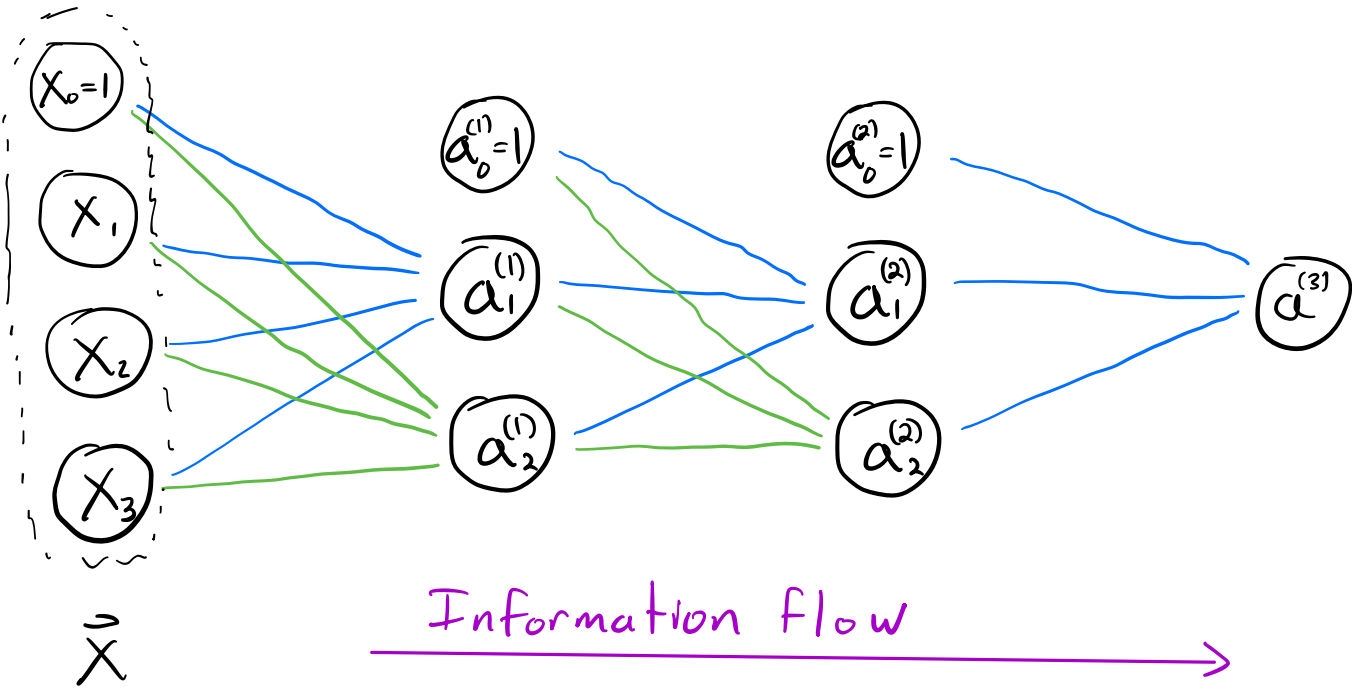


# Neural Networks (NN)

A NN is a function  $f(\vec{x})$

Ex:  $f(\vec{x}) = a^{(3)}$



## ERM with NNs (High Level)

$$\mathcal{D} = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n))$$

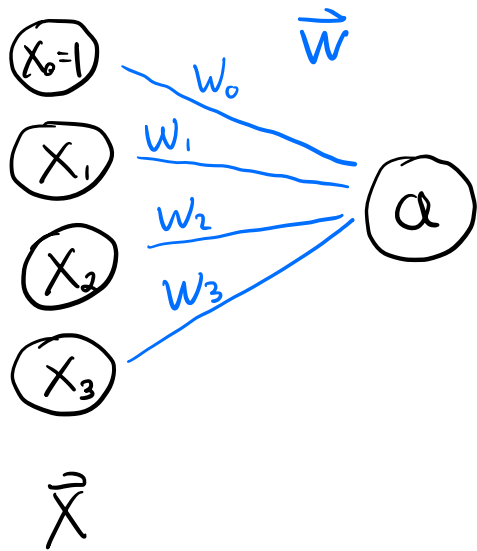
$$\mathcal{A}(\mathcal{D}) = \operatorname{argmin}_{f \in \mathcal{F}} \hat{L}(f)$$

$$\mathcal{F} = \{f \mid f: \mathcal{X} \rightarrow \mathcal{Y} \text{ and } f \text{ is a NN}\}$$

every  $f \in \mathcal{F}$  is defined by  $B$   
weight matrices  $W^{(1)}, \dots, W^{(B)}$

# Previously Seen functions as NNs

Linear:  $f(\vec{x}) = \vec{x}^T \vec{w}$ ,  $d=3$

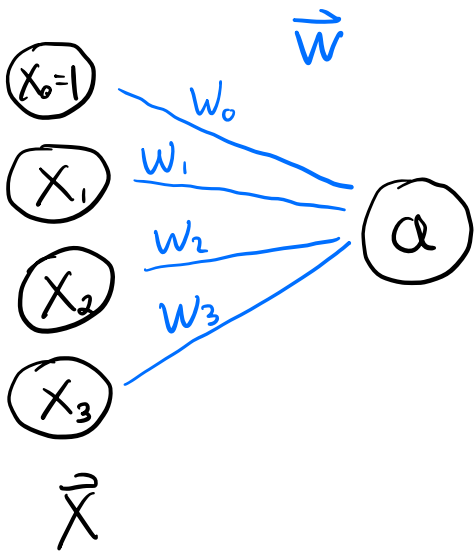


$$a = x_0 w_0 + x_1 w_1 + \dots + x_3 w_3$$

$$= \vec{x}^T \vec{w}$$

$$= h(\vec{x}^T \vec{w}) \quad \text{where } h(z) = z$$

Sigmoid:  $f(\vec{x}) = \sigma(\vec{x}^T \vec{w})$



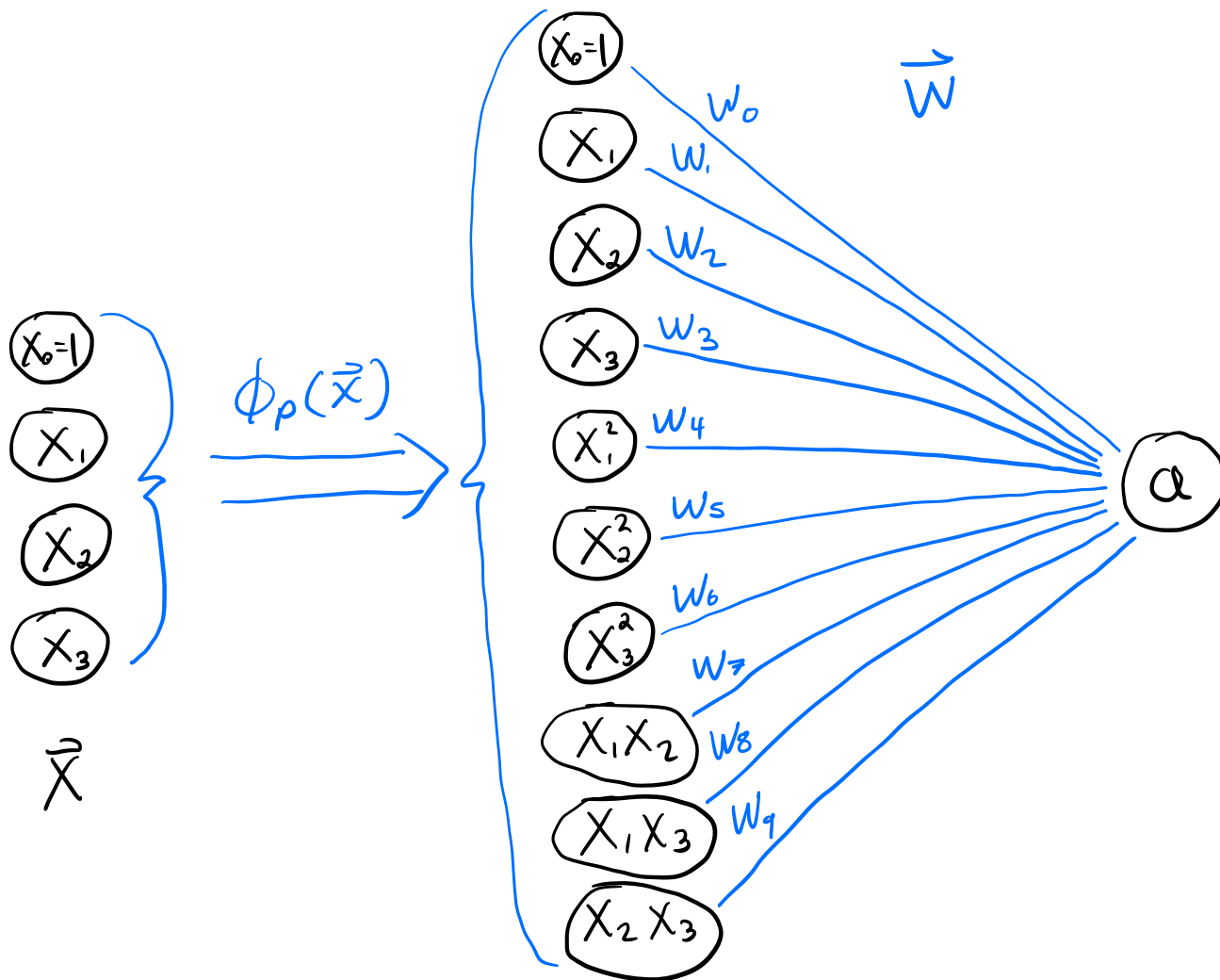
$$a = \sigma(x_0 w_0 + x_1 w_1 + \dots + x_3 w_3)$$

$$= \sigma(\vec{x}^T \vec{w})$$

$$= h(\vec{x}^T \vec{w}) \quad \text{where } h = \sigma$$

Linear (with polynomial features):

$$f(\vec{x}) = \Phi_p(\vec{x})^T \vec{w} = a, \quad d=3, p=2$$

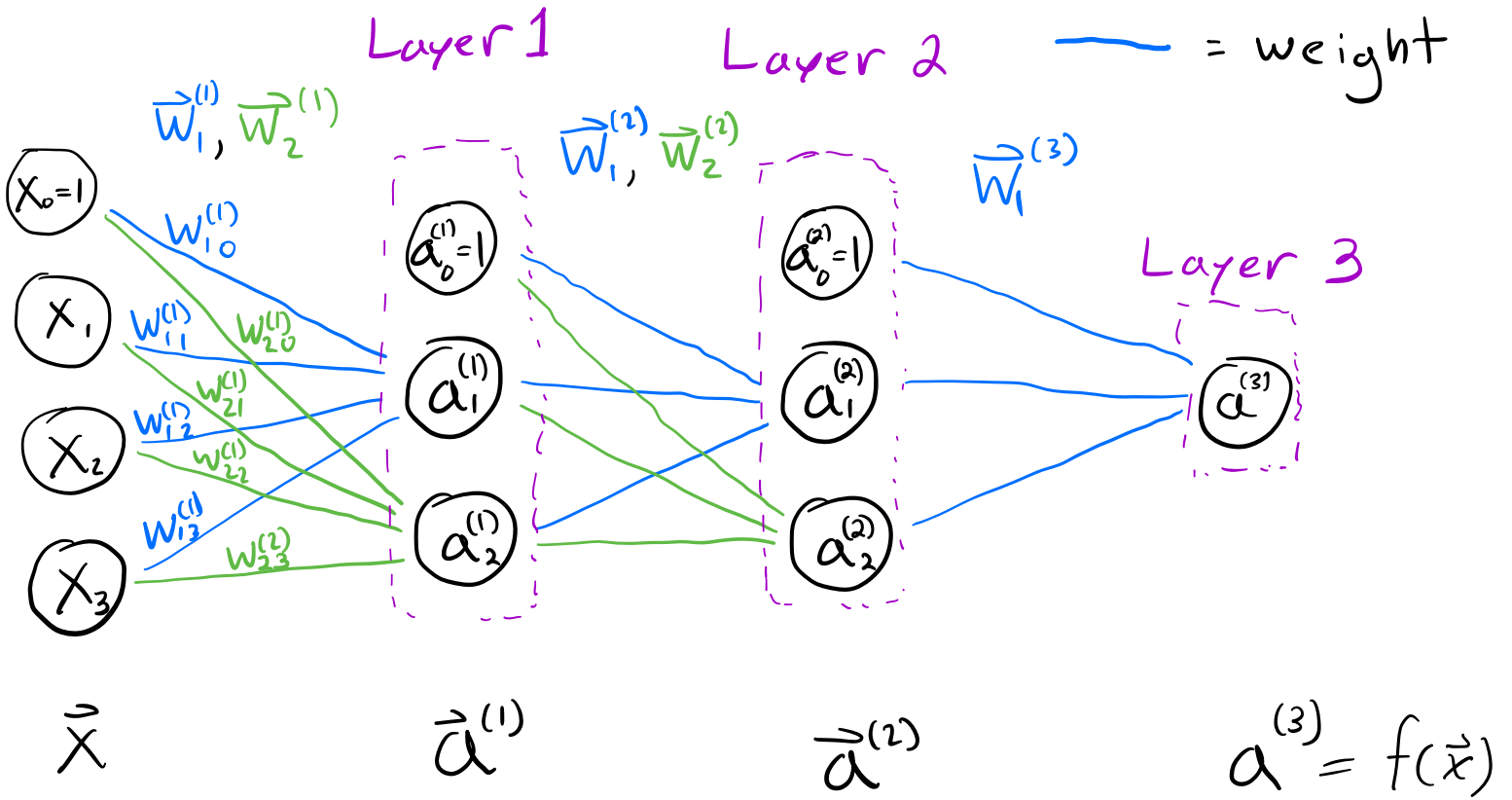


Many of the new features in  $\Phi_p(\vec{x})$  may not be useful

When  $p$  is large this is very wasteful

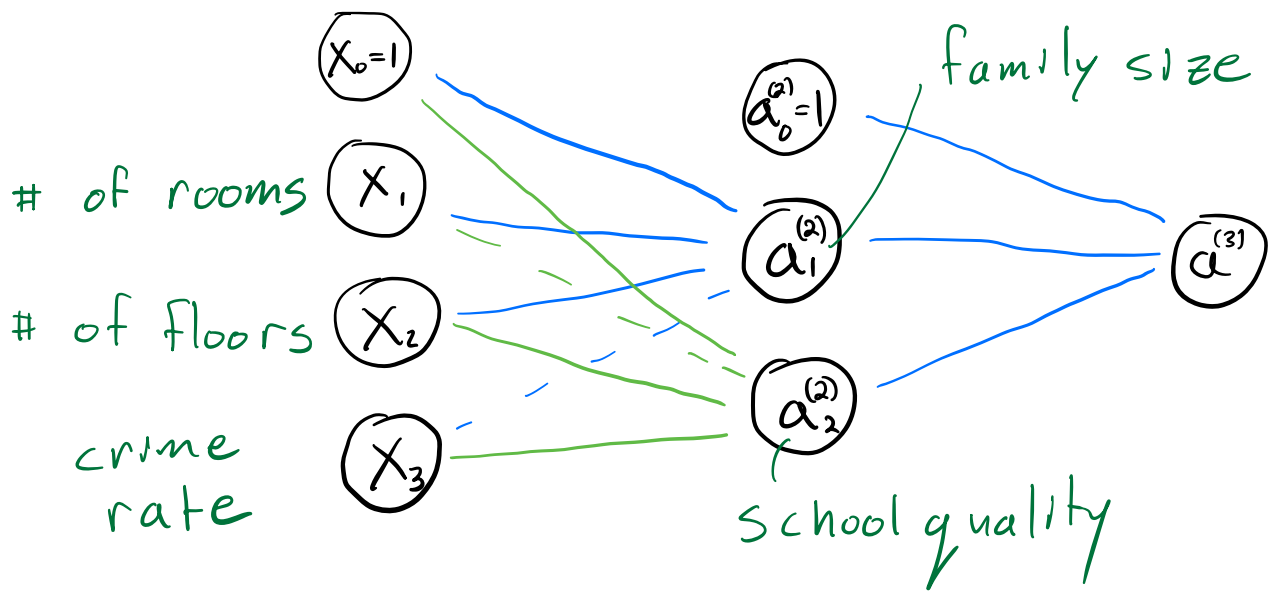
# NN Terminology

○ = Neuron  
— = weight



Ex: Regression (House price prediction)

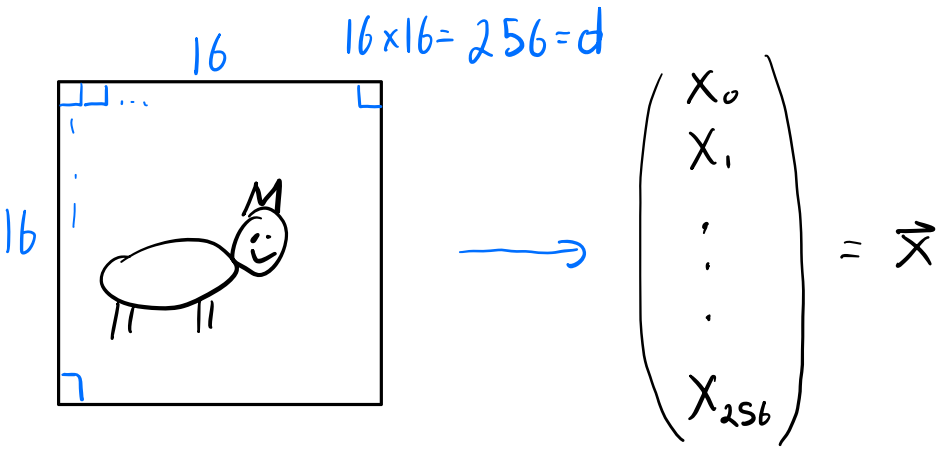
$$y = \mathbb{R}, \mathcal{X} = \mathbb{R}^{d+1}, d=3$$



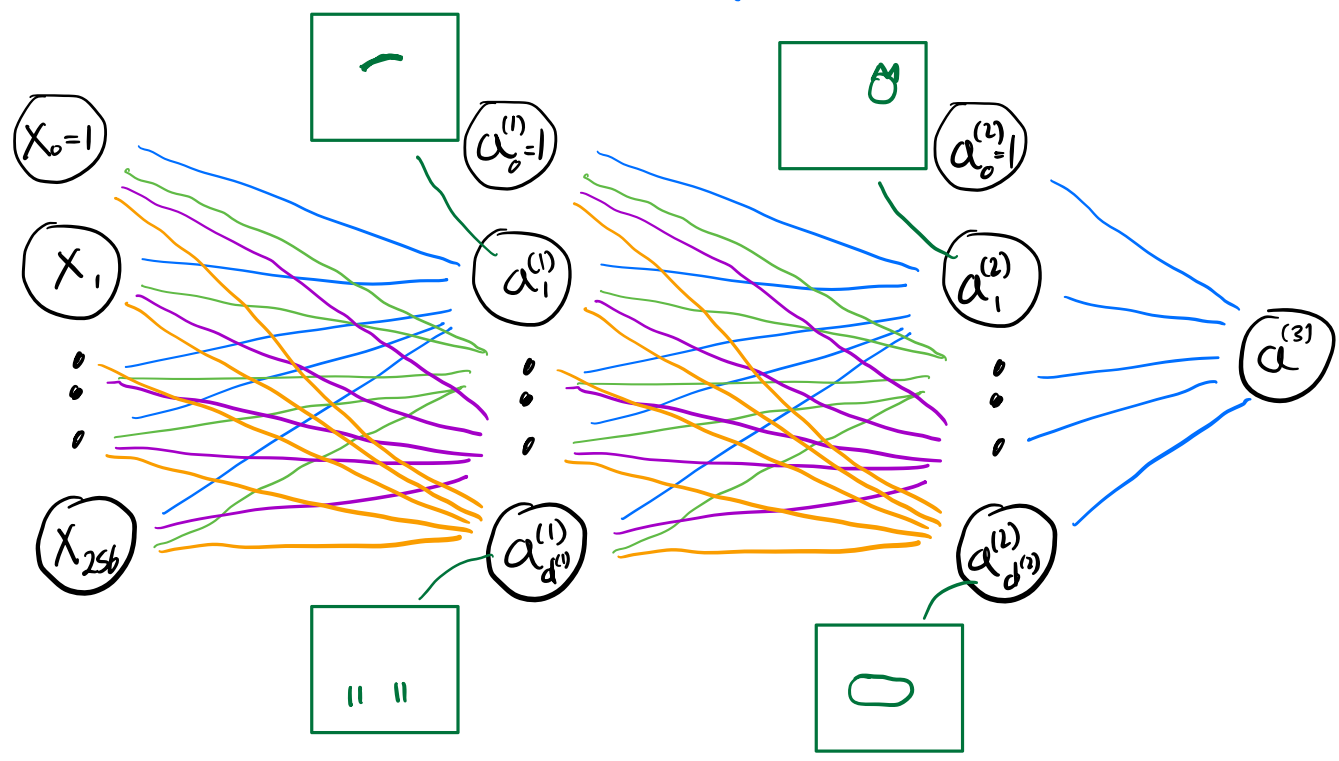
# Ex: Binary Classification (Cat or No Cat)

$$Y = \{0, 1\} \quad X = \mathbb{R}^{256+1}$$

No cat  $\uparrow$   $\uparrow$  Cat



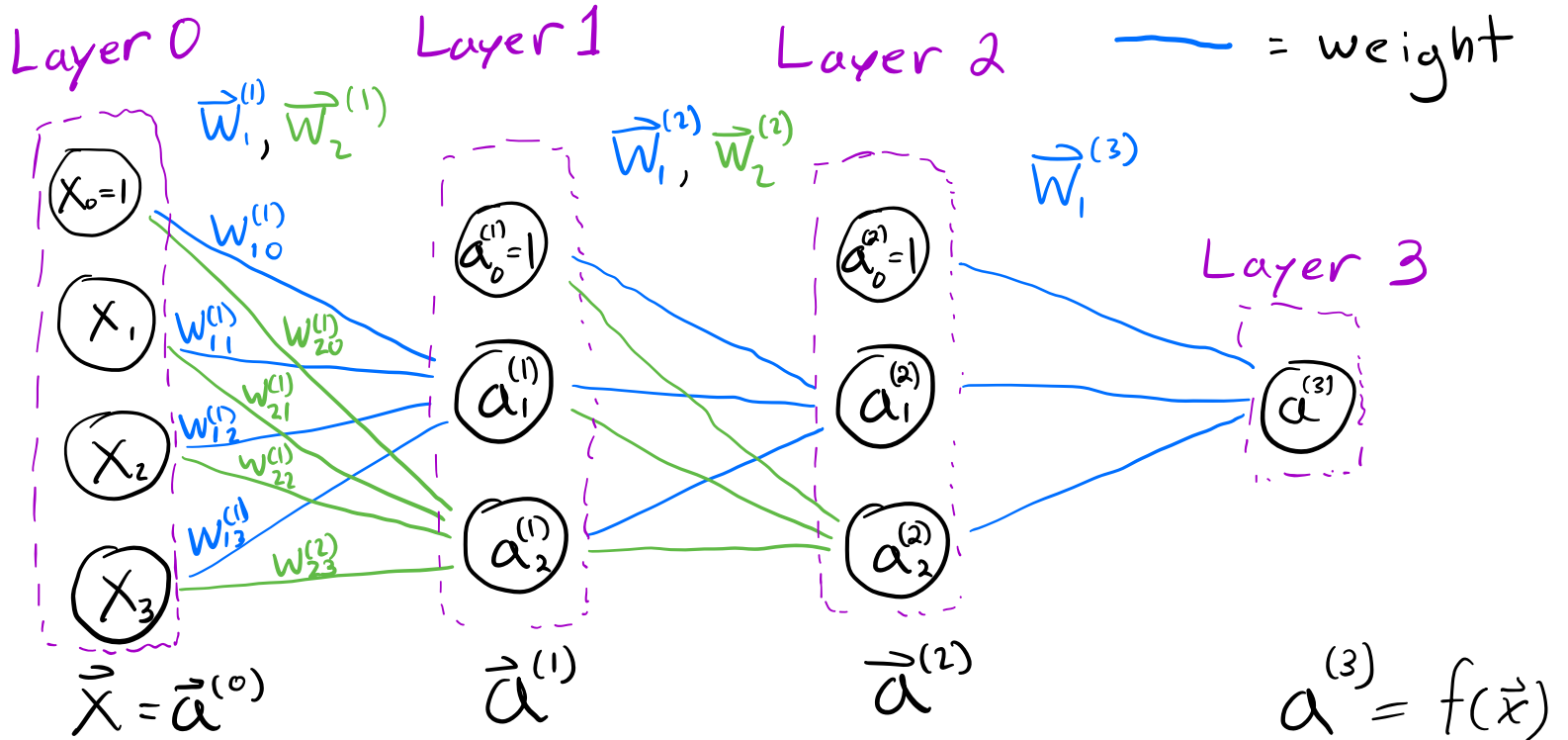
$$f(\vec{x}) = a^{(3)} \quad \text{probability of a cat}$$



# NN Definition (formal)

$d=3$

$\bigcirc$  = Neuron  
 $\text{---}$  = weight



activation

$$\vec{a}^{(1)} = (a_0^{(1)}=1, a_1^{(1)}, a_2^{(1)})$$

where  $a_1^{(1)} = h^{(1)}(z_1^{(1)})$ ,  $a_2^{(1)} = h^{(1)}(z_2^{(1)})$

activation function  $h^{(l)}: \mathbb{R} \rightarrow \mathbb{R}$  ex: sigmoid  $\sigma$

and  $z_1^{(1)} = (\vec{a}^{(0)})^T \vec{w}_1^{(1)}$ ,  $z_2^{(1)} = (\vec{a}^{(0)})^T \vec{w}_2^{(1)}$

pre-activation

$$\vec{a}^{(0)} \in \mathbb{R}^{d+1}, \vec{w}_1^{(1)}, \vec{w}_2^{(1)} \in \mathbb{R}^{d+1}, \vec{a}^{(1)} \in \mathbb{R}^{2+1}$$

$$\vec{a}^{(2)} = (a_0^{(2)}=1, a_1^{(2)}, a_2^{(2)})$$

$$\text{where } a_1^{(2)} = h^{(2)}(z_1^{(2)}), \quad a_2^{(2)} = h^{(2)}(z_2^{(2)})$$

$$\text{and } z_1^{(2)} = (\vec{a}^{(1)})^T \vec{w}_1^{(2)}, \quad z_2^{(2)} = (\vec{a}^{(1)})^T \vec{w}_2^{(2)}$$

$$\vec{w}_1^{(2)}, \vec{w}_2^{(2)} \in \mathbb{R}^{2+1}, \quad \vec{a}^{(2)} \in \mathbb{R}^{2+1}$$

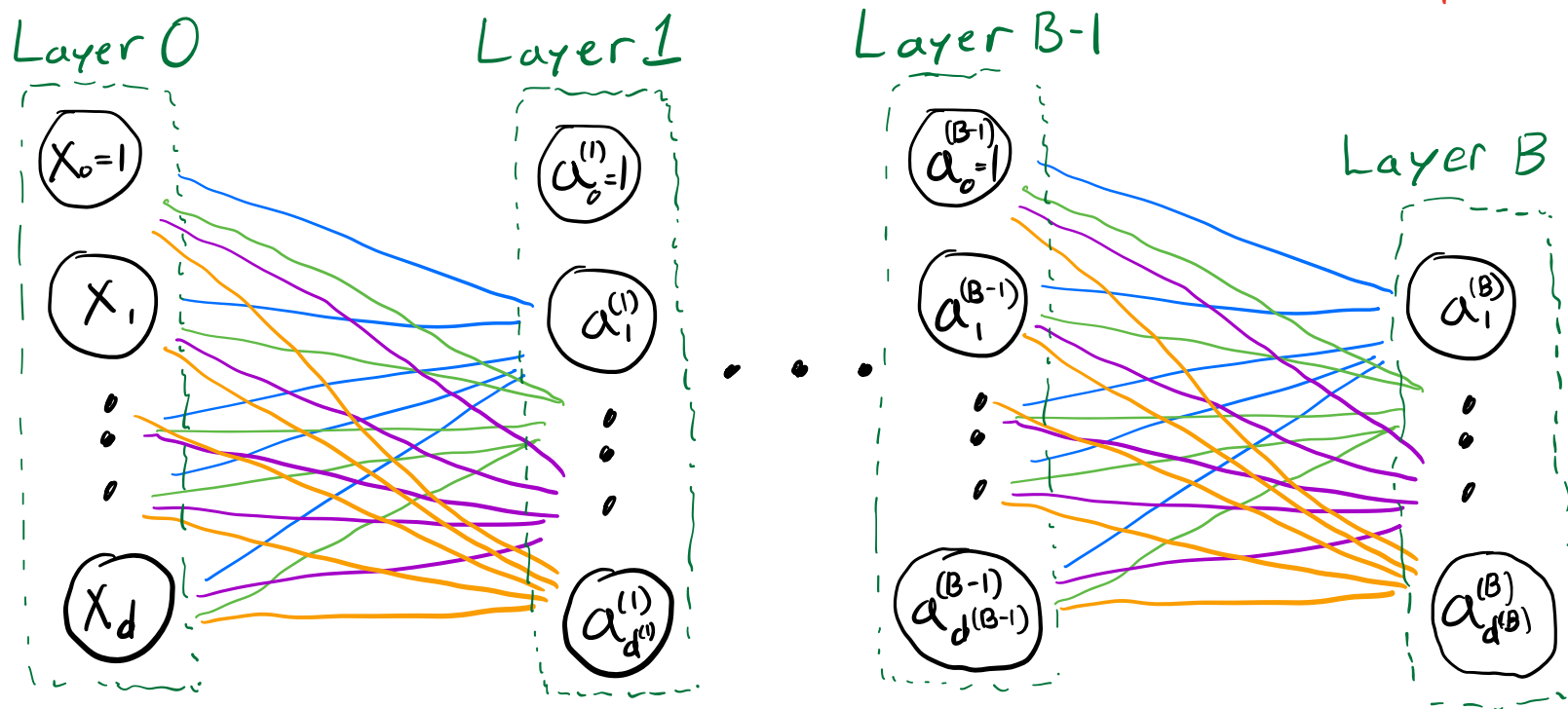
$$a^{(3)} = h^{(3)}(z^{(3)}) \quad \text{where} \quad z^{(3)} = (\vec{a}^{(2)})^T \vec{w}^{(3)}$$

$$\vec{w}_j^{(3)} \in \mathbb{R}^{2+1}, \quad a^{(3)} \in \mathbb{R}$$



In general:

No Bias  
in Layer B



For layer  $b \in \{1, \dots, B\}$

Weights

$$\vec{w}_1^{(b)}, \dots, \vec{w}_{d^{(b)}}^{(b)} \in \mathbb{R}^{d^{(b-1)}+1}$$

Pre-activations

$$\vec{z}^{(b)} = (z_1^{(b)}, \dots, z_{d^{(b)}}^{(b)}) \in \mathbb{R}^{d^{(b)}}$$

$$\text{where } z_j^{(b)} = (\vec{a}^{(b-1)})^T \vec{w}_j \quad \text{for } j \in \{1, \dots, d^{(b)}\}$$

Activations

$$\vec{a}^{(0)} = \vec{x} \in \mathbb{R}^{d+1} = \mathbb{R}^{d^{(0)}+1}, \quad d = d^{(0)}$$

$$\vec{a}^{(b)} = (a_0^{(b)}=1, a_1^{(b)}, \dots, a_{d^{(b)}}^{(b)}) \in \mathbb{R}^{d^{(b)}+1}$$

except  $b=B$

$$\vec{a}^{(B)} = (a_1^{(B)}, \dots, a_{d^{(B)}}^{(B)}) \in \mathbb{R}^{d^{(B)}}$$

where  $a_j^{(b)} = h(z_j^{(b)})$  for  $j \in \{1, \dots, d^{(b)}\}$

## Activation function

$$h: \mathbb{R} \rightarrow \mathbb{R}$$

$$\underline{NN}: f(\vec{x}) = \vec{a}^{(B)}$$

Defining the number of layers  $B$

and the number of neurons in each layer:

$d^{(1)}, \dots, d^{(B)}$ , and the activation functions

$h^{(1)}, \dots, h^{(B)}$  defines a "NN architecture"

# Matrix Notation

Represent  $w_1^{(b)}, \dots, w_d^{(b)}$  as a matrix:

$$W^{(b)} = \begin{bmatrix} | & | & & | \\ w_1^{(b)} & w_2^{(b)} & \dots & w_d^{(b)} \\ | & | & & | \end{bmatrix}$$

dimension:  
 $(d^{(b-1)} + 1)$  by  $d^{(b)}$

← applied element-wise

$$\vec{a}^{(b)} = h^{(b)}\left(\left(\vec{a}^{(b-1)}\right)^T W^{(b)}\right)$$

$$\begin{aligned} h^{(b)}(\vec{v}) &= h^{(b)}(v_1, \dots, v_d) \\ &= (h^{(b)}(v_1), \dots, h^{(b)}(v_d)) \end{aligned}$$

$$f(\vec{x}) = \vec{a}^{(B)} = h^{(B)}\left(h^{(B-1)}\left(\dots h^{(2)}\left(h^{(1)}\left(\vec{x}\right)^T W^{(1)}\right) W^{(2)}\right) \dots W^{(B-1)}\right) W^{(B)}{}^T$$

# ERM with Neural Networks

$$\mathcal{D} = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n))$$

$$\mathcal{A}(\mathcal{D}) = \hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \hat{L}(f) \text{ where } \hat{L}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(\vec{x}_i), y_i)$$

$$\mathcal{F} = \{f \mid f: \mathcal{X} \rightarrow \mathcal{Y} \text{ and } f \text{ is a NN with a specific architecture}\}$$

every  $f_{W^{(1)}, \dots, W^{(B)}} \in \mathcal{F}$  is defined by B weight matrices  $W^{(1)}, \dots, W^{(B)}$

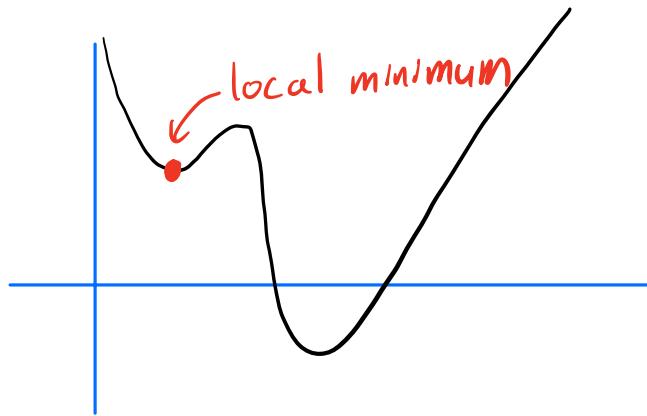
$$\mathcal{A}(\mathcal{D}) = \hat{f} \text{ where } \hat{f}(\vec{x}) = \vec{a}^{(B)} \text{ is defined by}$$

$$\hat{W}^{(1)}, \dots, \hat{W}^{(B)} = \operatorname{argmin}_{W^{(1)}, \dots, W^{(B)}} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(f_{W^{(1)}, \dots, W^{(B)}}(\vec{x}_i), y_i)}_{= \hat{L}(W^{(1)}, \dots, W^{(B)})}$$

Usually no closed form solution

$$\begin{pmatrix} \vec{W}_j^{(b)} \end{pmatrix}^{(t+1)} = \begin{pmatrix} \vec{W}_j^{(b)} \end{pmatrix}^{(t)} - \eta^{(t)} \nabla_{\vec{W}_j^{(b)}} \hat{L}((W^{(1)})^{(t)}, \dots, (W^{(B)})^{(t)})$$

$\hat{L}(W^{(1)}, \dots, W^{(B)})$  is usually not convex



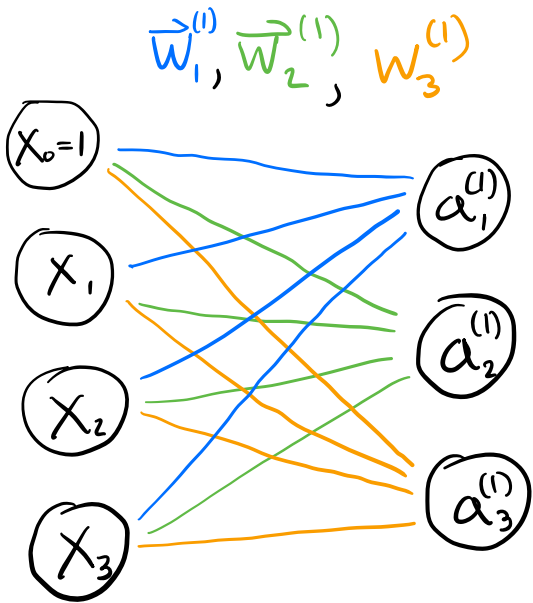
# Softmax

$$d=3, k=3,$$

$$\sigma(\vec{x}^T \vec{w}_1^{(1)}, \vec{x}^T \vec{w}_2^{(1)}, \vec{x}^T \vec{w}_3^{(1)})$$

$$= (\sigma_1(\vec{x}^T \vec{w}_1^{(1)}, \vec{x}^T \vec{w}_2^{(1)}, \vec{x}^T \vec{w}_3^{(1)}), \sigma_2(\vec{x}^T \vec{w}_1^{(1)}, \vec{x}^T \vec{w}_2^{(1)}, \vec{x}^T \vec{w}_3^{(1)}), \sigma_3(\vec{x}^T \vec{w}_1^{(1)}, \vec{x}^T \vec{w}_2^{(1)}, \vec{x}^T \vec{w}_3^{(1)}))^T$$

Softmax cannot be implemented  
by a single layer NN



$$a_1^{(1)} = h^{(1)}(\vec{x}^T \vec{w}_1^{(1)})$$

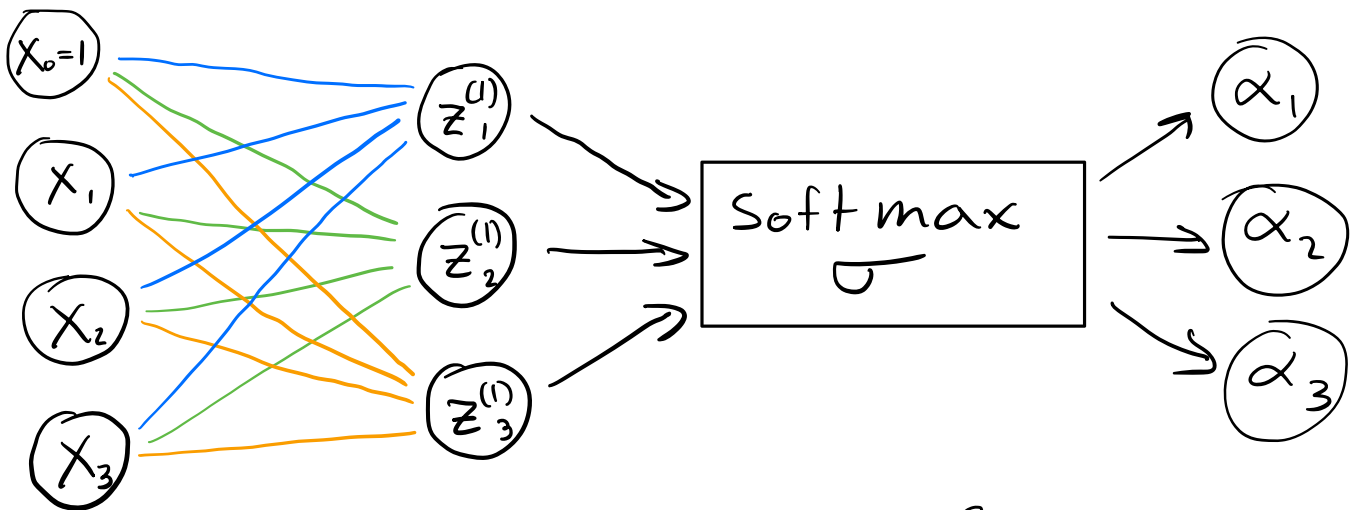
$$a_2^{(1)} = h^{(1)}(\vec{x}^T \vec{w}_2^{(1)})$$

$$a_3^{(1)} = h^{(1)}(\vec{x}^T \vec{w}_3^{(1)})$$

$$\text{if } h^{(1)}(z) = z$$

$$f(\vec{x}) = \vec{a}^{(1)} = (\vec{x}^T \vec{w}_1^{(1)}, \vec{x}^T \vec{w}_2^{(1)}, \vec{x}^T \vec{w}_3^{(1)})^T$$

$$\sigma(f(\vec{x})) \quad \text{post-processing step}$$



$$\sum_{q=1}^3 \alpha_q = 1, \alpha_q \in [0, 1]$$

$$\sigma(z_1^{(1)}, z_2^{(1)}, z_3^{(1)})$$

$$= \left( \frac{\exp(z_1^{(1)})}{\sum_{q=1}^3 \exp(z_q^{(1)})}, \frac{\exp(z_2^{(1)})}{\sum_{q=1}^3 \exp(z_q^{(1)})}, \frac{\exp(z_3^{(1)})}{\sum_{q=1}^3 \exp(z_q^{(1)})} \right)^T$$