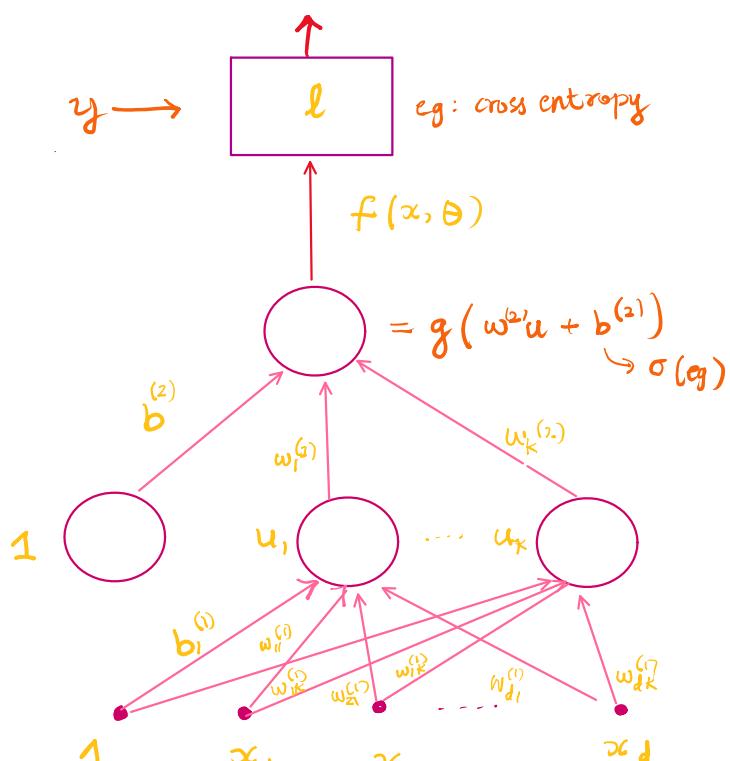


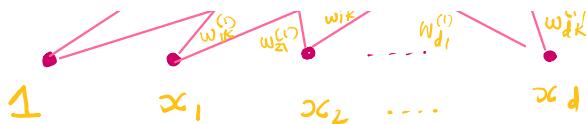
The neural networks we are addressing so far is "fit forward"

### NN training algorithm

Input: NN( $x; \theta$ ) , training data  $(x_i, y_i)_{i=1}^n$  , loss function,  $l$ .

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





$$g\left(\sum_{k=1}^d w_k^{(l)} x_k + b^{(l)}\right) = u \rightarrow u_i = g\left(\sum_{j=1}^d w_{ij}^{(l)} x_j + b_i^{(l)}\right)$$

Goal: learn the weights  $w, b$  that minimises  $\ell$ .

Gradient descent:

$$\frac{\partial \ell}{\partial w_{ij}^{(k)}} \leftarrow \begin{array}{l} \text{large} \\ \text{output node} \\ \text{input node of that layer} \end{array}$$

$$w_{ij}^{(k)} \leftarrow w_{ij}^{(k)} - \eta \frac{\partial \ell}{\partial w_{ij}^{(k)}}$$

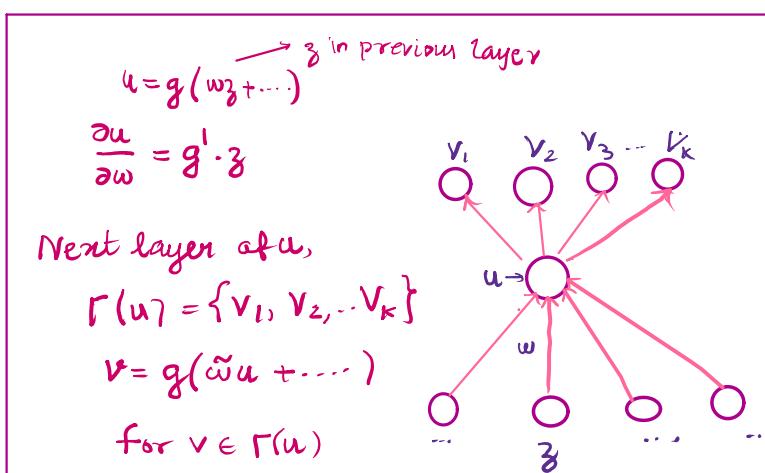
How will we compute these gradient values efficiently?

Back Propagation:

1. Compute  $\frac{\partial \ell}{\partial u} \Big|_{\theta \leftarrow \text{current } \theta}$  for every node  $u$  in the NN  
 $\theta = (w, b)$

2. Compute  $\frac{\partial \ell}{\partial w} = \frac{\partial \ell}{\partial u} \cdot \frac{\partial u}{\partial w}$

$$\frac{\partial \ell}{\partial u} = \sum_{v \in \Gamma(u)} \frac{\partial \ell}{\partial v} \cdot \frac{\partial v}{\partial u}$$



Back propagation algorithm

Two stages:

(a) **Forward Pass:** Calculate the value of each node given an input  $(x_i, y_i)$  at  $\theta$

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

(b) Backward Pass: Base case:  $\frac{\partial l}{\partial z} = 1$

For each  $u$  in a given layer:

For each  $v \in \Gamma(u)$ :

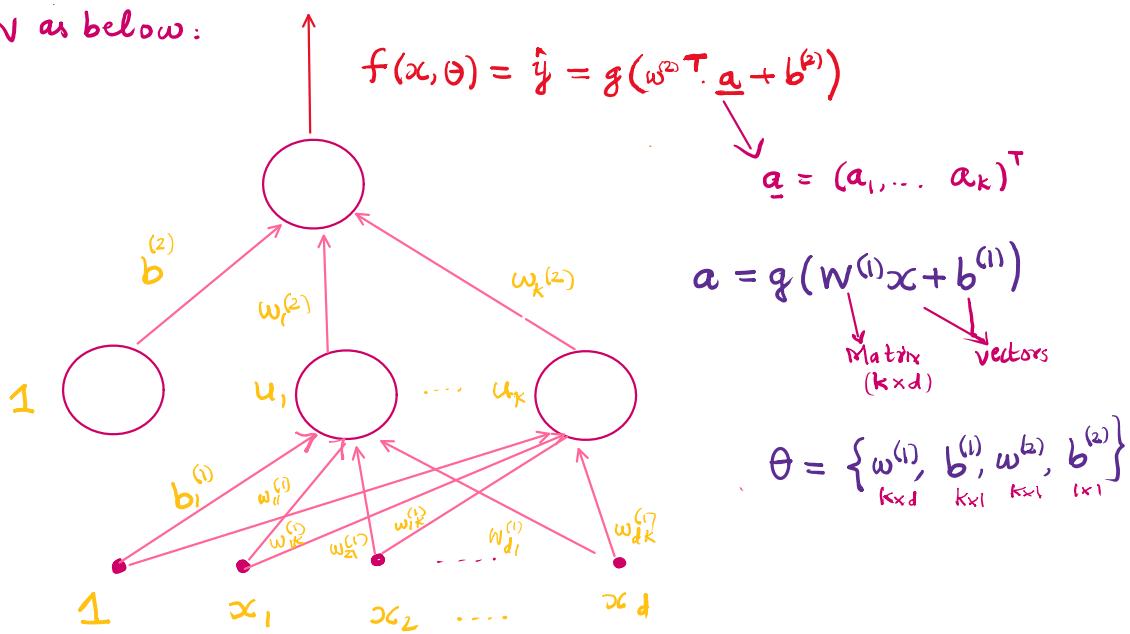
- Use already computed  $\frac{\partial l}{\partial v}$  (known by induction hypothesis)
- compute  $\frac{\partial v}{\partial u}$
- $\frac{\partial l}{\partial v} \cdot \frac{\partial v}{\partial u}$  (sum this)

$$\text{get } \frac{\partial l}{\partial u} \left( \sum_{v \in \Gamma(u)} \frac{\partial l}{\partial v} \cdot \frac{\partial v}{\partial u} \right)$$

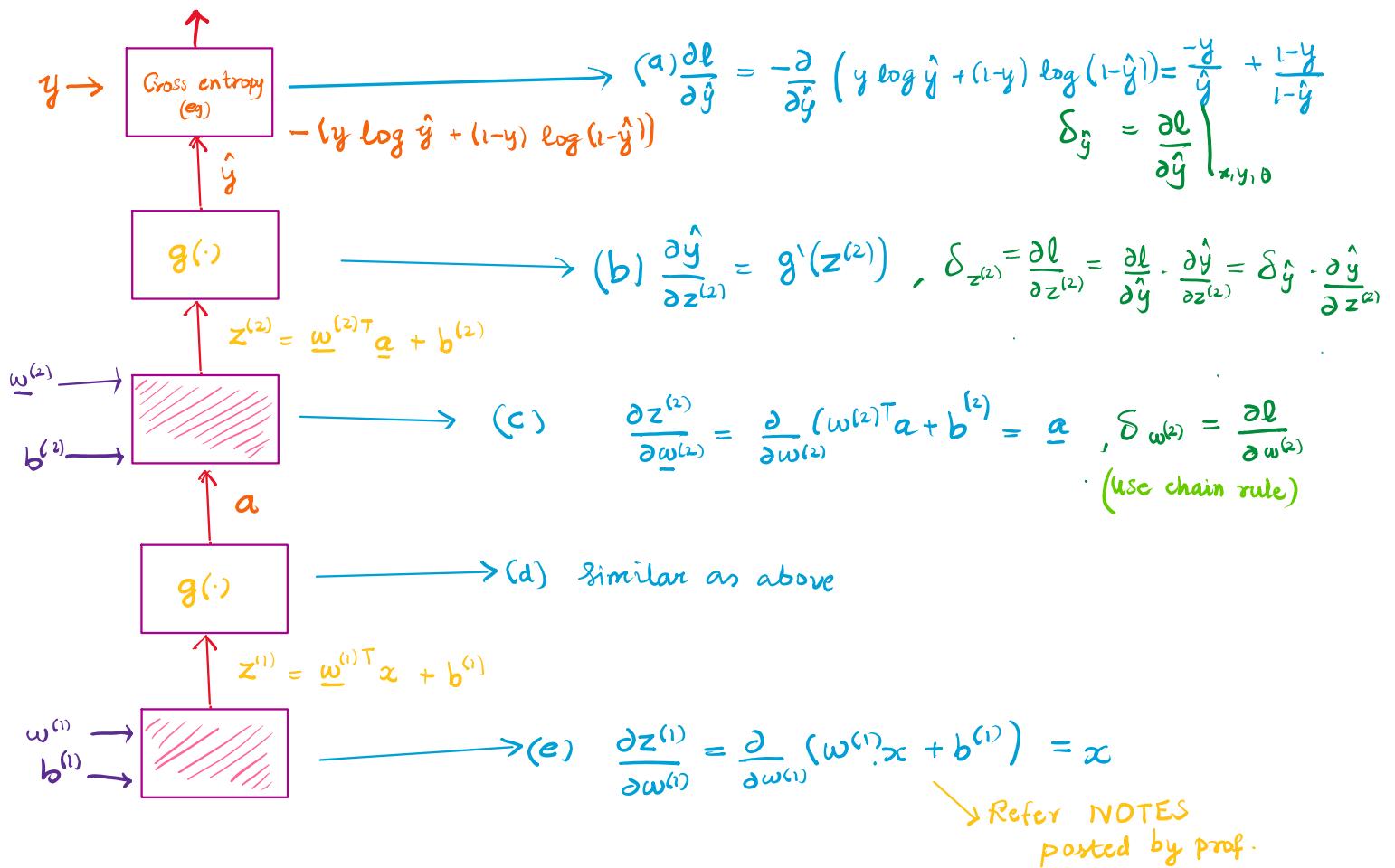
$$\text{Compute } \frac{\partial l}{\partial w} = \frac{\partial l}{\partial u} \cdot \frac{\partial u}{\partial w}, \quad \frac{\partial l}{\partial b}$$

Typically these derivatives are more efficient to compute in vector form. We use a Computational graph for this.

Consider an NN as below:



The computational graph for above NN is as below:



**Explainable AI:** Tries to understand what the weights,  $a_i$ 's in intermediate layers mean.

## Regularization in NN:

To prevent overfitting in NN.

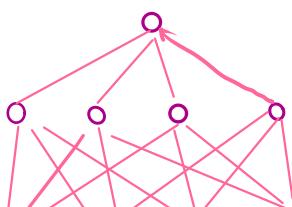
The commonly used regularization functions:

$$\text{i) L}_2 \text{ regularization: } l_{\text{Reg}} = l_{\text{CE}} + \lambda \|w\|_2^2$$

$\sum_{ij,k} (w_{ij}^{(k)})^2$

ii) Drop out regularization:

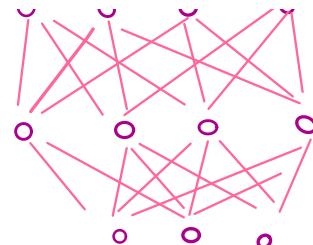
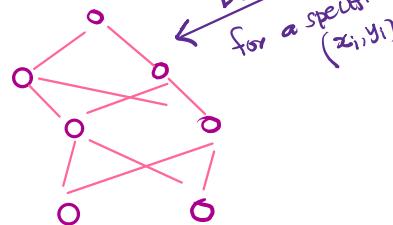
During training keep a neuron active with probability  $P$ , delete otherwise



During training keep a neuron active  
with probability  $P$ , delete otherwise

hyper parameter

Expected weight  
as an output.



During test time, use the entire network with expected weights

### (iii) Early stopping:

Stopping when performance on validation set stops improving.

