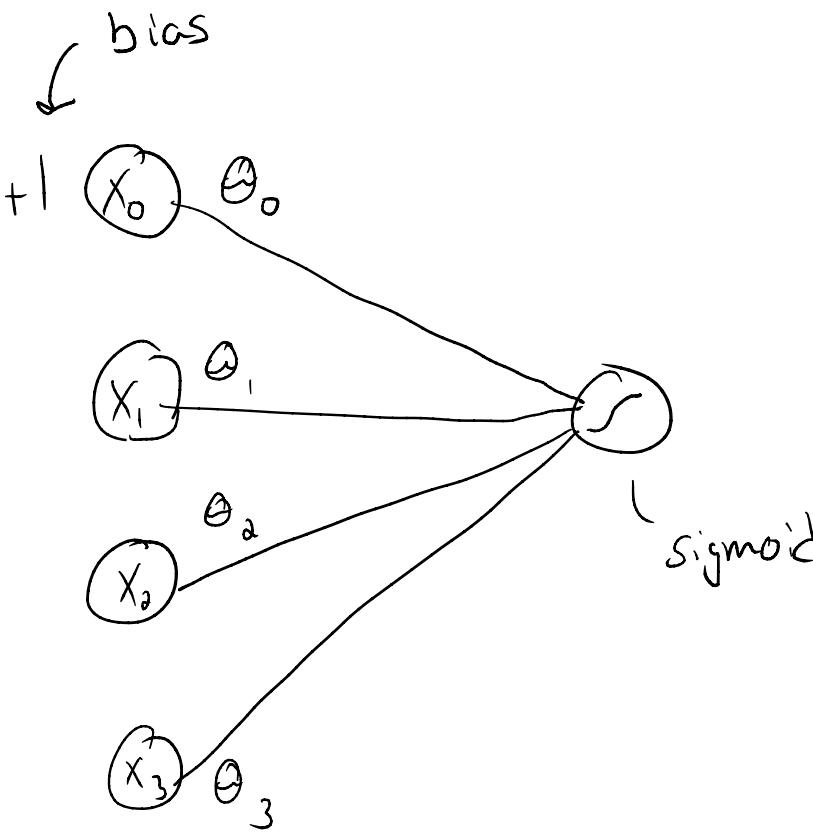


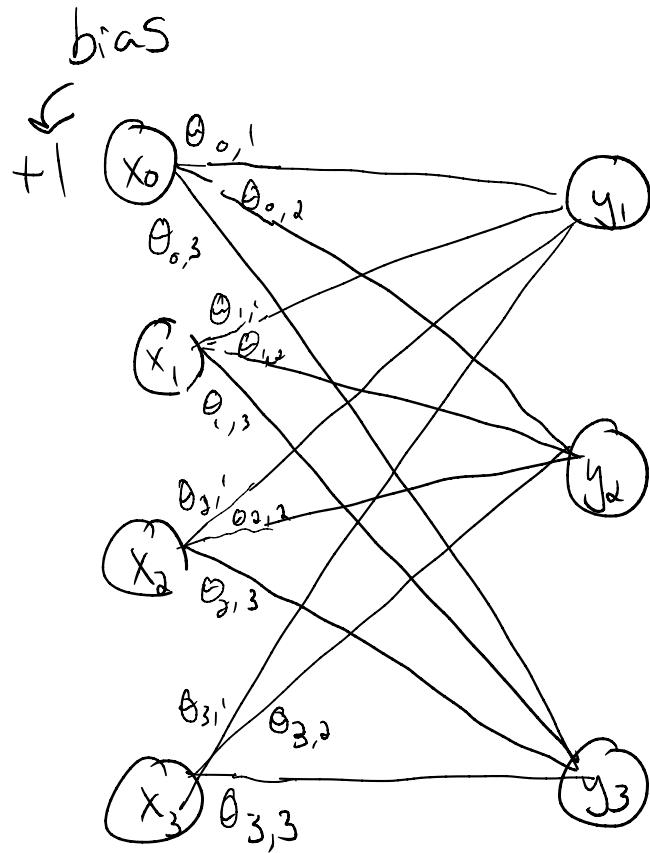
Recall from last class:

single layer neural network



$$y_i = \frac{1}{1 + \exp \left[ - \sum_{j=1}^d \theta_j x_j \right]}$$

for  $R > 2$ , add more outputs and use  
softmax instead of sigmoid



$$y_i = \frac{\exp(\theta_i^T \cdot x)}{\sum_k \exp(\theta_k^T \cdot x)}$$

↳ gives probabilities  
for each

- Alternative, each is option as well (this will be what you do in next project)

## Training Perceptron:

- update weights until convergence with some form of gradient descent
- update = learning-factor  $(\text{desiredout} - \text{actualout})$  input

$$\text{SGD} = \Delta \theta_{ij}^t = \alpha (r_j^t - y_j^t) x_i^t \quad \begin{matrix} \leftarrow \text{current sample} \\ \text{consider} \end{matrix}$$

$$\text{GD} = \Delta \theta_{ij} = \alpha \sum_t (r_j^t - y_j^t) x_i^t$$

# Power / Problems

- could implement AND, OR, NOT
- could not implement XOR
- could only represent linear functions
- solution → add intermediate ("hidden") layers between input and output layers

First thought:

$$\text{let } z_h = \sum_j \theta_{jh} x_j \quad h=1, \dots, H$$

$$\text{Then, let } y_i = \sum_{h=0}^H \psi_{h,i} z_h$$

Then apply softmax ( $h > 2$ ) or sigmoid ( $h = 2$ )

- doesn't actually gain us anything

↳ linear combination of linear combination  
is still just a linear combination

Need output of hidden layer to  
be nonlinear

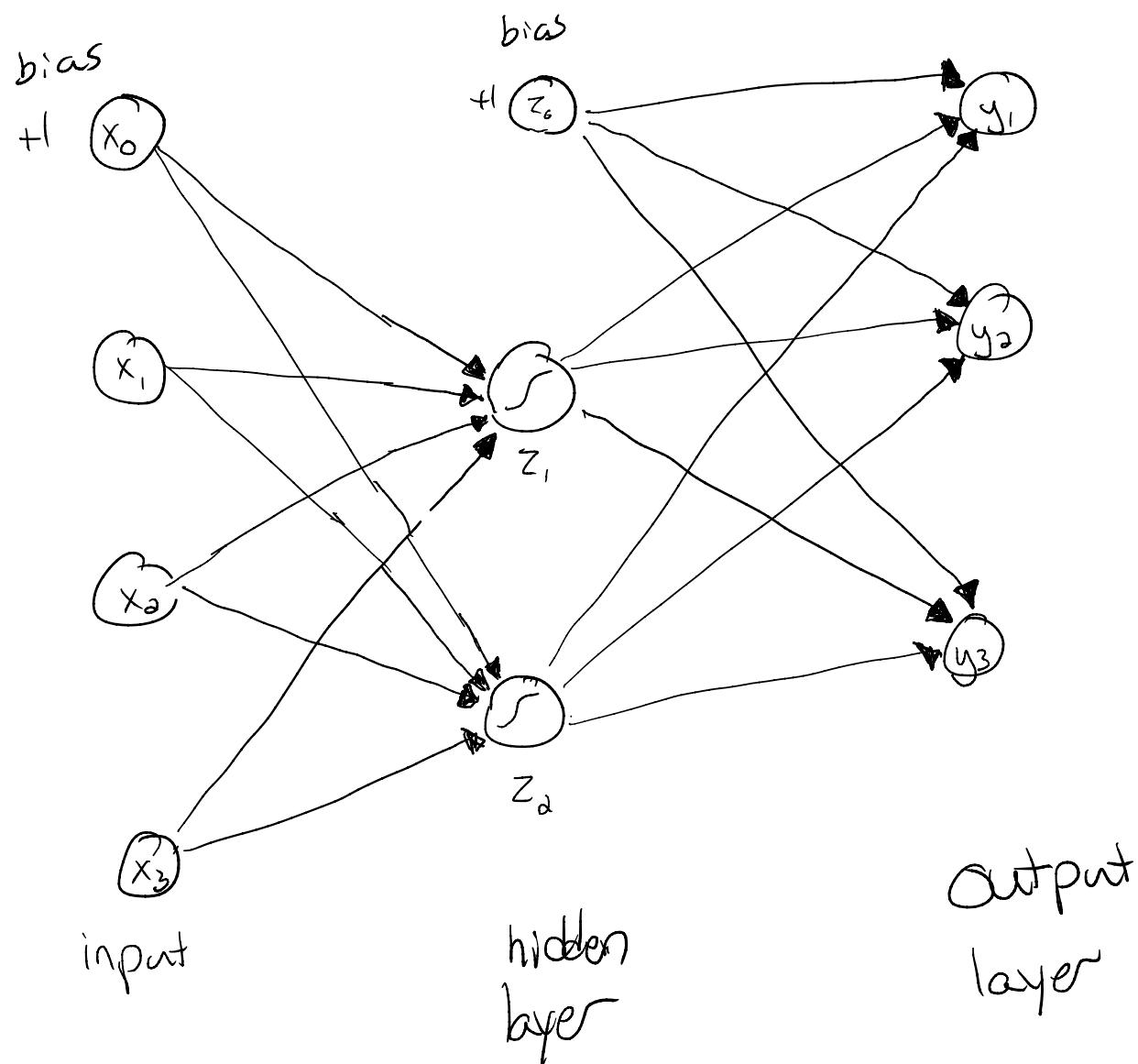
↳ get linear combination of  
nonlinear basis functions  
↳ nonlinear

Still want differentiability so we can use  
gradient based update

Options: sigmoid, tanh, etc.

Hidden units  $\rightarrow$  nonlinear transformation  
from  $d$ -dimensional input space to  
 $H$ -dimensional hidden layer space

Example : 2-layer



- Can add multiple hidden layers
- Training a multilayer neural network
  - error propagates from output  $y$  back to inputs  
 $\Rightarrow$  back propagation
  - gradient descent  $\rightarrow$  negative gradient points downhill
  - partial derivatives w/ respect to weight

$$\frac{\partial E}{\partial \theta_{jh}} = \frac{\partial E}{\partial y_i} \underbrace{\frac{\partial y_i}{\partial z_h}}_{\text{chain rule}} \underbrace{\frac{\partial z_h}{\partial \theta_{jh}}}$$

Basically two steps for  
each training instance

1) feed forward

plug in inputs and compute  
outputs

2) update weights based on how  
far off outputs were

via GD, SGD, MBGD

↳ multilayer must use back propagation

↳ need to update weights in  
earlier layers, not just  
weights for last layer

Multilayer: Squared error (Quadratic Loss)  
with sigmoid on each output

$$E(\theta | x^t, r^t) = \frac{1}{2} (r^t - y^t)^2$$

derivative of

$$y = \sigma(g) = \frac{1}{1 + \exp(-g)} \quad g = \theta^\top x$$

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

For each output  $i$ :

$$E_i = (r_i - y_i)y_i(1 - y_i)$$

because of chain  
rule

$$\Delta \gamma_{hi} = \alpha E_i z_h$$

For each hidden  $h$ :

$$E_h = z_h(1 - z_h) \sum_{i=1}^n E_i$$

$$\Delta \gamma_{ih} = \alpha E_h x_j$$

# Multilayer : Two Classes - Cross Entropy Cost

- still use one output with sigmoid activation function

$$y^t = \text{sigmoid} \left( \sum_{h=0}^{+} \psi_h z_h^t \right)$$

$$E(\theta, \psi | X) = - \sum_t r^t \log(y^t) + (1-r^t) \log(1-y^t)$$

⇒

$$\begin{cases} \Delta \psi_h = \alpha \sum_t (r^t - y^t) z_h^t \\ \Delta \theta_{jh} = \alpha \sum_t (r^t - y^t) \psi_h z_h^t (1 - z_h^t) x_j^t \end{cases}$$

weight update  
equations

Multilayer: > 2 classes cross entropy loss

$$y_i^t = \frac{\exp\left(\sum_{h=0}^H \gamma_{hi} z_h^t\right)}{\sum_h \exp\left(\sum_h \gamma_{hi} z_h\right)}$$

$$E(\theta, \gamma | x) = - \sum_t \sum_i r_i^t \log(y_i^t)$$

update equations

$$\begin{cases} \Delta \gamma_{hi} = \alpha \sum_t (r_i^t - y_i^t) z_h^t \\ \Delta \theta_{jh} = \alpha \sum_t \left[ \sum_i (r_i^t - y_i^t) \bar{\gamma}_{hi} \right] z_h^t (1 - z_h^t) x_j^t \end{cases}$$