}$ for every node u in the NN $\Theta = (W, b)$
 $\Theta \leftarrow$ current Θ

② Compute $\frac{\partial \mathcal{L}}{\partial w} = \left(\frac{\partial \mathcal{L}}{\partial u} \right) \frac{\partial u}{\partial w}$

$$\frac{\partial \mathcal{L}}{\partial u} = \sum_{v \in T(u)} \frac{\partial \mathcal{L}}{\partial v} \frac{\partial v}{\partial u}$$

$$u = g(wz + \dots)$$

$$\frac{\partial u}{\partial w} = g' \cdot z$$

$$v = g(\tilde{w}u + \dots)$$

Backpropagation: two stage algorithm

Forward pass: Calculate the value of each node given an input (x_i, y_i) at θ

$$u_i = g \left(w_{1i}^{(k)} z_1 + w_{2i}^{(k)} z_2 + \dots + b_i^{(k)} \right)$$

Backward pass: Base case $\frac{\partial \mathcal{L}}{\partial \mathcal{L}} = 1$

For each u in a given layer:

For each $v \in T(u)$

- use already computed $\frac{\partial \mathcal{L}}{\partial v}$

- compute $\frac{\partial v}{\partial u}$ | x_i, y_i, θ

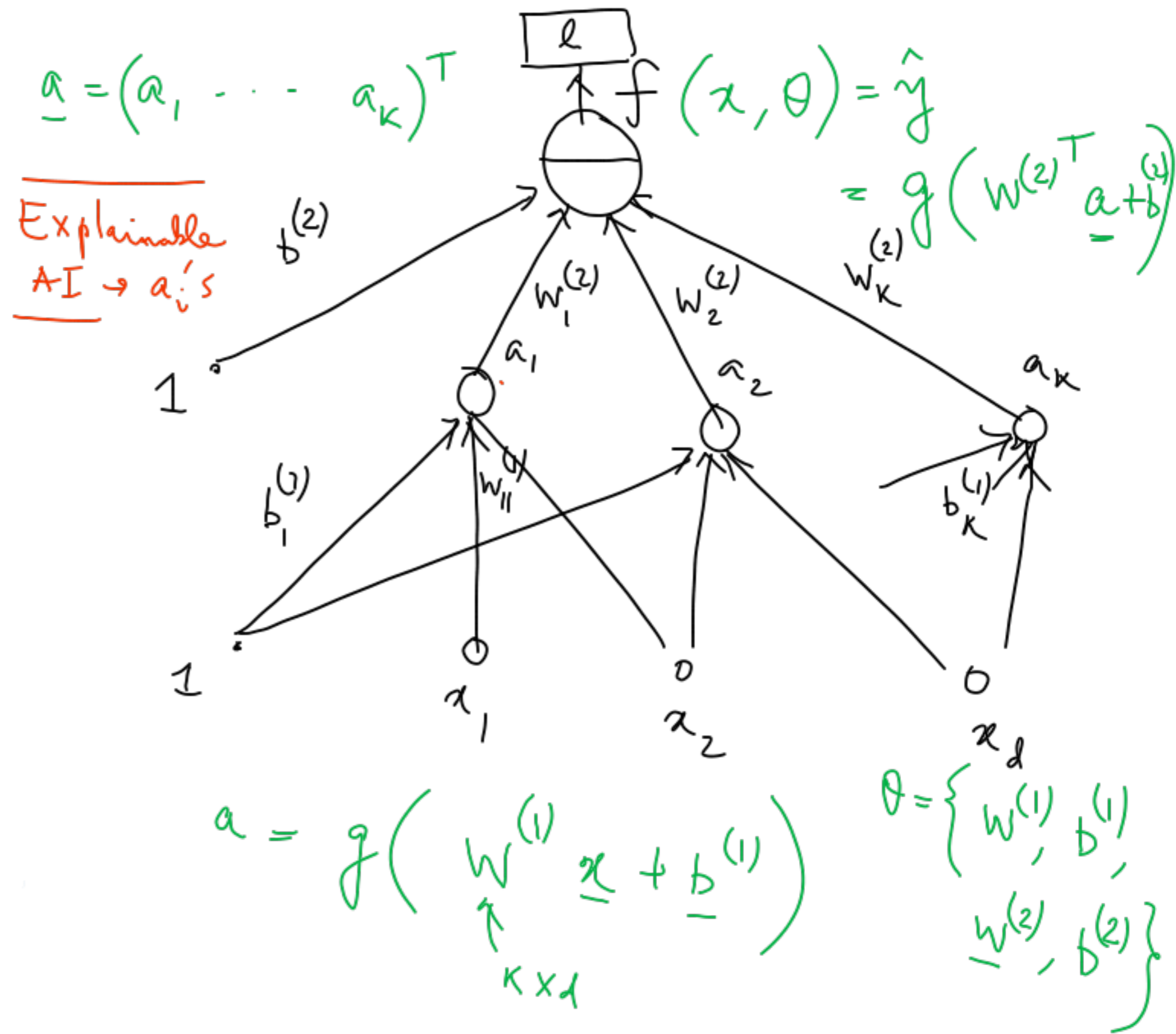
- $\frac{\partial \mathcal{L}}{\partial v} \cdot \frac{\partial v}{\partial u}$

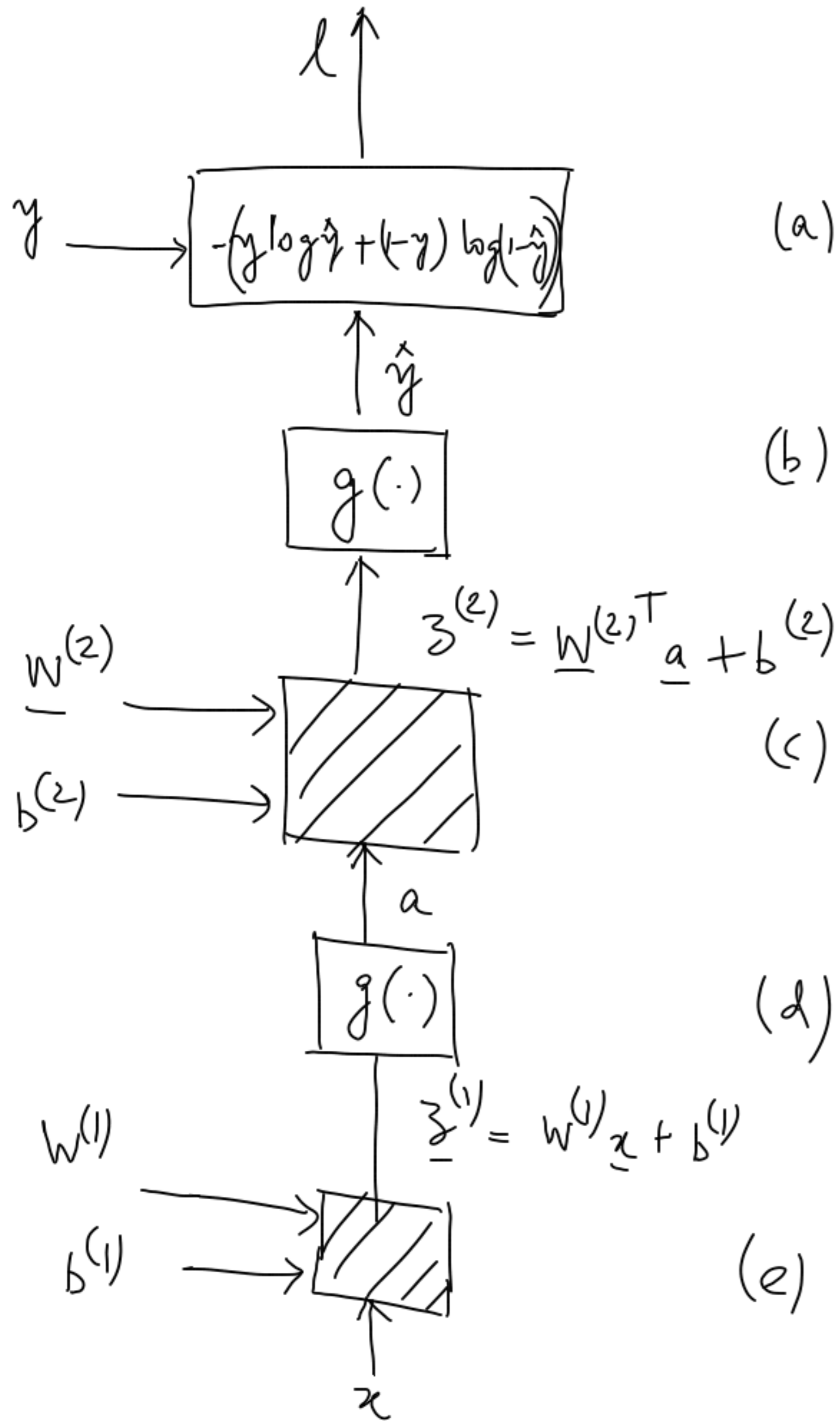
get $\frac{\partial \mathcal{L}}{\partial u}$, compute $\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial u} \cdot \frac{\partial u}{\partial w} \Rightarrow \frac{\partial \mathcal{L}}{\partial b}$

One example in the webpage:

Typical these derivatives are more efficient to compute in vector form.

Computational graph





(a) \rightarrow

$$\frac{\partial L}{\partial \hat{y}} = - \frac{\partial}{\partial \hat{y}} \left[y \log \hat{y} + (1-y) \log(1-\hat{y}) \right] = - \frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$$

$$\delta_{\hat{y}} = \frac{\partial L}{\partial \hat{y}} \Big|_{x, y, \theta}$$

(b) \rightarrow

$$\frac{\partial \hat{y}}{\partial z^{(2)}} = g'(z^{(2)}), \quad \delta_{z^{(2)}} = \frac{\partial L}{\partial z^{(2)}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z^{(2)}}$$

(c) \rightarrow

$$\frac{\partial z^{(2)}}{\partial \underline{w}^{(2)}} = \frac{\partial}{\partial \underline{w}^{(2)}} \left[\underline{w}^{(2)T} \underline{a} + b^{(2)} \right] = \underline{a} \quad \delta_{\underline{w}^{(2)}} = \frac{\partial L}{\partial \underline{w}^{(2)}} = \delta_{\hat{y}} \cdot \frac{\partial \hat{y}}{\partial z^{(2)}}$$

(d) \rightarrow

use chain rule.

(e) \rightarrow

$$\frac{\partial z^{(1)}}{\partial w^{(1)}} = \frac{\partial}{\partial w^{(1)}} \left[w^{(1)} x + b^{(1)} \right] = x$$

\rightarrow HW

$\delta_{w^{(1)}} \rightarrow$

Regularization:

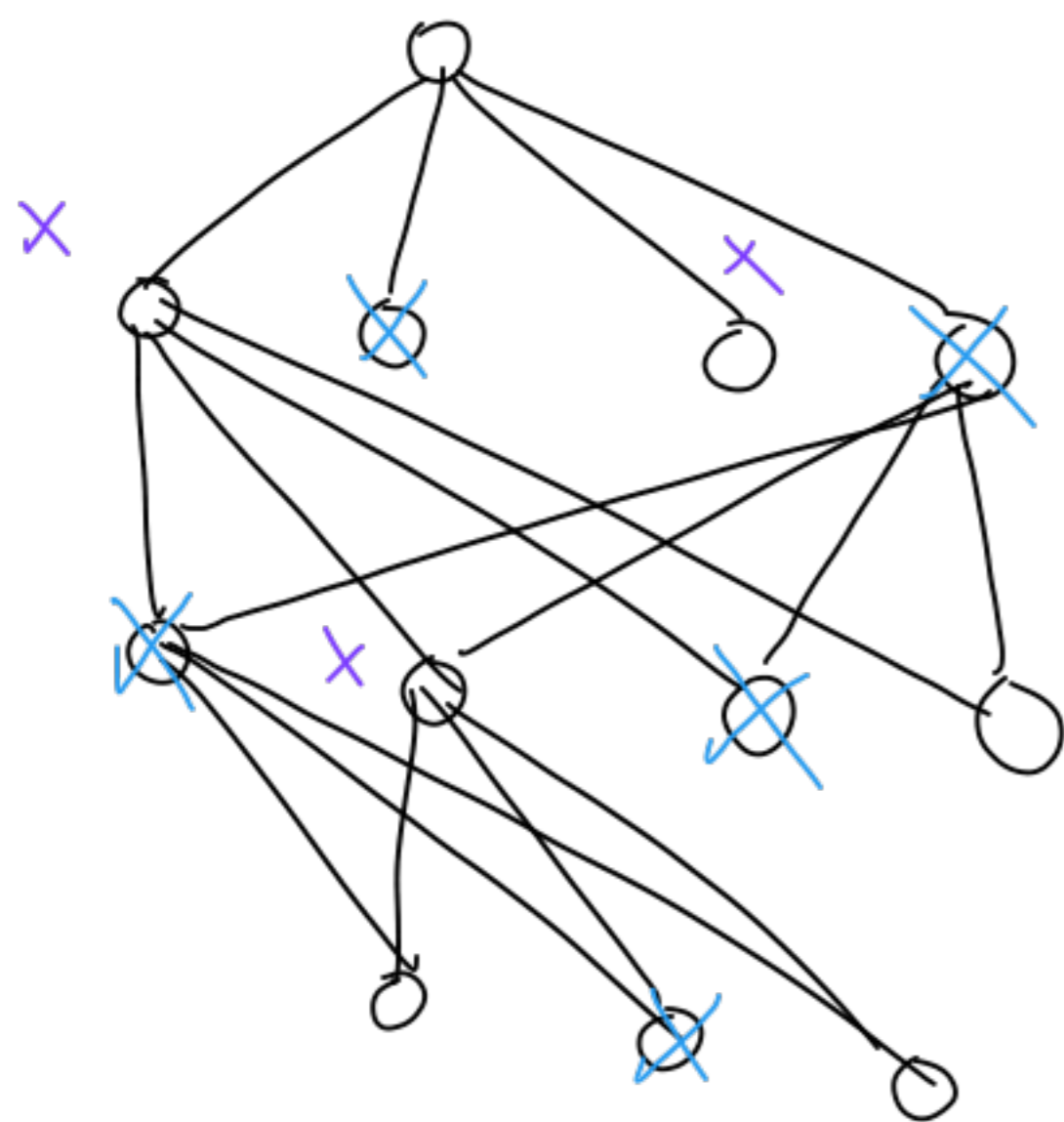
o L_2 regularization:

$$l_{\text{Reg}} = l_{\text{CE}} + \lambda \underbrace{\|W\|_2^2}_{\sum_{ijk} (W_{ij}^{(k)})^2}$$

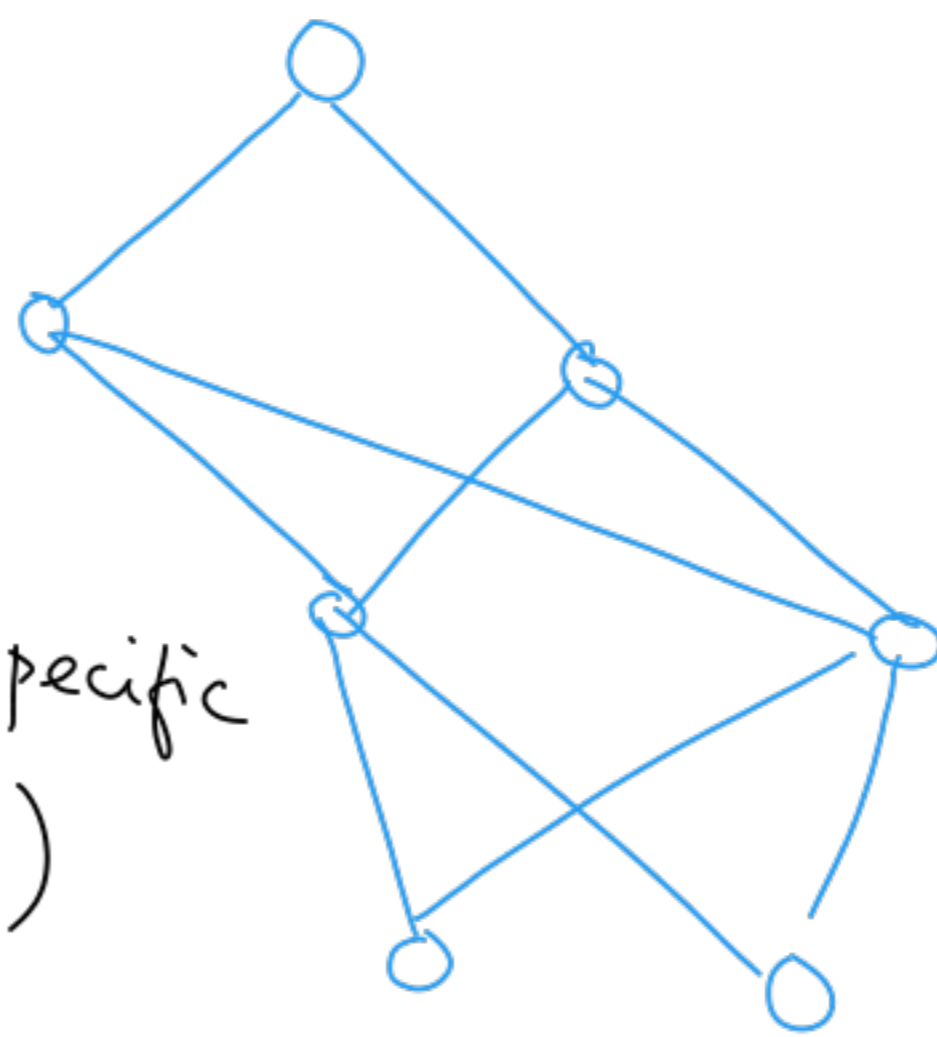
o Dropout regularization:

during training, keep a neuron active w.p. p ,
delete otherwise

↓ hyperparameter



before dropout



for a specific
 (x_i, y_i)

after dropout

expected weights as an output

at test time → use the
entire network with expected
weights.

- Early stopping: Stopping when performance on validation set stops improving.

