

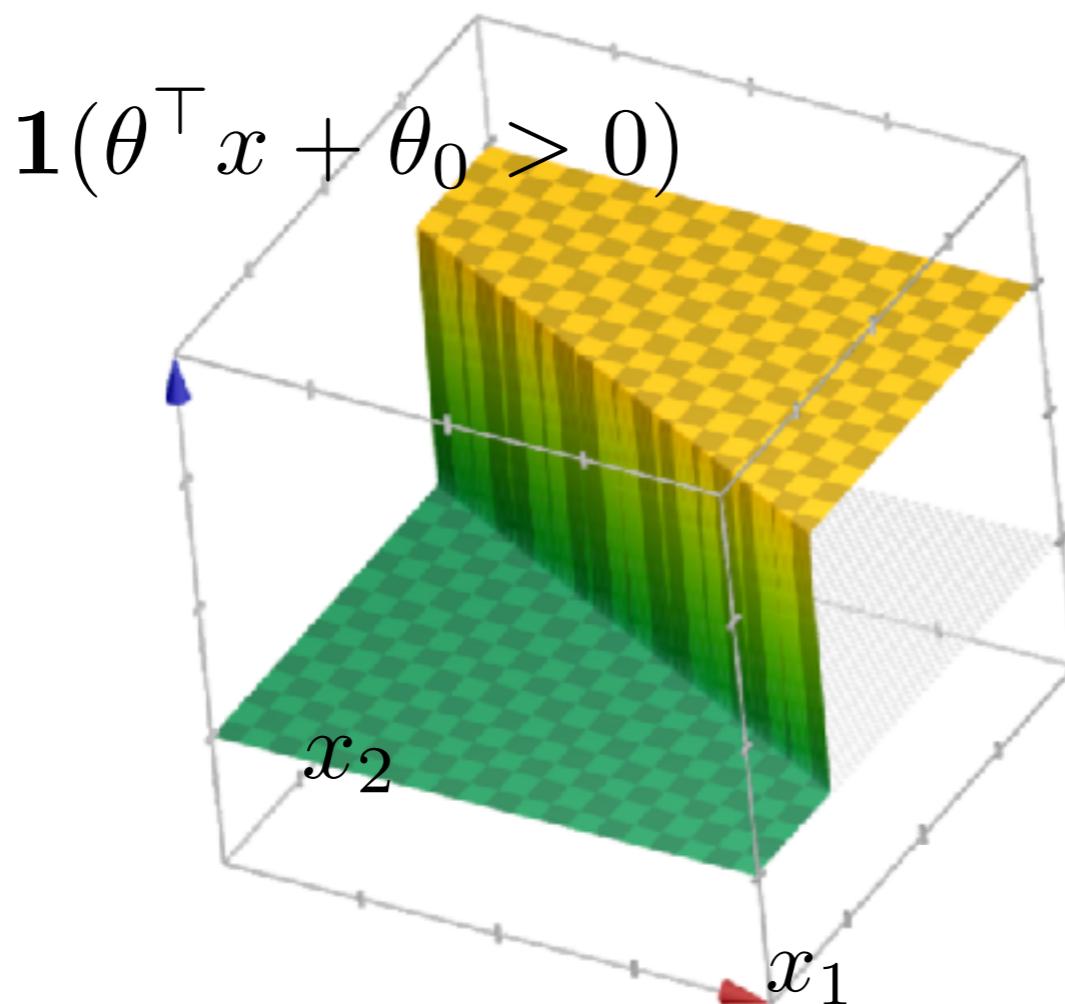
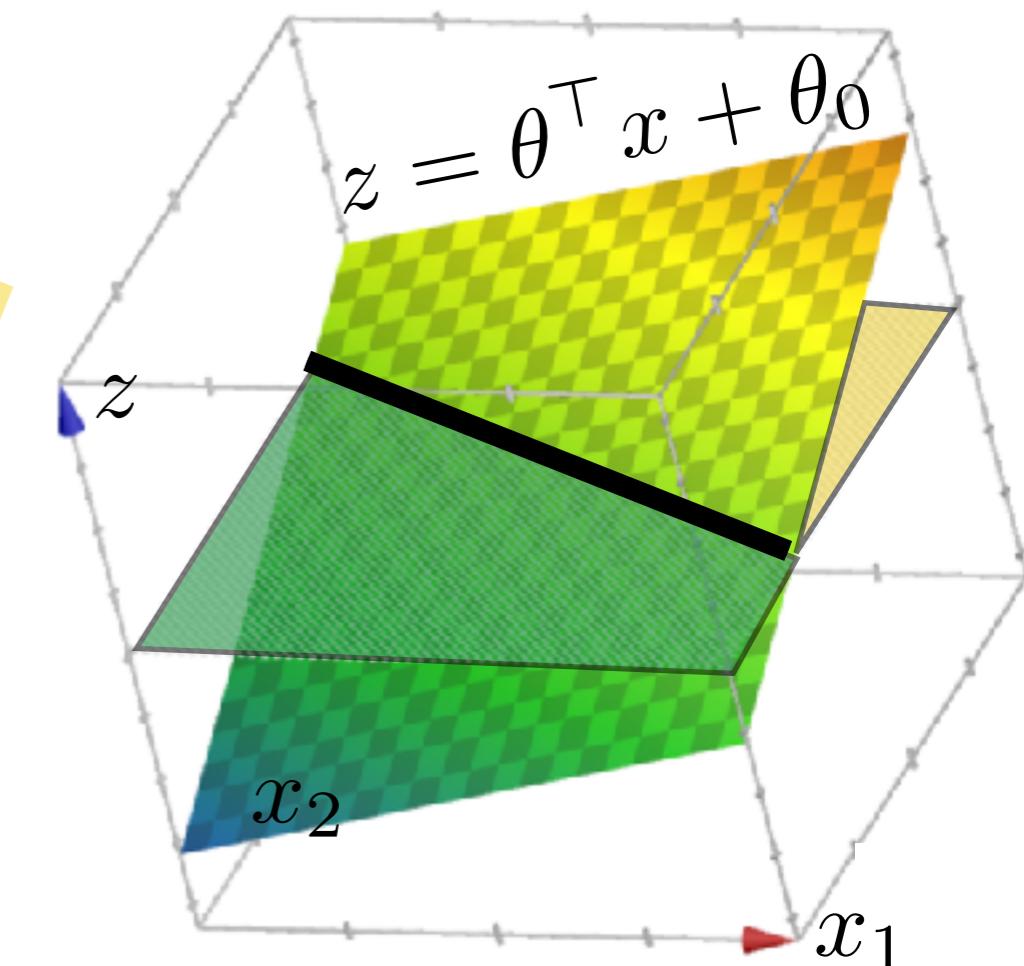
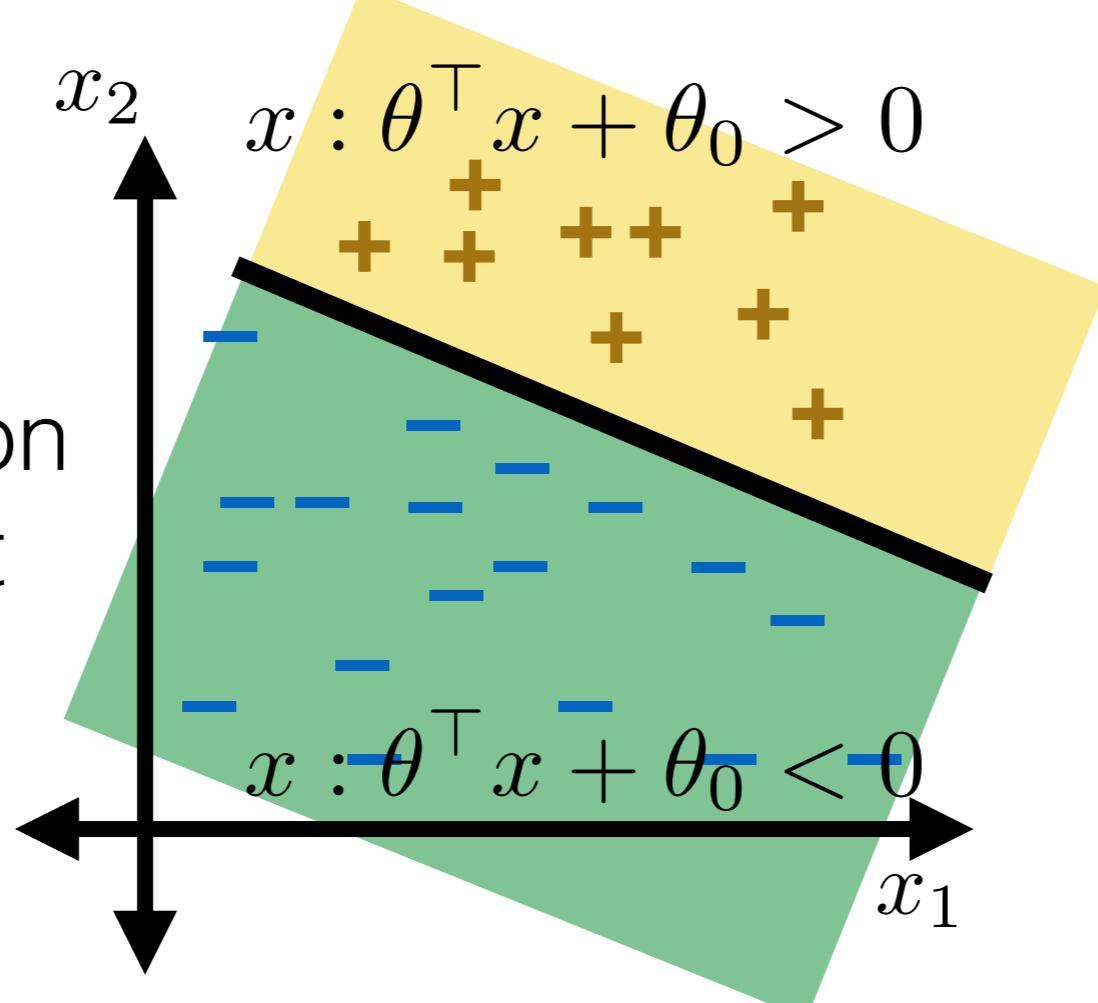
# **Neural Networks**

**Prof. Tamara Broderick**

**Edited From 6.036 Fall21 Offering**

# Recall

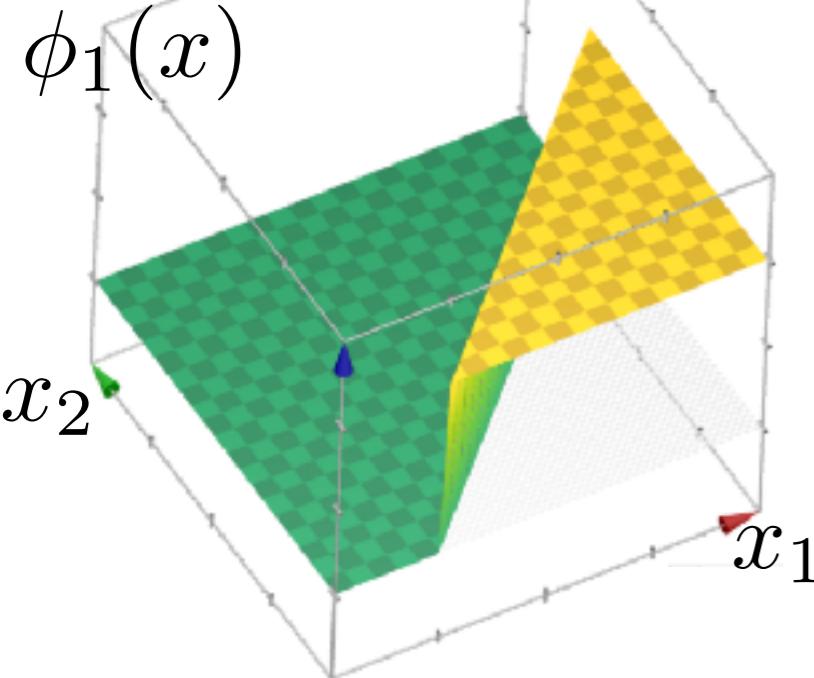
- Linear classification with default features:



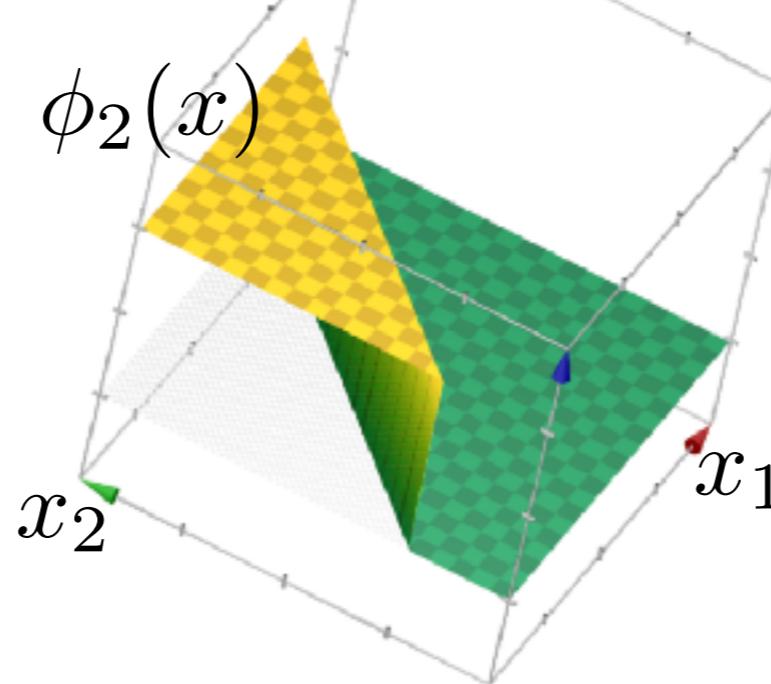
- We're used to using step functions to classify
- New idea today: we'll use step functions as *features*, with their own parameters

# New features: step functions!

$$\phi_1(x) = \mathbf{1}\{w^\top x + w_0 \geq 0\}$$

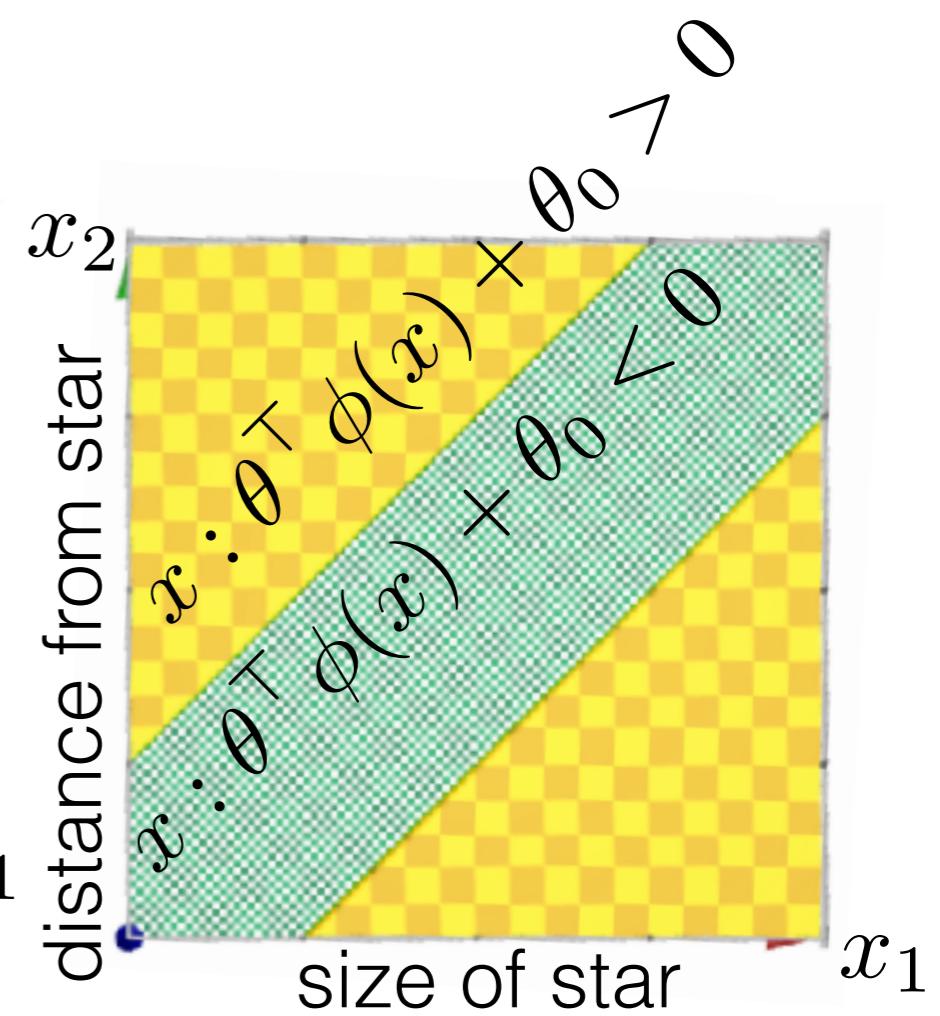
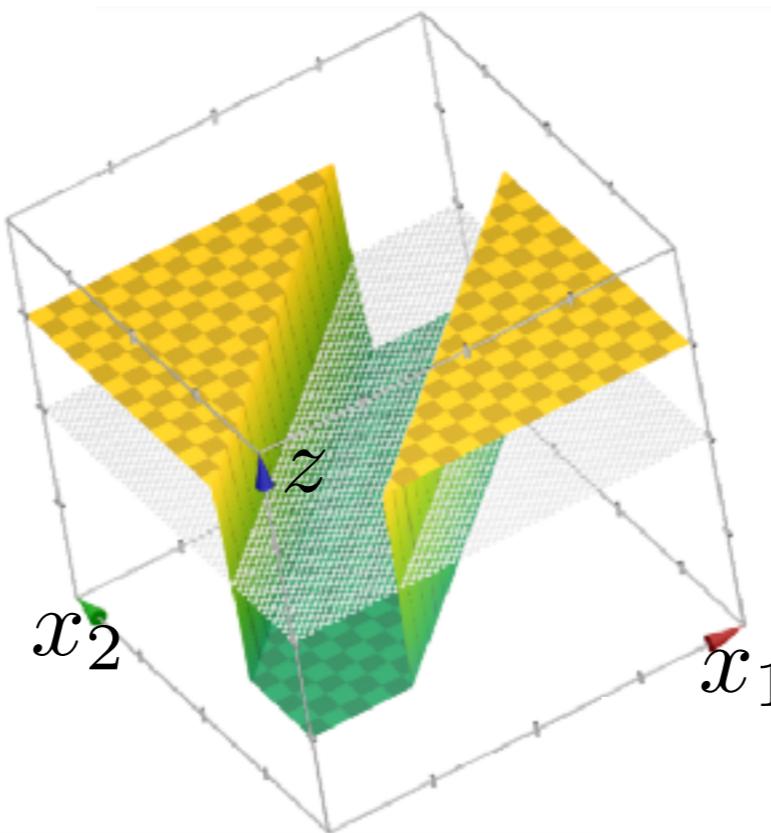


$$\phi_2(x) = \mathbf{1}\{\tilde{w}^\top x + \tilde{w}_0 \geq 0\}$$



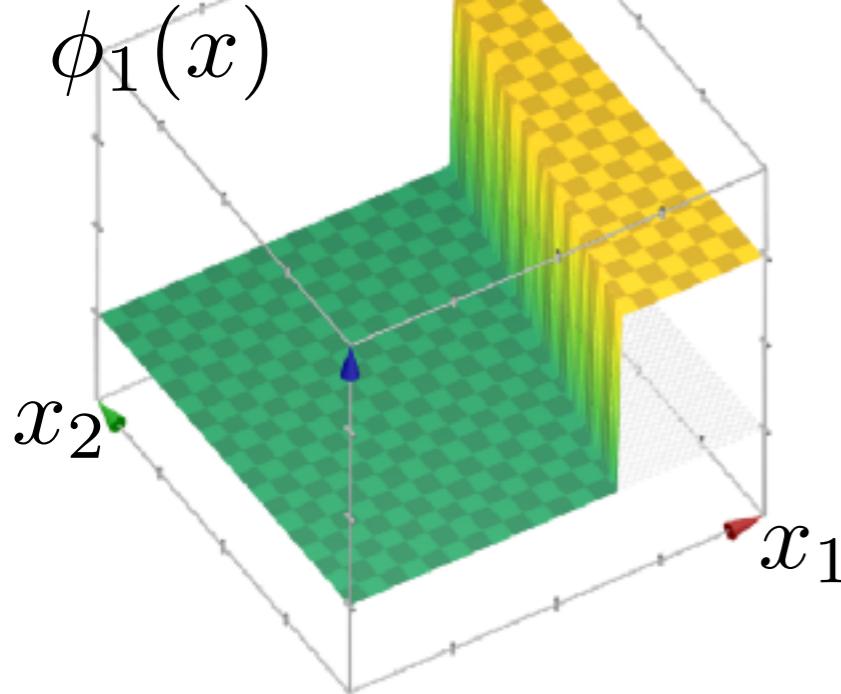
Is an exoplanet outside the habitable zone?

$$\begin{aligned} z &= \theta^\top \phi(x) + \theta_0 \\ &= \theta_1 \phi_1(x) + \theta_2 \phi_2(x) \\ &\quad + \theta_0 \\ &= 1 \cdot \phi_1(x) + 1 \cdot \phi_2(x) \\ &\quad + (-0.5) \end{aligned}$$

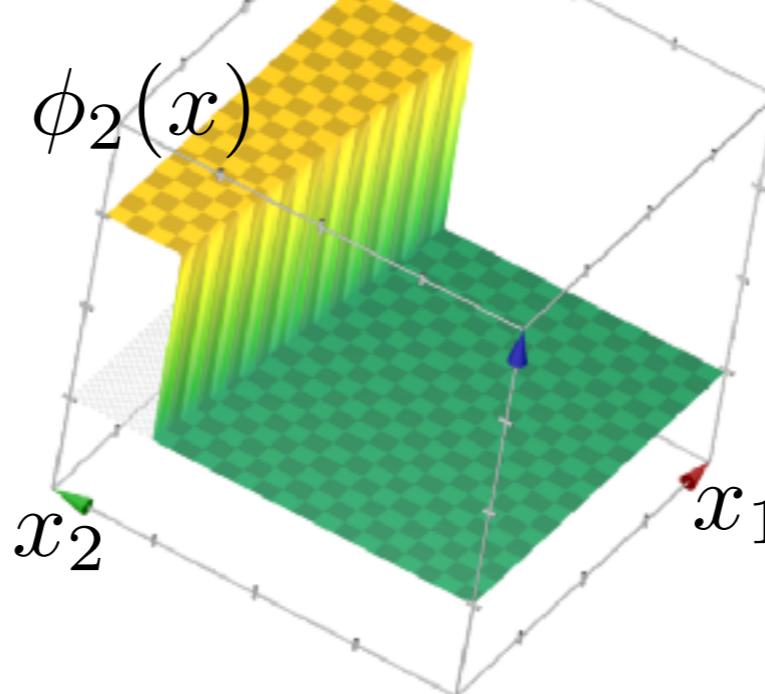


# New features: step functions!

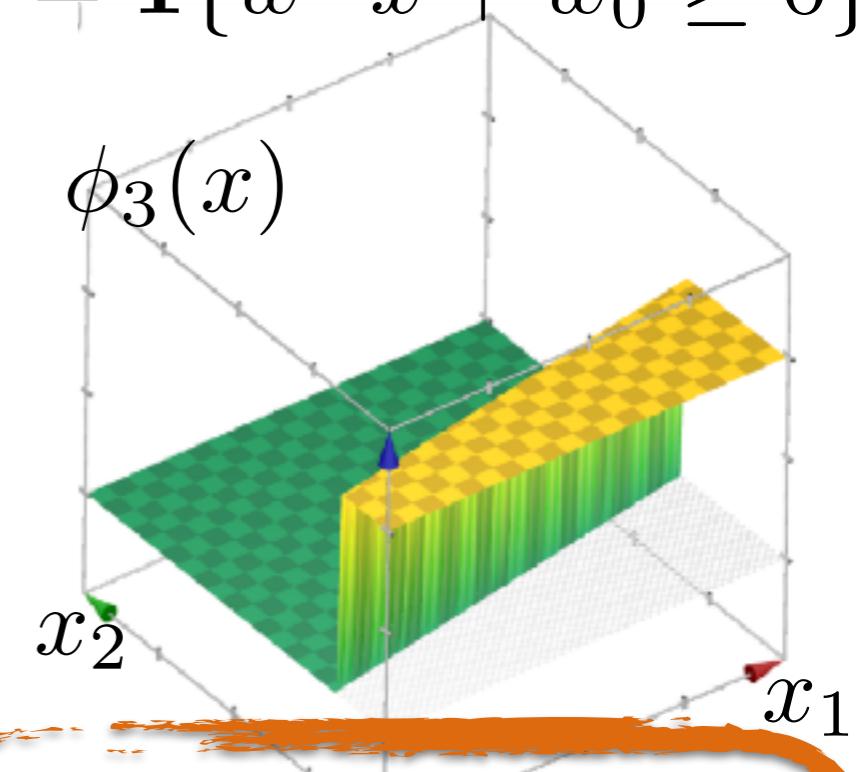
$$\phi_1(x) = \mathbf{1}\{w^\top x + w_0 \geq 0\}$$



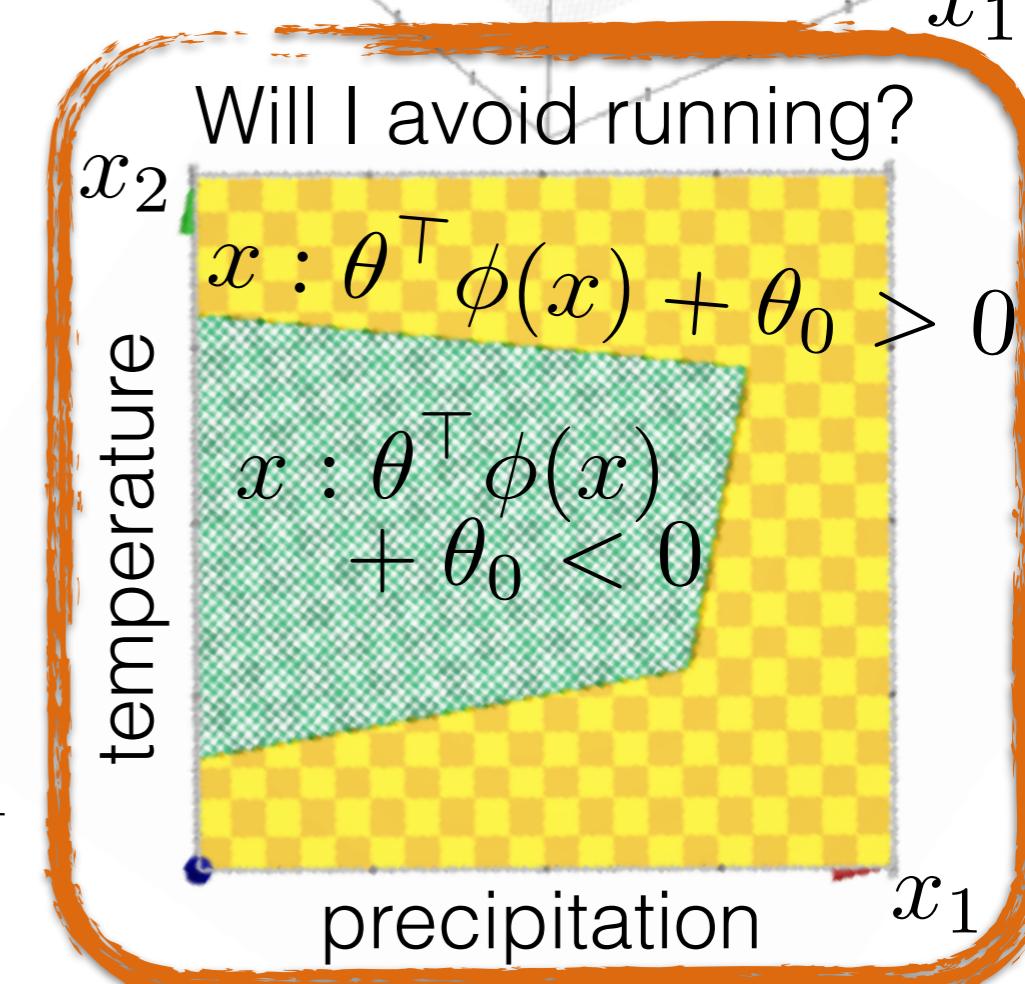
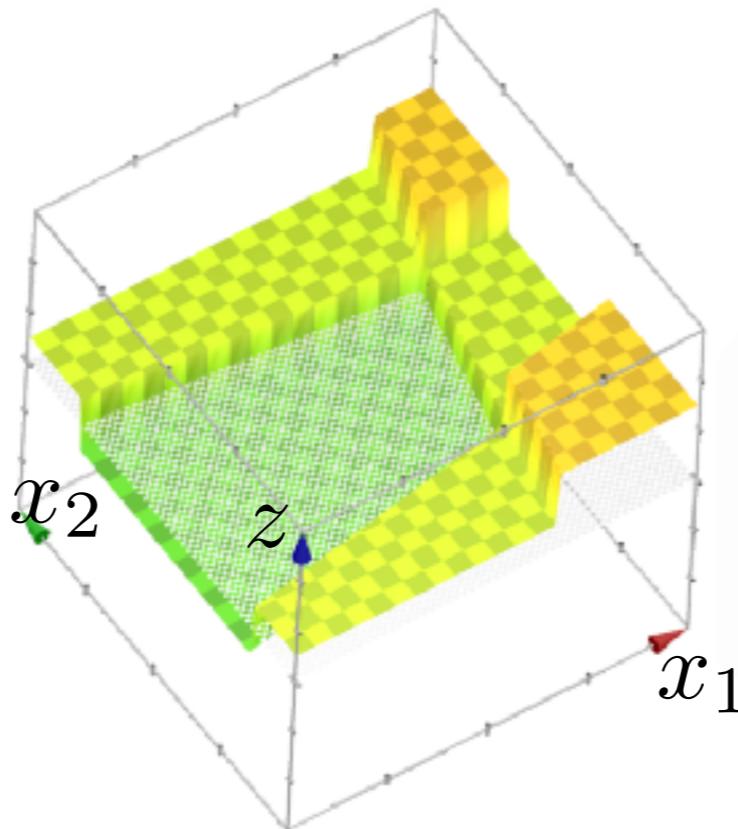
$$\phi_2(x) = \mathbf{1}\{\tilde{w}^\top x + \tilde{w}_0 \geq 0\}$$



$$\phi_3(x) = \mathbf{1}\{\tilde{\tilde{w}}^\top x + \tilde{\tilde{w}}_0 \geq 0\}$$



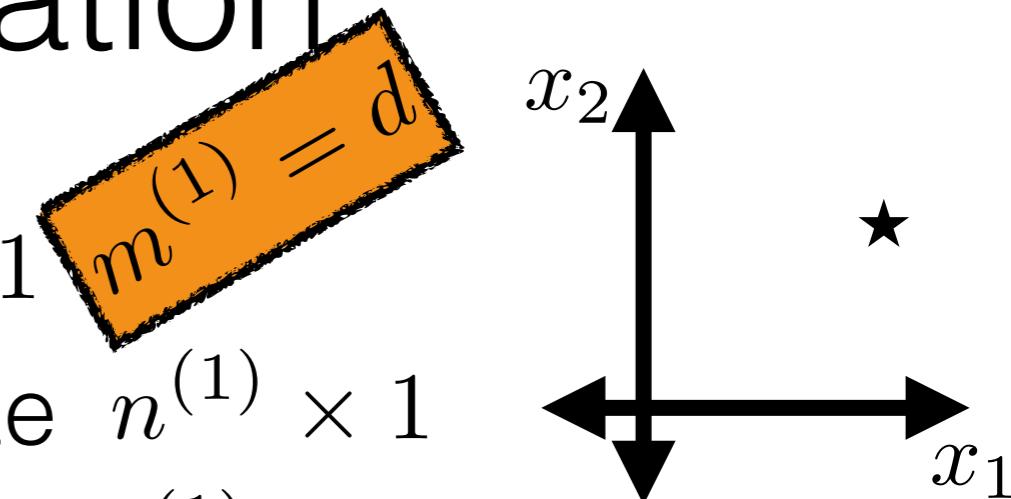
$$\begin{aligned} z &= \theta^\top \phi(x) + \theta_0 \\ &= \theta_1 \phi_1(x) + \theta_2 \phi_2(x) \\ &\quad + \theta_3 \phi_3(x) + \theta_0 \\ &= 1 \cdot \phi_1(x) + 1 \cdot \phi_2(x) \\ &\quad + 1 \cdot \phi_3(x) + (-0.5) \end{aligned}$$



# Let's get some new notation

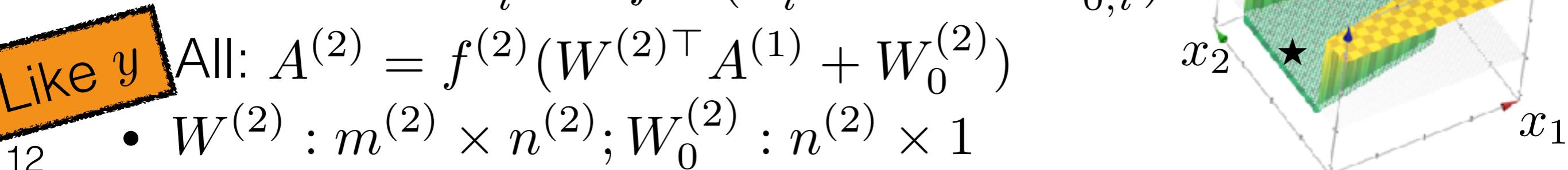
- 1st layer, constructing the features:

- Input  $x$  (a data point): size  $m^{(1)} \times 1$
- Output  $A^{(1)}$  (vector of features): size  $n^{(1)} \times 1$
- The  $i$ th feature:  $A_i^{(1)} = f^{(1)}(w_i^{(1)\top} x + w_{0,i}^{(1)})$
- All the features at once:
  - $A^{(1)} = f^{(1)}(W^{(1)\top} x + W_0^{(1)})$
  - $W^{(1)} : m^{(1)} \times n^{(1)}; W_0^{(1)} : n^{(1)} \times 1$



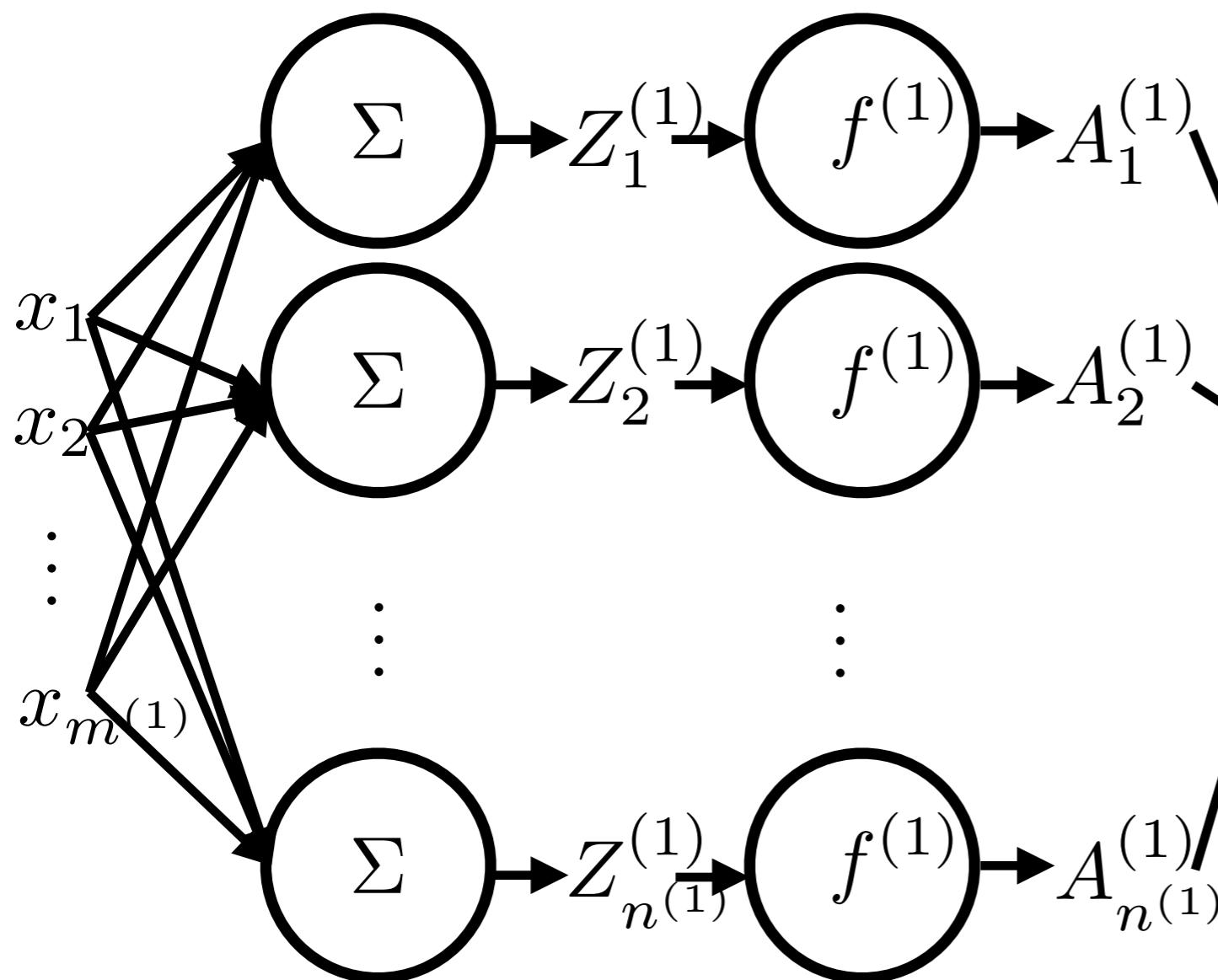
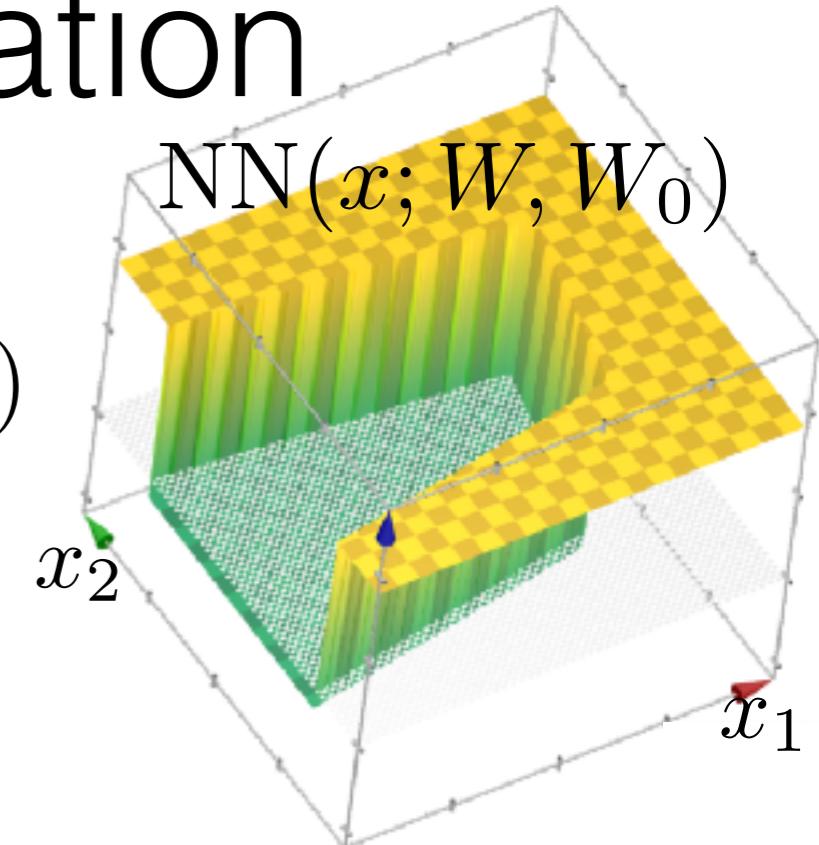
- 2nd layer, assigning a label (or labels):

- Input (the features): size  $m^{(2)} \times 1$
- Output  $A^{(2)}$  (vector of labels): size  $n^{(2)} \times 1$
- The  $i$ th label:  $A_i^{(2)} = f^{(2)}(w_i^{(2)\top} A^{(1)} + w_{0,i}^{(2)})$
- All:  $A^{(2)} = f^{(2)}(W^{(2)\top} A^{(1)} + W_0^{(2)})$
- $W^{(2)} : m^{(2)} \times n^{(2)}; W_0^{(2)} : n^{(2)} \times 1$



# Function graph representation

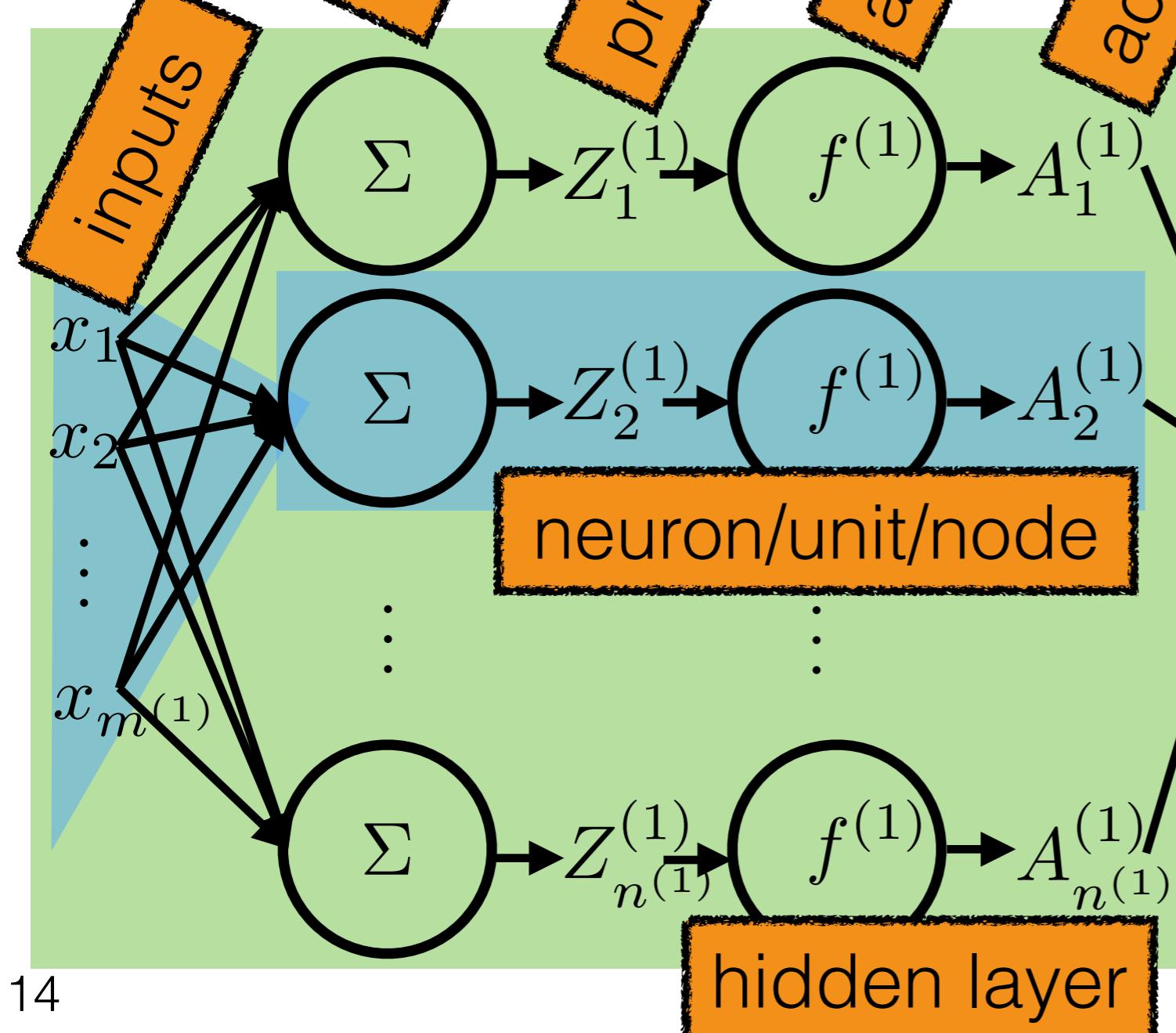
- 1st layer:  $A^{(1)} = f^{(1)}(W^{(1)\top} x + W_0^{(1)})$
- 2nd layer:  $A^{(2)} = f^{(2)}(W^{(2)\top} A^{(1)} + W_0^{(2)})$
- Whole thing:  $A^{(2)} = \text{NN}(x; W, W_0)$



- Circle: function evaluation
- Forward vs. backward
- A feed-forward neural network

# Function graph representation

- 1st layer:  $A_1^{(1)} = W^{(1)} \sigma(W^{(1)} x + b^{(1)})$
- 2nd layer:  $A_2^{(1)} = W^{(2)} \sigma(W^{(2)} A_1^{(1)} + b^{(2)})$
- Whole thing:  $A^{(2)} = \text{NN}(x; W, W_0)$



- Circle: function evaluation
- Box: neuron/unit/node
- Box: output layer
- Forward vs. backward
- A feed-forward neural network
- Fully connected

# Problem setup

- # layer  $\ell$  inputs  $m^{(\ell)}$
- # layer  $\ell$  outputs  $n^{(\ell)}$

1. Choose a hypothesis class. E.g.,

$$h(x; W, W_0) = \text{NN}(x; W, W_0)$$

- 1st layer:  $A^{(1)} = f^{(1)}(W^{(1)\top} x + W_0^{(1)})$
- 2nd layer:  $A^{(2)} = f^{(2)}(W^{(2)\top} A^{(1)} + W_0^{(2)})$

$n^{(\ell)}$  is NOT the number of data points

2. Choose a loss. (E.g. for classification: 0-1 loss, asymmetric, NLL.) Set up an objective.
3. Learn the parameters. E.g. gradient descent or SGD
4. Predict on new data using these parameters

# Problem setup

- # layer  $\ell$  inputs  $m^{(\ell)}$
- # layer  $\ell$  outputs  $n^{(\ell)}$

1. Choose a hypothesis class. E.g.,

$$h(x; W, W_0) = \text{NN}(x; W, W_0)$$

- 1st layer:  $A^{(1)} = f^{(1)}(W^{(1)\top} x + W_0^{(1)})$
- 2nd layer:  $A^{(2)} = f^{(2)}(W^{(2)\top} A^{(1)} + W_0^{(2)})$

2. Choose a loss. (E.g. for classification: 0-1 loss, asymmetric, NLL.) Set up an objective.

3. Learn the parameters. E.g. gradient descent or SGD

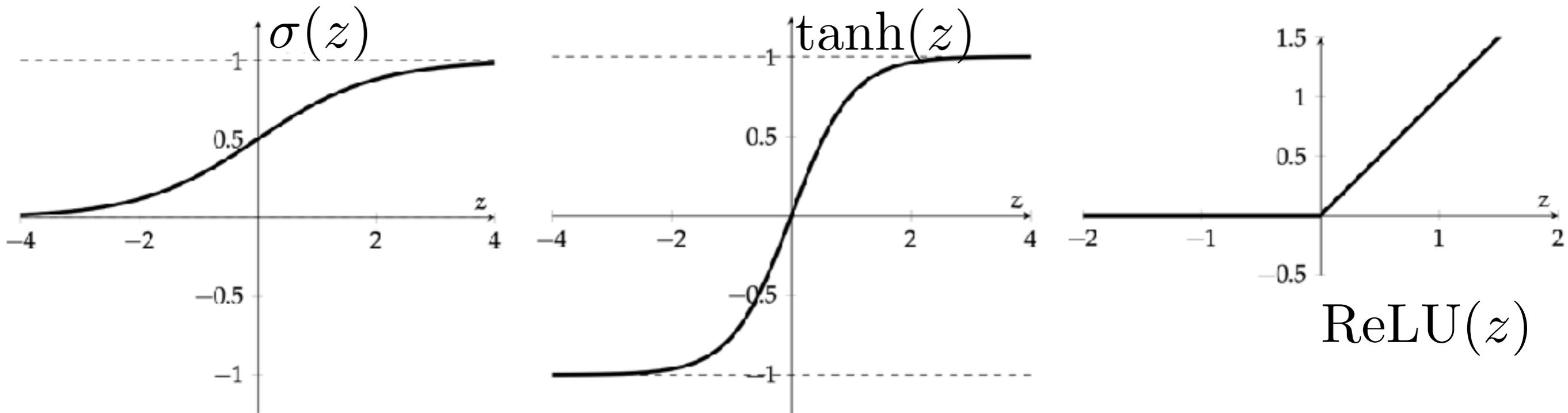
4. Predict on new data using these parameters

Issues:

- What if I want to do regression or use NLL loss?
- Derivatives (of loss with respect to parameters) are zero (or undefined) if we use step function activation
  - So (S)GD won't do what we want
- How to compute the derivatives in (S)GD?

# Different activation functions

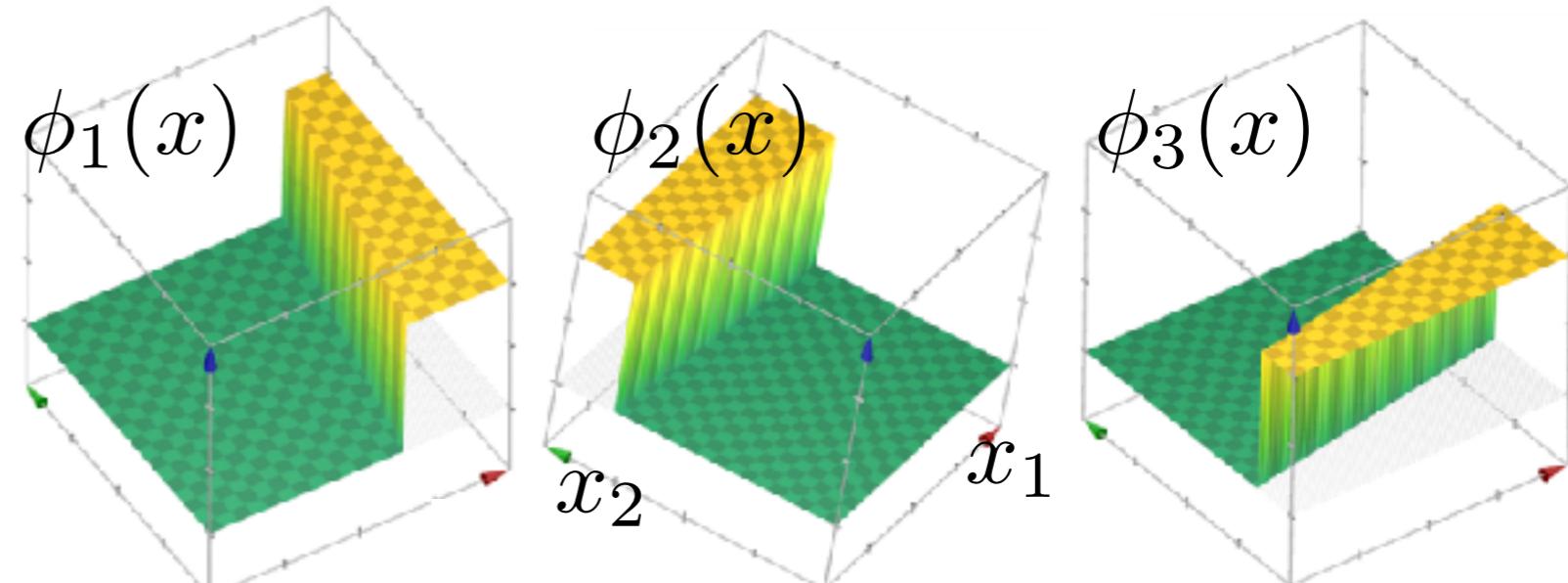
1. Hypotheses  $h(x; W, W_0) = \text{NN}(x; W, W_0)$ 
  - 1st layer:  $A^{(1)} = f^{(1)}(W^{(1)\top} x + W_0^{(1)})$
  - 2nd layer:  $A^{(2)} = f^{(2)}(W^{(2)\top} A^{(1)} + W_0^{(2)})$
- What if I want to do regression?  $f^{(2)}(z) = z$
- What if I want to use NLL loss?  $f^{(2)}(z) = \sigma(z)$
- Need non-zero derivatives for (S)GD: Above &  $f^{(1)}(z) =$



# Choices of activation function

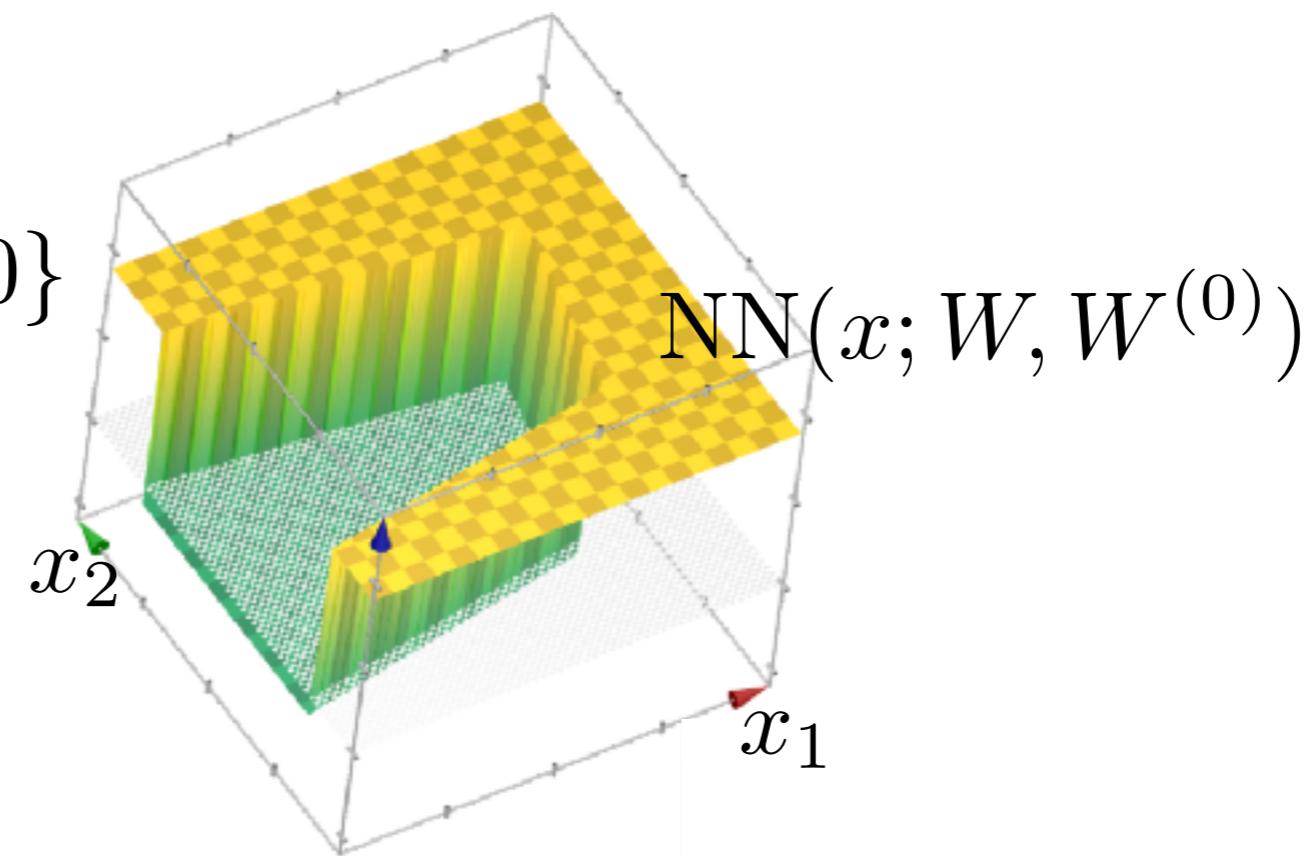
- 1st layer:  $A_i^{(1)} = f^{(1)}(w_i^{(1)\top} x + w_0^{(1)})$

- Choose  
 $f^{(1)}(z) = \mathbf{1}\{z \geq 0\}$



- 2nd layer:  $A_i^{(2)} = f^{(2)}(w_i^{(2)\top} A^{(1)} + w_0^{(2)})$

- Choose  
 $f^{(2)}(z) = \mathbf{1}\{z \geq 0\}$

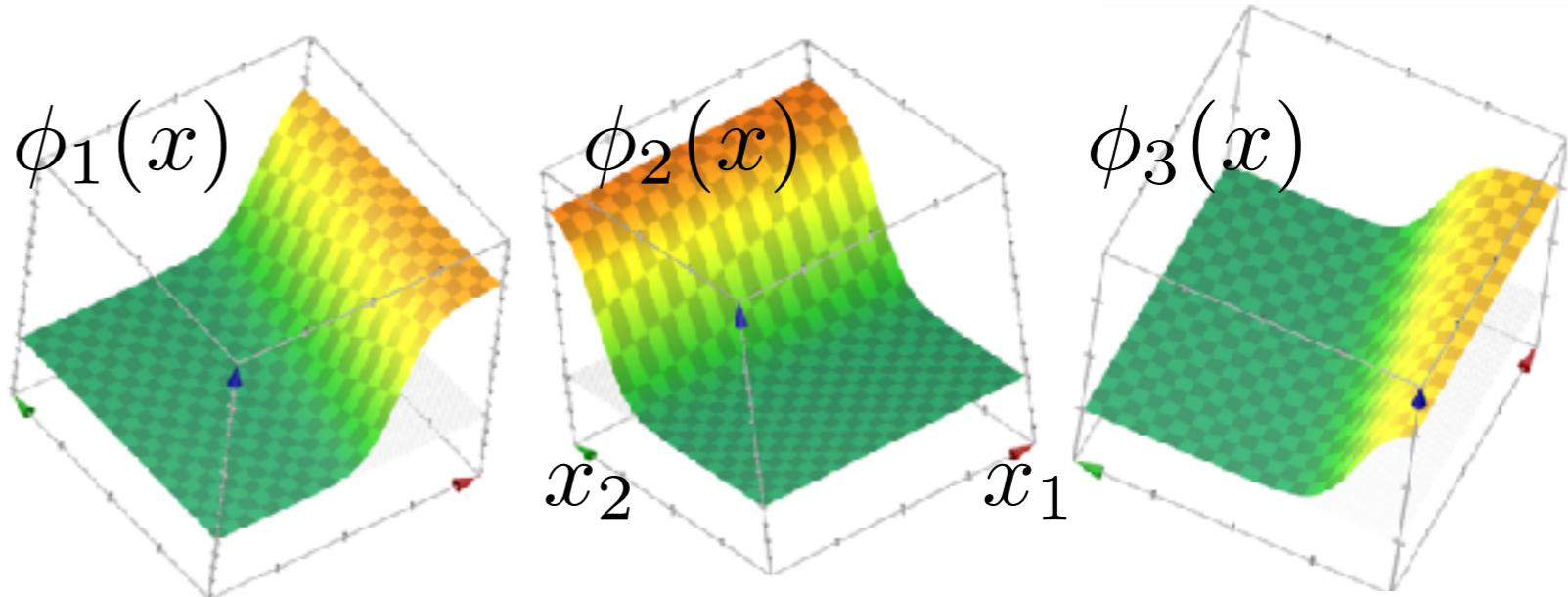


# Choices of activation function

- 1st layer:  $A_i^{(1)} = f^{(1)}(w_i^{(1)\top} x + w_0^{(1)})$

- Choose

$$f^{(1)}(z) = \sigma(z)$$



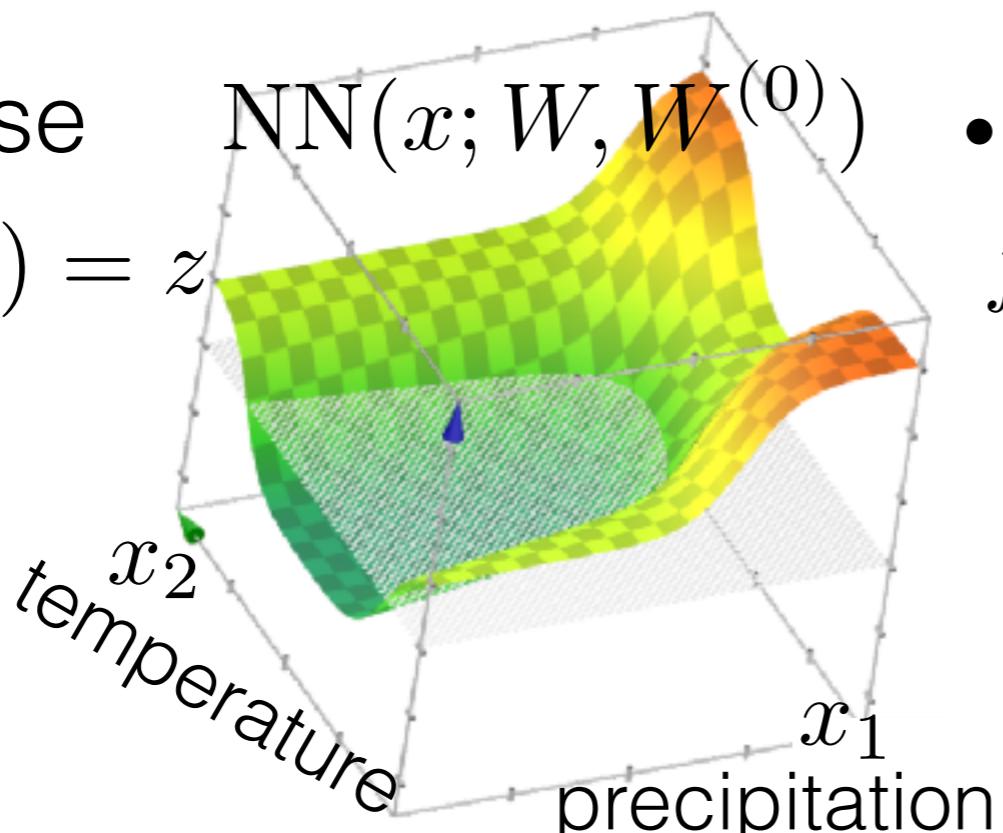
- 2nd layer:  $A_i^{(2)} = f^{(2)}(w_i^{(2)\top} A^{(1)} + w_0^{(2)})$

Probability that I  
won't run today

- Choose

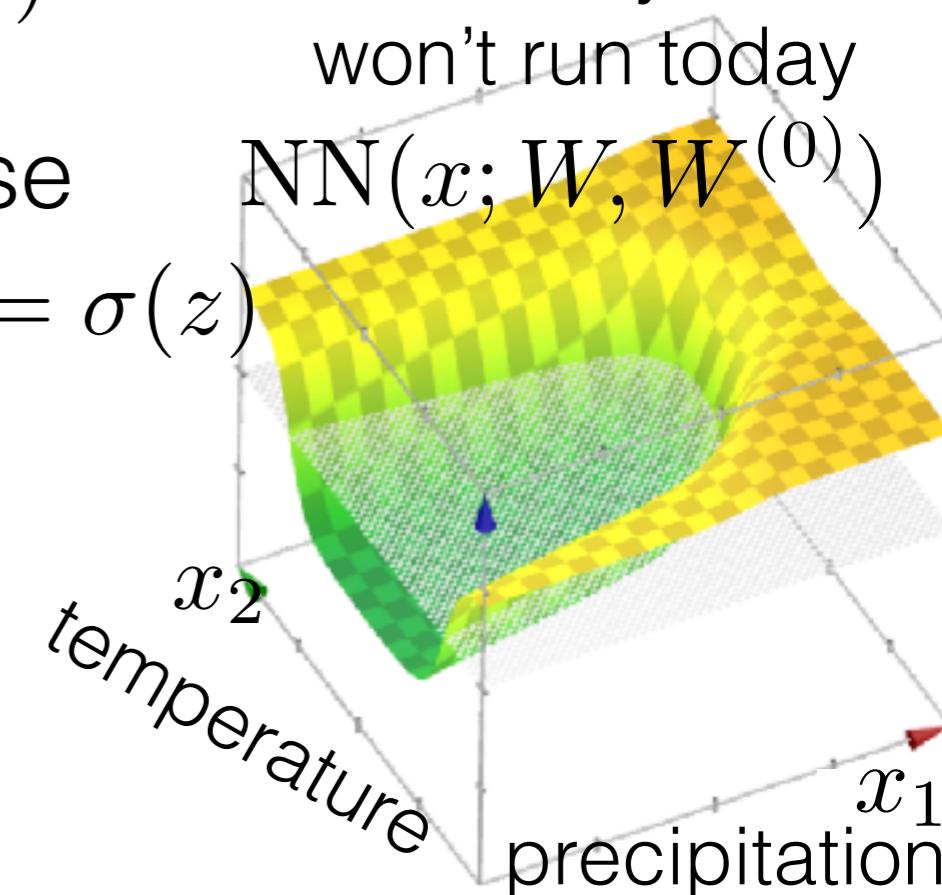
$$\text{NN}(x; W, W^{(0)})$$

How much  
will I dislike  
running  
today?



- Choose

$$f^{(2)}(z) = \sigma(z)$$



# Derivatives: Notation

- Objective:  $J(W, W_0) = \frac{1}{n_{\text{data}}} \sum_{i=1}^{n_{\text{data}}} \underbrace{L(\text{NN}(x^{(i)}; W, W_0), y^{(i)})}_{\text{loss}_i}$
- For SGD, we'll want derivatives of  $\text{loss}_i$  w.r.t. parameters
- Convention: for  $v$  of size  $1 \times 1$  and  $u$  of size  $1 \times 1$   
 $\partial v / \partial u$  is the (scalar) partial derivative of  $v$  w.r.t.  $u$ 
  - Example:  $\partial \text{loss}_i / \partial W_{1,1}^{(1)}$  1x1  
1x1
- Convention: for  $v$  of size  $1 \times 1$  and  $u$  of size  $b \times 1$   
 $\partial v / \partial u$  is a vector of size  $b \times 1$  with  $j$ th entry  $\partial v / \partial u_j$ 
  - Example:  $\partial \text{loss}_i / \partial W_0^{(1)}$  n<sup>1</sup>x1  
n<sup>1</sup>x1  $= \nabla_{W_0^{(1)}} \text{loss}_i$
- Convention: for  $v$  of size  $c \times 1$  and  $u$  of size  $b \times 1$   
 $\partial v / \partial u$  is a matrix of size  $b \times c$  with  $(j, k)$  entry  $\partial v_k / \partial u_j$ 
  - Example:  $\partial A^{(1)} / \partial Z^{(1)}$  n<sup>1</sup>xn<sup>1</sup>  
n<sup>1</sup>x1 n<sup>1</sup>x1
- Convention: for  $v$  of size  $1 \times 1$  and  $u$  of size  $b \times c$   
 $\partial v / \partial u$  is a matrix of size  $b \times c$  with  $(j, k)$  entry  $\partial v / \partial u_{j,k}$ 
  - Example:  $\partial \text{loss}_i / \partial W^{(1)}$  m<sup>1</sup>xn<sup>1</sup>  
1x1 m<sup>1</sup>xn<sup>1</sup>

# Taking derivatives for SGD

- Example 1-layer NN

- data dimension  $m^{(1)} = 4$  ; # outputs  $n^{(1)} = 3$   
$$A^{(1)} = \underbrace{W^{(1)\top}}_{4 \times 3} x + \underbrace{W_0^{(1)}}_{4 \times 1}$$
  
$$3 \times 1 \quad 4 \times 3 \quad 4 \times 1 \quad 3 \times 1$$

$$\text{loss} = L(A^{(1)}, y) = \sum_{k=1}^3 (A_k^{(1)} - y_k)^2$$

- Part of SGD step  $t$ : randomly choose  $(x, y)$  and update:  
value at end of step  $t$     
$$W_{0,j}^{(1)} \leftarrow W_{0,j}^{(1)} - \eta(t) \frac{\partial \text{loss}}{\partial W_{0,j}^{(1)}}$$
    notice: we're dropping  $i$  all values on this side are from step  $t-1$

- Use the scalar chain rule from calculus:

$$\frac{\partial \text{loss}}{\partial W_{0,j}^{(1)}} = \sum_{k=1}^3 \frac{\partial \text{loss}}{\partial A_k^{(1)}} \frac{\partial A_k^{(1)}}{\partial W_{0,j}^{(1)}} = \sum_{k=1}^3 \frac{\partial A_k^{(1)}}{\partial W_{0,j}^{(1)}} \frac{\partial \text{loss}}{\partial A_k^{(1)}}$$

$$\frac{\partial \text{loss}}{\partial W_{0,j}^{(1)}} = \frac{\partial A^{(1)}}{\partial W_{0,j}^{(1)}} \frac{\partial \text{loss}}{\partial A^{(1)}}$$

$$\frac{\partial \text{loss}}{\partial W_0^{(1)}} = \frac{\partial A^{(1)}}{\partial W_0^{(1)}} \frac{\partial \text{loss}}{\partial A^{(1)}}$$

# Taking derivatives for SGD

- Example 1-layer NN

- data dimension  $m^{(1)} = 4$ ; # outputs  $n^{(1)} = 3$   
$$A^{(1)} = \underbrace{W^{(1)\top}}_{4 \times 3} x + \underbrace{W_0^{(1)}}_{4 \times 1} \quad 3 \times 1$$

$$\text{loss} = L(A^{(1)}, y) = \sum_{k=1}^3 (A_k^{(1)} - y_k)^2$$

- Part of SGD step  $t$ : randomly choose  $(x, y)$  and update:  
value at end of step  $t$  
$$W_{0,j}^{(1)} \leftarrow W_{0,j}^{(1)} - \eta(t) \frac{\partial \text{loss}}{\partial W_{0,j}^{(1)}}$$
 notice: we're dropping  $i$  all values on this side are from step  $t-1$

- Use the scalar chain rule from calculus:

$$\frac{\partial \text{loss}}{\partial W_{0,j}^{(1)}} = \sum_{k=1}^3 \frac{\partial \text{loss}}{\partial A_k^{(1)}} \frac{\partial A_k^{(1)}}{\partial W_{0,j}^{(1)}} = \sum_{k=1}^3 \frac{\partial A_k^{(1)}}{\partial W_{0,j}^{(1)}} \frac{\partial \text{loss}}{\partial A_k^{(1)}}$$

$$\frac{\partial \text{loss}}{\partial W_{0,j}^{(1)}} = \frac{\partial A^{(1)}}{\partial W_{0,j}^{(1)}} \frac{\partial \text{loss}}{\partial A^{(1)}}$$

$$\frac{\partial \text{loss}}{\partial W_0^{(1)}} = \frac{\partial A^{(1)}}{\partial W_0^{(1)}} \frac{\partial \text{loss}}{\partial A^{(1)}}$$

Exercise:  
compute  
this matrix

# Taking derivatives for SGD

- Example 1-layer NN

- data dimension  $m^{(1)} = 4$ ; # outputs  $n^{(1)} = 3$   
$$A^{(1)} = \underbrace{W^{(1)\top}}_{4 \times 3} x + \underbrace{W_0^{(1)}}_{4 \times 1} \quad 3 \times 1$$

$$\text{1x1 loss} = L(A^{(1)}, y) = \sum_{k=1}^3 (A_k^{(1)} - y_k)^2$$

- Part of SGD step  $t$ : randomly choose  $(x, y)$  and update:  
value at end of step  $t$  
$$W_{0,j}^{(1)} \leftarrow W_{0,j}^{(1)} - \eta(t) \frac{\partial \text{loss}}{\partial W_{0,j}^{(1)}}$$
 notice: we're dropping  $i$  all values on this side are from step  $t-1$

- Use the scalar chain rule from calculus:

$$\frac{\partial \text{loss}}{\partial W^{(1)}} = \sum_{j=1}^3 \frac{\partial \text{loss}}{\partial A_j^{(1)}} \frac{\partial A_j^{(1)}}{\partial W^{(1)}} = \sum_{k=1}^3 \frac{\partial A_k^{(1)}}{\partial W_{0,j}^{(1)}} \frac{\partial \text{loss}}{\partial A_k^{(1)}}$$

Notice: We can't just substitute the weight **matrix** in the final formula. (Why not?)

Exercise:  
compute  
this matrix

$$\frac{\partial \text{loss}}{\partial W_0^{(1)}} = \frac{\partial A^{(1)}}{\partial W_0^{(1)}} \frac{\partial \text{loss}}{\partial A^{(1)}}$$

# Taking derivatives for SGD

- Example 1-layer NN

- data dimension  $m^{(1)} = 4$ ; # outputs  $n^{(1)} = 3$   
$$A^{(1)} = \underbrace{W^{(1)\top}}_{4 \times 3} x + \underbrace{W_0^{(1)}}_{4 \times 1}$$
  
$$3 \times 1 \quad 4 \times 3 \quad 4 \times 1 \quad 3 \times 1$$

$$\text{loss} = L(A^{(1)}, y) = \sum_{k=1}^3 (A_k^{(1)} - y_k)^2$$

- Our last derivation doesn't work for a full weight matrix.  
(Exercise: why not?) So what about the full matrix case?

$$\frac{\partial \text{loss}}{\partial W_{j,k}^{(1)}} = \sum_{p=1}^3 \frac{\partial A_p^{(1)}}{\partial W_{j,k}^{(1)}} \frac{\partial \text{loss}}{\partial A_p^{(1)}} = \frac{\partial A_k^{(1)}}{\partial W_{j,k}^{(1)}} \frac{\partial \text{loss}}{\partial A_k^{(1)}} = x_j \frac{\partial \text{loss}}{\partial A_k^{(1)}}$$

- Observe: if  $u$  is a  $c \times 1$  vector and  $v$  is a  $b \times 1$  vector,  $vu^\top$  is a matrix of size  $b \times c$  with  $(j, k)$  entry  $v_j u_k$
- By our convention,  $\partial \text{loss} / \partial W^{(1)}$  is a matrix with  $(j, k)$  entry  $\partial \text{loss} / \partial W_{j,k}^{(1)}$

- So: 
$$\frac{\partial \text{loss}}{\partial W^{(1)}} = x \left( \frac{\partial \text{loss}}{\partial A^{(1)}} \right)^\top$$
  
$$4 \times 3 \quad 4 \times 1 \quad 3 \times 1$$

# Taking derivatives for SGD

- General case: NN with  $L$  layers

- Layer  $\ell$  has  $m^{(\ell)}$  inputs and  $n^{(\ell)} = m^{(\ell+1)}$  outputs

$$A^{(\ell)} = f^\ell(Z^{(\ell)}), Z^{(\ell)} = W^{(\ell)\top} A^{(\ell-1)} + W_0^{(\ell)} \quad \text{with} \quad A^{(0)} = x$$

- To find  $\partial \text{loss} / \partial W^{(\ell)}$ , break off the “final” weight derivative

$$\frac{\partial \text{loss}}{\partial W_{j,k}^{(\ell)}} = \frac{\partial Z_k^{(\ell)}}{\partial W_{j,k}^{(\ell)}} \frac{\partial \text{loss}}{\partial Z_k^{(\ell)}} = A_j^{(\ell-1)} \frac{\partial \text{loss}}{\partial Z_k^{(\ell)}} \rightarrow \frac{\partial \text{loss}}{\partial W^{(\ell)}} = A^{(\ell-1)} \underbrace{\left( \frac{\partial \text{loss}}{\partial Z^{(\ell)}} \right)^\top}_{m^\ell \times n^\ell}$$

$1 \times 1 \quad 1 \times 1 \quad 1 \times 1 \quad 1 \times 1 \quad 1 \times 1 \quad m^\ell \times n^\ell \quad m^\ell \times 1 \quad 1 \times n^\ell$

- It remains to find  $\partial \text{loss} / \partial Z^{(\ell)}$

$$\frac{\partial \text{loss}}{\partial Z^{(\ell)}} = \frac{\partial A^{(\ell)}}{\partial Z^{(\ell)}} \frac{\partial Z^{(\ell+1)}}{\partial A^{(\ell)}} \frac{\partial A^{(\ell+1)}}{\partial Z^{(\ell+1)}} \cdots \frac{\partial A^{(L-1)}}{\partial Z^{(L-1)}} \frac{\partial Z^{(L)}}{\partial A^{(L-1)}} \frac{\partial A^{(L)}}{\partial Z^{(L)}} \frac{\partial \text{loss}}{\partial A^{(L)}}$$

$n^\ell \times 1 \quad n^\ell \times n^\ell \quad n^\ell \times n^{\ell+1} \quad n^{\ell+1} \times n^{\ell+1} \quad n^{L-1} \times n^{L-1} \quad n^{L-1} \times n^L \quad n^L \times n^L \quad n^L \times 1$

- Lots of this computation will be shared across layers
- More efficient to compute the shared parts only once

- Can similarly show:  $\partial \text{loss} / \partial W_0^{(\ell)} = \partial \text{loss} / \partial Z^{(\ell)}$

# Taking derivatives for SGD

- Exercise: use these formulas to show that the final derivatives w.r.t.  $W^\ell$  are 0 if any layer above  $\ell$  has step function activations

$$\frac{\partial \text{loss}}{\partial W_{j,k}^{(\ell)}} = \frac{\partial Z_k^{(\ell)}}{\partial W_{j,k}^{(\ell)}} \frac{\partial \text{loss}}{\partial Z_k^{(\ell)}} = A_j^{(\ell-1)} \frac{\partial \text{loss}}{\partial Z_k^{(\ell)}}$$

1x1      1x1      1x1

Exercise: Check that you agree with these formulas!

$$A^{(\ell-1)} + W_0^{(\ell)} \quad \text{with} \quad A^{(0)} = x$$

iff the “final” weight derivative

$$\frac{\partial \text{loss}}{\partial W^{(\ell)}} = A^{(\ell-1)} \underbrace{\left( \frac{\partial \text{loss}}{\partial Z^{(\ell)}} \right)^T}_{\substack{m^\ell \times n^\ell \\ m^\ell \times 1 \\ 1 \times n^\ell}}$$

- It remains to find  $\partial \text{loss} / \partial Z^{(\ell)}$

$$\frac{\partial \text{loss}}{\partial Z^{(\ell)}} = \frac{\partial A^{(\ell)}}{\partial Z^{(\ell)}} \frac{\partial Z^{(\ell+1)}}{\partial A^{(\ell)}} \frac{\partial A^{(\ell+1)}}{\partial Z^{(\ell+1)}} \cdots \frac{\partial A^{(L-1)}}{\partial Z^{(L-1)}} \frac{\partial Z^{(L)}}{\partial A^{(L-1)}} \frac{\partial A^{(L)}}{\partial Z^{(L)}} \frac{\partial \text{loss}}{\partial A^{(L)}}$$

$n^\ell \times 1$        $n^\ell \times n^\ell$        $n^\ell \times n^{\ell+1}$        $n^{\ell+1} \times n^{\ell+1}$        $n^{L-1} \times n^{L-1}$        $n^{L-1} \times n^L$        $n^L \times n^L$        $n^L \times 1$

- Lots of this computation will be shared across layers
- More efficient to compute the shared parts only once

- Can similarly show:  $\partial \text{loss} / \partial W_0^{(\ell)} = \partial \text{loss} / \partial Z^{(\ell)}$