

#### Outline

#### Last time: Neural Networks

- Three-neuron network + optimization
- Neural network structure (adding more neurons)
- Forward pass and Backpropagation (scalar version)

Rec: multivariate

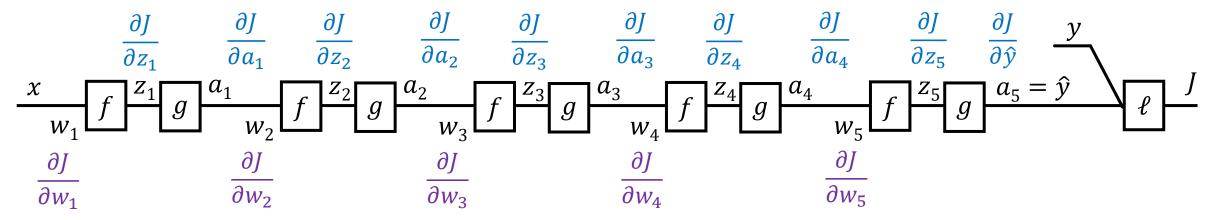
**Today: Neural Networks** 

- Calculus: multi-variate chain rule
- Backpropagation (vector version)
- Neural Network Properties and Intuition

## Backpropagation

$$\hat{y} = h_{\theta}(\mathbf{x}) = g\left(w_5 \cdot g\left(w_4 \cdot g\left(w_3 \cdot g\left(w_2 \cdot g\left(w_1 \cdot x\right)\right)\right)\right)\right)$$

Width 1 deep network (no bias) (dumb but will help with calculus)



To do gradient descent we need the partial derivative of the objective with respect to each parameter,  $\mathbf{w}_\ell \leftarrow \mathbf{w}_\ell - \alpha \ \partial J/\partial \mathbf{w}_\ell$ 

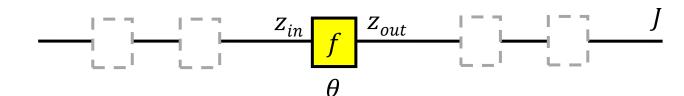
The backward pass propagates the change in the objective with respect to intermediate values  $(\partial J/\partial z_{\ell})$  and  $\partial J/\partial a_{\ell}$  back through the network to produce each  $\partial J/\partial w_{\ell}$ 

## Generic Layer Implementation (so-far)

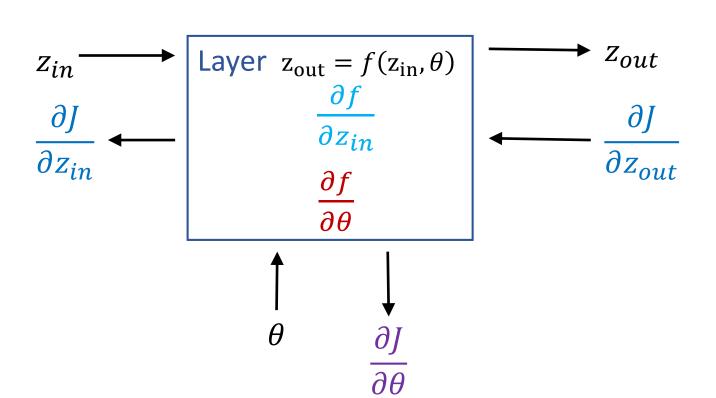
#### Compute derivatives per layer, utilizing previous derivatives

Objective:  $J(\theta)$ 

Arbitrary layer:  $z_{out} = f(z_{in}, \theta)$ 



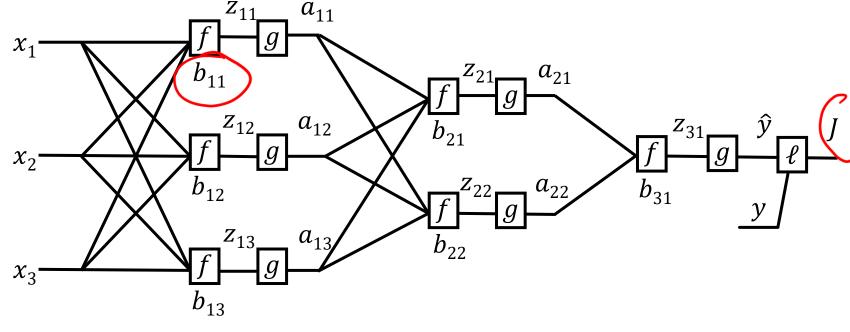
$$\frac{\partial J}{\partial \theta} = \frac{\partial J}{\partial z_{out}} \frac{\partial f}{\partial \theta}$$



## Optimization

$$\hat{y} = h_{\theta}(\mathbf{x}) = g \left( b_{3,1} + \sum_{k} w_{3,1,k} \ g \left( b_{2,k} + \sum_{i} w_{2,k,i} \ g \left( b_{1,i} + \sum_{j} w_{1,i,j} \ x_{j} \right) \right) \right)$$

Tons of repeated partial derivatives



$$\frac{\partial J}{\partial b_{3,1}} = \frac{\partial \ell}{\partial \hat{y}} \frac{\partial a_{3,1}}{\partial z_{3,1}} \frac{\partial z_{3,1}}{\partial b_{3,1}}$$

$$\frac{\partial J}{\partial b_{2,i}} = \frac{\partial \ell}{\partial \hat{y}} \frac{\partial a_{3,1}}{\partial z_{3,1}} \frac{\partial z_{3,1}}{\partial a_{2,i}} \frac{\partial a_{2,i}}{\partial z_{2,i}} \frac{\partial z_{2,i}}{\partial b_{2,i}}$$

$$\frac{\partial J}{\partial b_{1,i}} = \frac{\partial \ell}{\partial \hat{y}} \frac{\partial a_{3,1}}{\partial z_{3,1}} \frac{\partial z_{3,1}}{\partial a_{2,i}} \frac{\partial a_{2,i}}{\partial z_{2,i}} \frac{\partial z_{2,i}}{\partial b_{2,i}}$$

$$\frac{\partial J}{\partial b_{1,i}} = \frac{\partial \ell}{\partial \hat{y}} \frac{\partial a_{3,1}}{\partial z_{3,1}} \frac{\partial z_{3,1}}{\partial a_{2}} \frac{\partial a_{2,i}}{\partial z_{2}} \frac{\partial z_{2,i}}{\partial a_{1,i}} \frac{\partial a_{1,i}}{\partial z_{1,i}} \frac{\partial z_{1,i}}{\partial b_{1,i}}$$

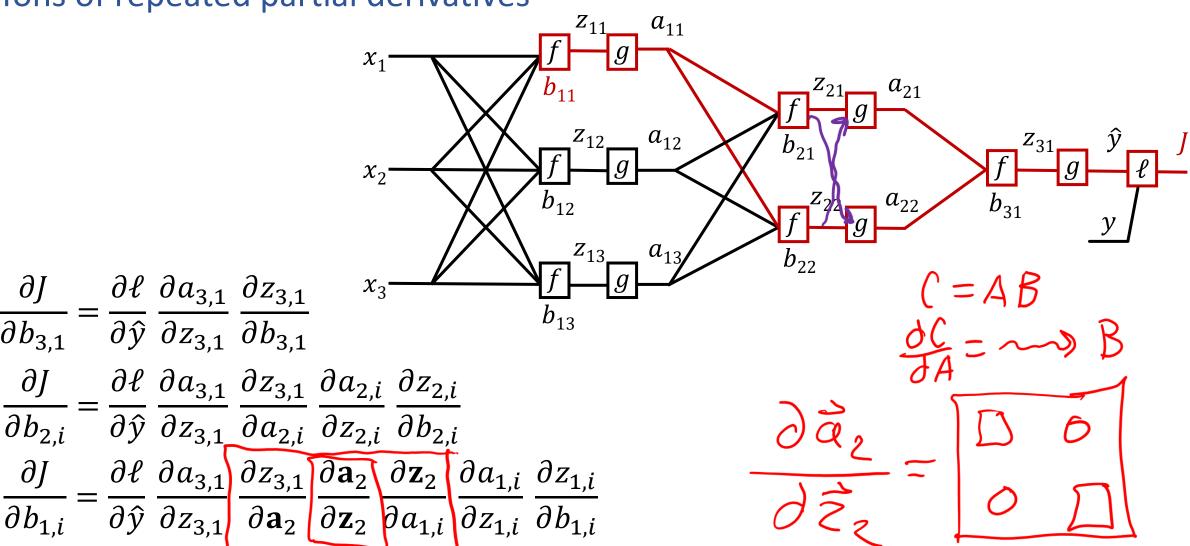
## Optimization

 $\partial \ell \ \partial a_{3,1} \ \partial z_{3,1}$ 

 $\partial b_{3,1} - \partial \hat{y} \partial z_{3,1} \partial b_{3,1}$ 

$$\hat{y} = h_{\theta}(\mathbf{x}) = g \left( b_{3,1} + \sum_{k} w_{3,1,k} \ g \left( b_{2,k} + \sum_{i} w_{2,k,i} \ g \left( b_{1,i} + \sum_{j} w_{1,i,j} \ x_{j} \right) \right) \right)$$

Tons of repeated partial derivatives



## Calculus Time-out

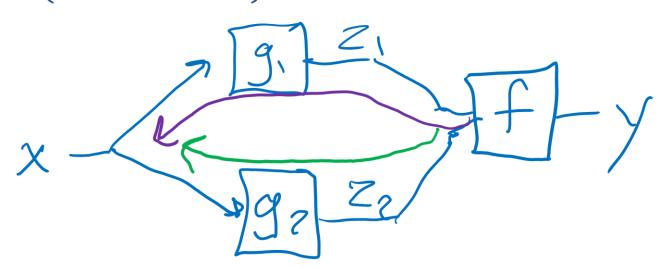
## Multivariate-chain rule

## Multivariable Chain Rule

(Note: math notation)

$$g_1(x) = 3x$$
  
 $g_2(x) = 5x$   
 $f(z_1, z_2) = 2z_1 + 7z_2$   
 $y = f(g_1(x), g_2(x))$ 

dx = dx dz + dy dzz dx, dx dz dx



### Exercise: Multivariable Chain Rule

$$z_1 = g_1(x) = \sin(x)$$

$$z_2 = g_2(x) = x^3$$

$$y = f(z_1, z_2) = z_1^4 e^{z_2} + 5z_1 + 7z_2$$

$$\frac{df}{dx} = \frac{\partial f}{\partial g_1} \frac{\partial g_1}{\partial x} + \frac{\partial f}{\partial g_2} \frac{\partial g_2}{\partial x}$$

## Exercise: Multivariable Chain Rule

$$\frac{df}{dx} = \frac{\partial f}{\partial g_1} \frac{\partial g_1}{\partial x} + \frac{\partial f}{\partial g_2} \frac{\partial g_2}{\partial x}$$

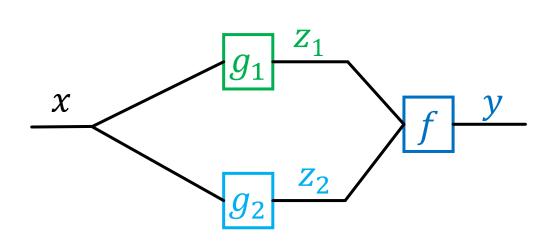
$$z_1 = g_1(x) = \sin(x)$$
  $dg_1/dx = \cos(x)$   
 $z_2 = g_2(x) = x^3$   $dg_2/dx = 3x^2$   
 $y = f(z_1, z_2) = z_1^4 e^{z_2} + 5z_1 + 7z_2$   $\partial f/\partial z_1 = 4z_1^3 e^{z_2} + 5$ 

$$dg_1/dx = \cos(x)$$

$$dg_2/dx = 3x^2$$

$$\partial f/\partial z_1 = 4z_1^3 e^{z_2} + 5$$

$$\partial f/\partial z_2 = z_1^4 e^{z_2} + 7$$



$$\frac{df}{dx} = \frac{\partial f}{\partial g_1} \frac{\partial g_1}{\partial x} + \frac{\partial f}{\partial g_2} \frac{\partial g_2}{\partial x}$$
$$= (4z_1^3 e^{z_2} + 5) \cos(x)$$
$$+ (z_1^4 e^{z_2} + 7) 3x^2$$

### Calculus Chain Rule

#### Scalar:

$$y = f(z)$$

$$z = g(x)$$

$$\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx}$$

#### Multivariate:

$$y = f(\mathbf{z})$$

$$\mathbf{z} = g(x)$$

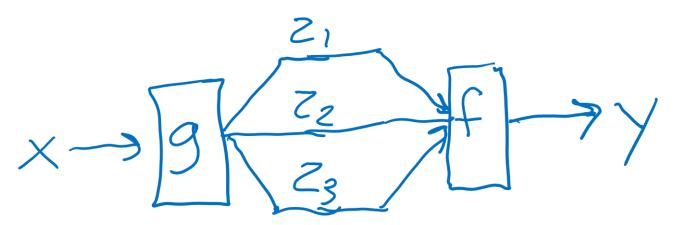
$$\frac{dy}{dx} = \sum_{j} \frac{\partial y}{\partial z_{j}} \frac{\partial z_{j}}{\partial x}$$

#### Multivariate:

$$\mathbf{y} = f(\mathbf{z})$$

$$\mathbf{z} = g(\mathbf{x})$$

$$\frac{dy_{i}}{dx_{k}} = \sum_{j} \frac{\partial y_{i}}{\partial z_{j}} \frac{\partial z_{j}}{\partial x_{k}}$$



## Calculus Chain Rule



#### Scalar:

$$y = f(z)$$

$$z = g(x)$$

$$\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx}$$

#### Multivariate:

$$y = f(\mathbf{z})$$

$$\mathbf{z} = g(x)$$

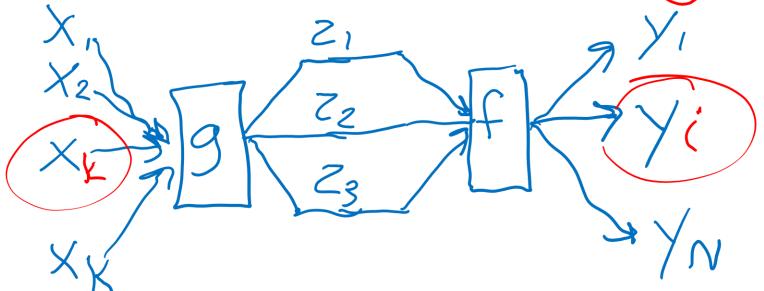
$$\frac{dy}{dx} = \sum_{j} \frac{\partial y}{\partial z_{j}} \frac{\partial z_{j}}{\partial x}$$

#### Multivariate:

$$\mathbf{y} = f(\mathbf{z})$$

$$\mathbf{z} = g(\mathbf{x})$$

$$\frac{dy_i}{dx_k} = \sum_j \frac{\partial y_i}{\partial z_j} \frac{\partial z_j}{\partial x_k}$$



## Multivariable Chain Rule

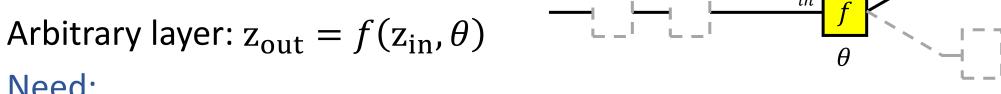
	Numerator layout	Denominator layout	Notes
$\frac{d}{dt} f(g(t), h(t))$	$\frac{df}{dg}\frac{dg}{dt} + \frac{df}{dh}\frac{dh}{dt}$	Same	$f: (\mathbb{R} \times \mathbb{R}) \to \mathbb{R}, \ t \in \mathbb{R}$ $g: \mathbb{R} \to \mathbb{R}, \ h: \mathbb{R} \to \mathbb{R}$
$\frac{d}{dt} f(g_1(t), \dots, g_N(t))$	$\sum_{i=1}^{N} \frac{df}{dg_i} \frac{dg_i}{dt}$	Same	$h: (\mathbb{R} \times \cdots \times \mathbb{R}) \to \mathbb{R}, \ t \in \mathbb{R}$ $f_i: \mathbb{R} \to \mathbb{R} \ \forall i \in \{1, \dots, N\}$
$\frac{d}{dt} f(\mathbf{g}(t))$	$\frac{\partial f}{\partial t} = \frac{\partial f}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial t}$ $= \square \square$	$\frac{\partial \mathbf{f}}{\partial t} \frac{\partial \mathbf{g}}{\partial \mathbf{g}}$ $= \square$	$f: \mathbb{R}^N \to \mathbb{R}, \ \mathbf{g}: \mathbb{R} \to \mathbb{R}^N$ $t \in \mathbb{R}$
$\frac{\partial}{\partial \mathbf{v}} f(\mathbf{g}(\mathbf{v}))$	$\frac{\partial f}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial \mathbf{v}}$	$\frac{\partial f}{\partial \mathbf{v}} = \frac{\partial \mathbf{g}}{\partial \mathbf{v}} \frac{\partial f}{\partial \mathbf{g}}$ $H = HHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHH$	$f: \mathbb{R}^N \to \mathbb{R}, \ \mathbf{g}: \mathbb{R}^M \to \mathbb{R}^N$ $\mathbf{v} \in \mathbb{R}^M$
$\frac{\partial}{\partial \mathbf{v}} \ f(\mathbf{g}(\mathbf{v}), \mathbf{h}(\mathbf{v}))$	$\frac{\partial f}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial \mathbf{v}} + \frac{\partial f}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{v}}$	$\frac{\partial \mathbf{g}}{\partial \mathbf{v}} \frac{\partial f}{\partial \mathbf{g}} + \frac{\partial \mathbf{h}}{\partial \mathbf{v}} \frac{\partial f}{\partial \mathbf{h}}$	$h: \mathbb{R}^K \times \mathbb{R}^N \to \mathbb{R},  \mathbf{v} \in \mathbb{R}^M$ $f: \mathbb{R}^M \to \mathbb{R}^K,  \mathbf{g}: \mathbb{R}^M \to \mathbb{R}^N$

# Backpropagation (vector version)

## Generic Layer Implementation (updated)

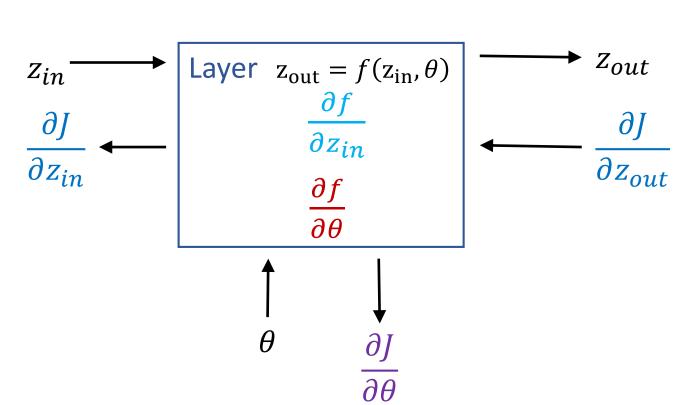
Compute derivatives per layer, utilizing previous derivatives

Objective:  $I(\theta)$ 



$$= \frac{\partial J}{\partial z_{in}} + = \frac{\partial J}{\partial z_{out}} \frac{\partial f}{\partial z_{in}}$$

$$= \frac{\partial J}{\partial \theta} + = \frac{\partial J}{\partial z_{out}} \frac{\partial f}{\partial \theta}$$



## Generic Layer Implementation (vector output version)

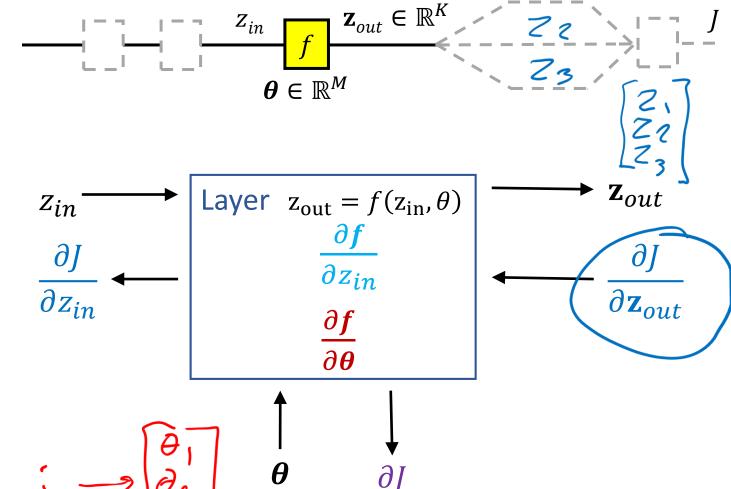
Compute derivatives per layer, utilizing previous derivatives

Objective: 
$$J(\theta)$$

Arbitrary layer: 
$$z_{out} = f(z_{in}, \theta)$$

$$\frac{\partial J}{\partial z_{in}} = \sum_{k=1}^{3} \frac{\partial J}{\partial z_{k}} \frac{\partial Z_{k}}{\partial z_{in}}$$

$$\frac{\partial J}{\partial \theta_{j}} = \sum_{k=1}^{3} \frac{\partial J}{\partial Z_{k}} \frac{\partial Z_{k}}{\partial \theta_{j}}$$

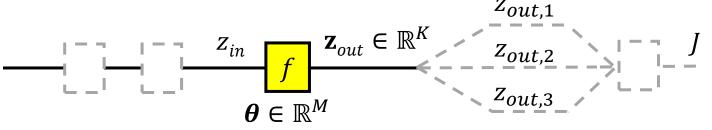


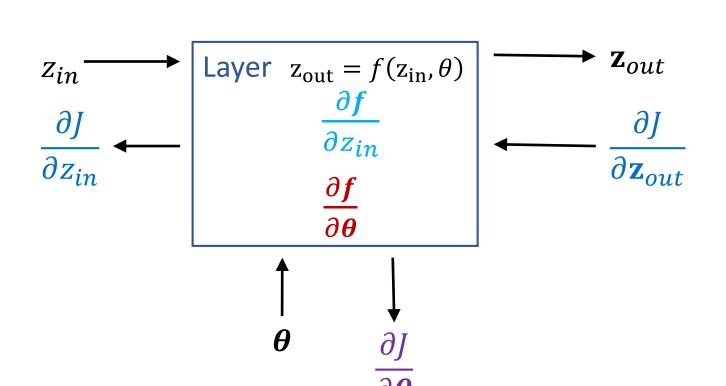
## Generic Layer Implementation (vector output version)

Compute derivatives per layer, utilizing previous derivatives

Objective:  $J(\theta)$ 

Arbitrary layer:  $z_{out} = f(z_{in}, \theta)$ 





# $\mathbf{z} \in \mathbb{R}^{K} \quad \mathbf{z} \in \mathbb{R}^{K}$ $W \in \mathbb{R}^{K \times M}, \mathbf{b} \in \mathbb{R}^{K}$

#### Compute derivatives per layer

Objective: 
$$J(\theta)$$

Layer: 
$$\mathbf{z} = f(\mathbf{a}, W, \mathbf{b}) = W\mathbf{a} + \mathbf{b}$$

Scalar version (no formatting)

$$\frac{\partial J}{\partial a_{j}} = \begin{cases} \begin{cases} \frac{\partial J}{\partial z_{k}} & \frac{\partial Z}{\partial a_{j}} \\ \frac{\partial J}{\partial w_{i,j}} & = \\ \frac{\partial J}{\partial b_{i}} & = \\ \end{cases} \begin{cases} \frac{\partial J}{\partial z_{k}} & \frac{\partial Z}{\partial z_{k}} \\ \frac{\partial J}{\partial z_{k}} & \frac{\partial J}{\partial z_{k}} \end{cases} \begin{cases} \frac{\partial J}{\partial z_{k}} & \frac{\partial Z}{\partial z_{k}} \\ \frac{\partial J}{\partial z_{k}} & \frac{\partial J}{\partial z_{k}} & \frac{\partial J}{\partial z_{k}} \end{cases}$$

Layer 
$$\mathbf{a} = f(\mathbf{a}, W, \mathbf{b})$$

$$\frac{\partial f}{\partial \mathbf{a}}$$

$$\frac{\partial f}{\partial W} \frac{\partial f}{\partial \mathbf{b}}$$

$$\uparrow \qquad \qquad \downarrow \qquad \qquad \downarrow$$

$$W \mathbf{b} \qquad \frac{\partial J}{\partial W} \frac{\partial J}{\partial \mathbf{b}}$$

$$\mathbf{z} = W\mathbf{a} + \mathbf{b}$$

$$z_k = b_k + \sum_j W_{k,j} a_j$$

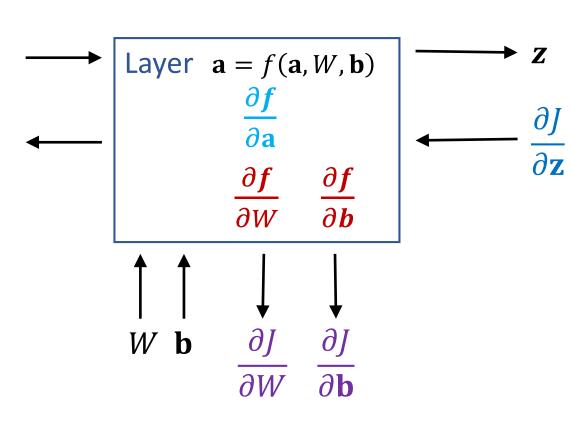
# $\mathbf{z} \in \mathbb{R}^{K} \quad \mathbf{z} \in \mathbb{R}^{K}$ $W \in \mathbb{R}^{K \times M}, \mathbf{b} \in \mathbb{R}^{K}$

#### Compute derivatives per layer

Objective: 
$$J(\theta)$$

Layer: 
$$\mathbf{z} = f(\mathbf{a}, W, \mathbf{b}) = W\mathbf{a} + \mathbf{b}$$

Scalar version (no formatting)



$$\mathbf{z} = W\mathbf{a} + \mathbf{b}$$

$$z_k = b_k + \sum_j W_{k,j} a_j$$

# $\mathbf{z} \in \mathbb{R}^{K}$ $W \in \mathbb{R}^{K \times M}, \mathbf{b} \in \mathbb{R}^{K}$

#### Compute derivatives per layer

Objective: 
$$J(\theta)$$

Layer: 
$$\mathbf{z} = f(\mathbf{a}, W, \mathbf{b}) = W\mathbf{a} + \mathbf{b}$$

#### (Denominator format)

$$\frac{\partial J}{\partial a} = \frac{\partial \hat{z}}{\partial \hat{a}} \frac{\partial \hat{y}}{\partial \hat{z}}$$

$$\frac{\partial J}{\partial W} = \frac{\partial \dot{z}}{\partial \dot{w}} \frac{\partial \dot{J}}{\partial \dot{z}}$$

$$\frac{\partial J}{\partial \mathbf{b}} = \frac{\partial \dot{z}}{\partial \dot{b}} \frac{\partial \dot{z}}{\partial \dot{z}}$$

$$\frac{\partial J}{\partial a}$$

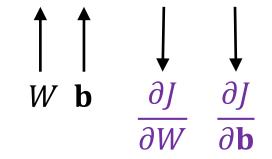




Layer 
$$\mathbf{z} = f(\mathbf{a}, W, \mathbf{b})$$

$$\frac{\partial f}{\partial \mathbf{a}}$$

$$\frac{\partial f}{\partial \mathbf{z}}$$



$$W \leftarrow W - \alpha \frac{\partial M}{\partial A}$$

#### 

#### Compute derivatives per layer

Objective: 
$$J(\theta)$$

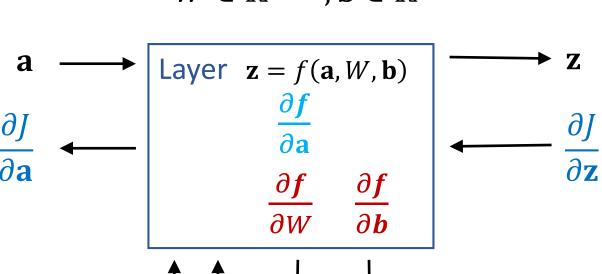
Layer: 
$$\mathbf{z} = f(\mathbf{a}, W, \mathbf{b}) = W\mathbf{a} + \mathbf{b}$$

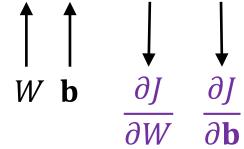
(Numerator format)

$$\frac{\partial J}{\partial a}$$

$$\blacksquare \frac{\partial J}{\partial W}$$

$$\frac{\partial J}{\partial \mathbf{b}}$$





#### Outline

#### Last time: Neural Networks

- Three-neuron network + optimization
- Neural network structure (adding more neurons)
- Forward pass and Backpropagation (scalar version)

#### Today: Neural Networks

- Calculus: multi-variate chain rule
- Backpropagation (vector version)
  - Neural Network Properties and Intuition

# Neural Network Properties

## Neural Networks Properties

#### **Practical considerations**

- Large number of neurons
  - Danger for overfitting
- Modeling assumptions vs data assumptions trade-off
- Gradient descent can easily get stuck local optima

#### What if there are no non-linear activations?

 A deep neural network with only linear layers can be reduced to an exactly equivalent single linear layer

#### **Universal Approximation Theorem:**

 A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.

## Neural Networks Properties

Non-linearity

Fitting complex functions

Fitting any! function (universal approximation theorem)

Overfitting

# Neural Network Properties

Non-linearity

## Non-linearity

Neural network prediction function,  $\hat{y} = h(x)$ , is definitely not linear. The non-linear activation functions provide this non-linearity

#### What if there are no non-linear activations?

 A deep neural network with only linear layers can be reduced to an exactly equivalent single linear layer

#### Example

What happens with you add to linear functions together?

- Adding two lines?  $h(x) = f_1(x, w_1, b_1) + f_2(x, w_2, b_2)$
- Adding two planes?

## Objective Function is Not Convex

Objective function for	Convex?	Closed-form solution?
Linear regression	Yes	Ves
Logistic regression	Yes	$N_{o}$
Neural networks	$N_{0}$	No

## Optimization

#### Convex function

If f(x) is convex, then:

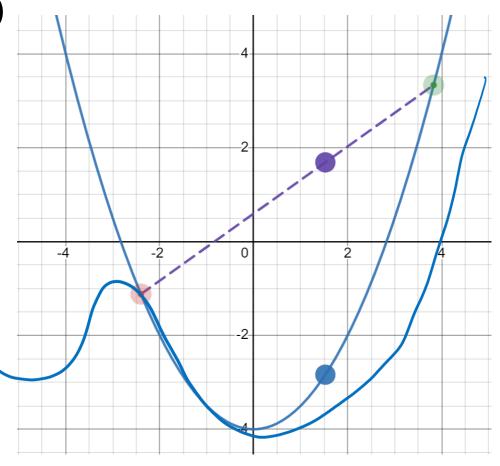
 $f(\alpha \mathbf{x} + (1 - \alpha)\mathbf{z}) \le \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{z})$   $\forall 0 \le \alpha \le 1$ 

#### Convex optimization

If second derivative is  $\geq 0$  everywhere then function is convex

If f(x) is convex, then:

■ Every local minimum is also a global minimum ©



## Non-linearity

#### **Prediction function**

Neural network prediction function,  $\hat{y} = h(x)$ , is definitely not linear. The non-linear activation functions provide this non-linearity

#### Objective function

Neural network prediction function,  $\hat{y} = h(x)$ , is definitely not linear. The non-linear activation functions provide this non-linearity

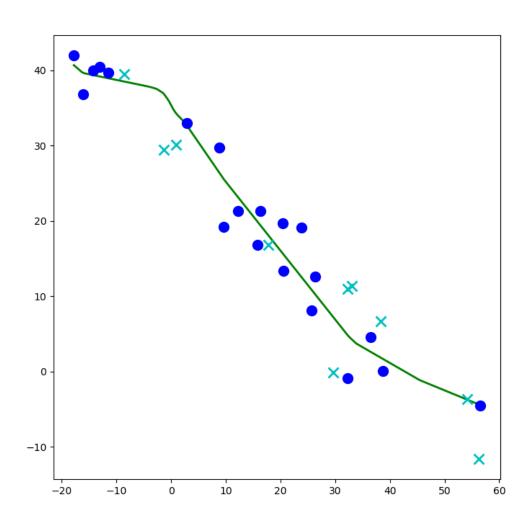
#### Just because it can fit any function will it?

- Objective function is non-convex
  - → it will get stuck in local minima
- Stochastic gradient decent take pseudorandom steps
  - → Helps pop out of local minima
- Is getting stuck at a local minima necessarily a bad thing??

# Neural Network Properties

Fitting complex functions (with 1-D input)

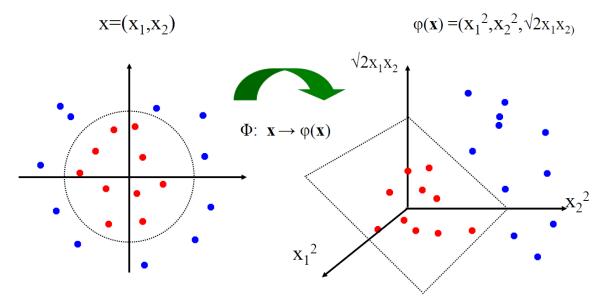
## Network to Approximate a 1-D Function



#### Linear Classifiers for Nonlinear Data

Linear classifiers have linear decision boundaries

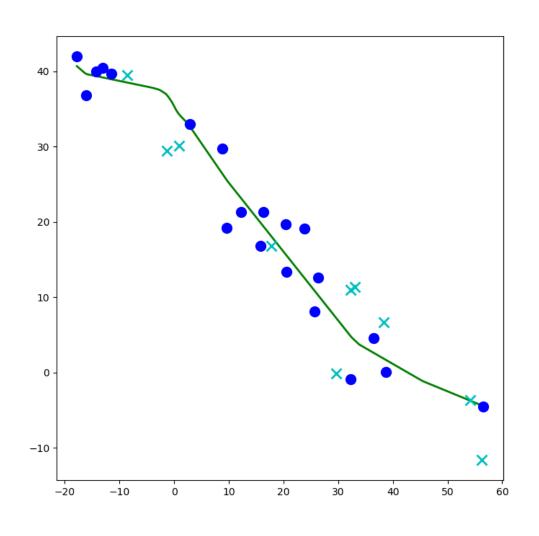
Feature mapping can convert nonlinear data to higher dimensions

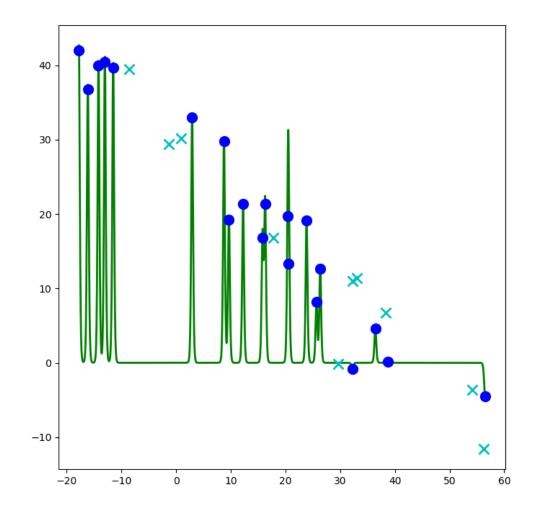


This slide is courtesy of www.iro.umontreal.ca/~pift6080/documents/papers/svm tutorial.ppt

Today, instead of choosing a feature mapping function, we'll use neural networks to learn nonlinear decision boundaries and regression functions

## Network to Approximate a 1-D Function

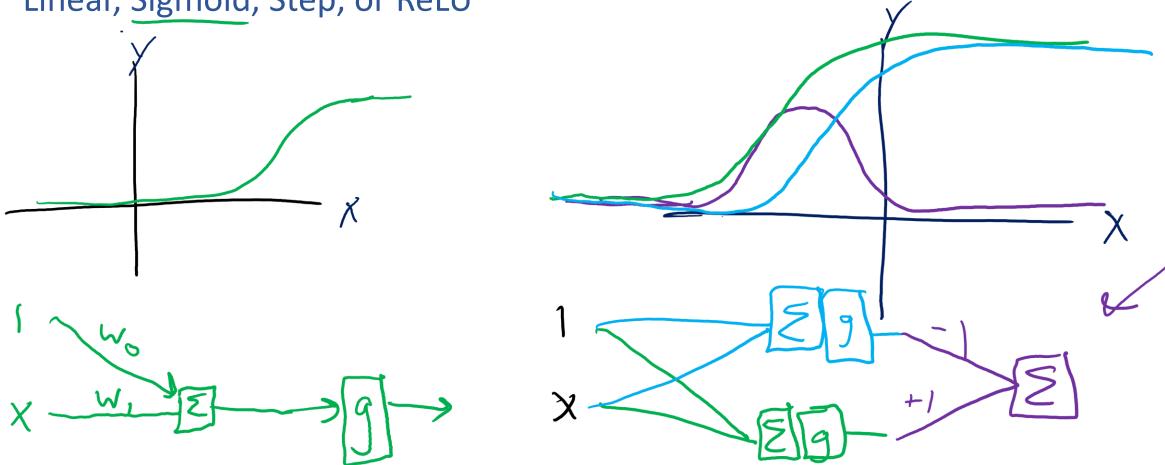




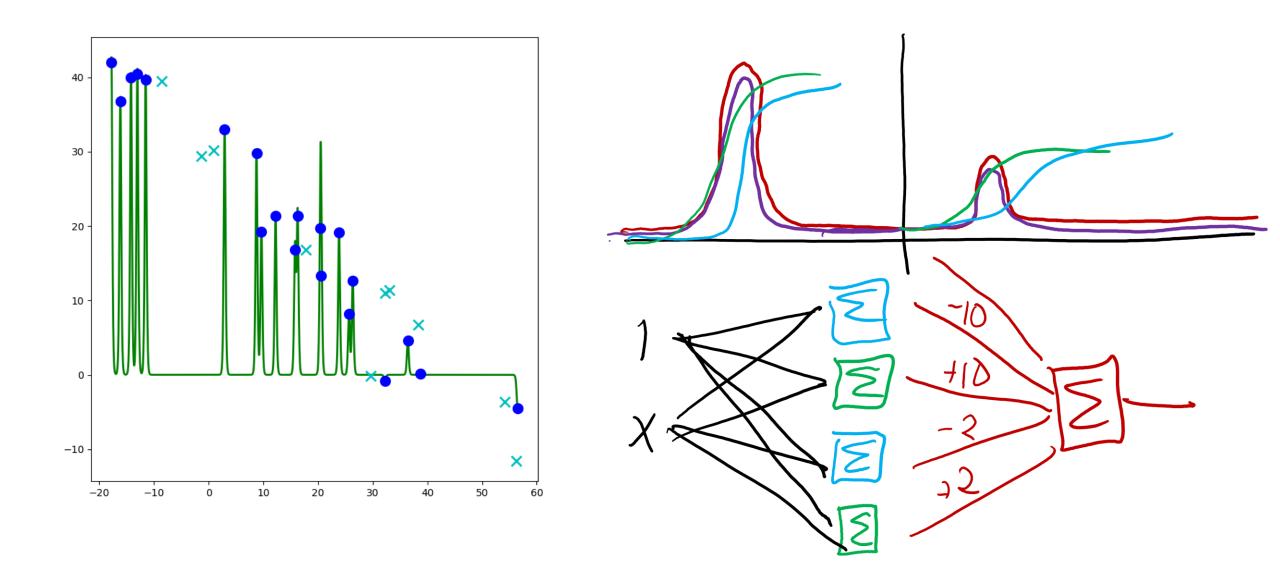
## Network to Approximate a 1-D Function

Design a network to approximate this function using:

Linear, Sigmoid, Step, or ReLU



## Network to Approximate a 1-D Function

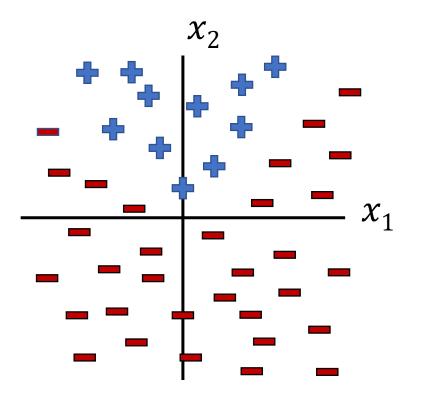


## Neural Network Properties

Fitting complicated functions (with 2-D input)

### Classification Design Challenge

How could you configure three specific perceptrons to classify this data?



$$h_A(\mathbf{x}) = sign(\mathbf{w}_A^{\mathsf{T}}\mathbf{x} + b_A)$$

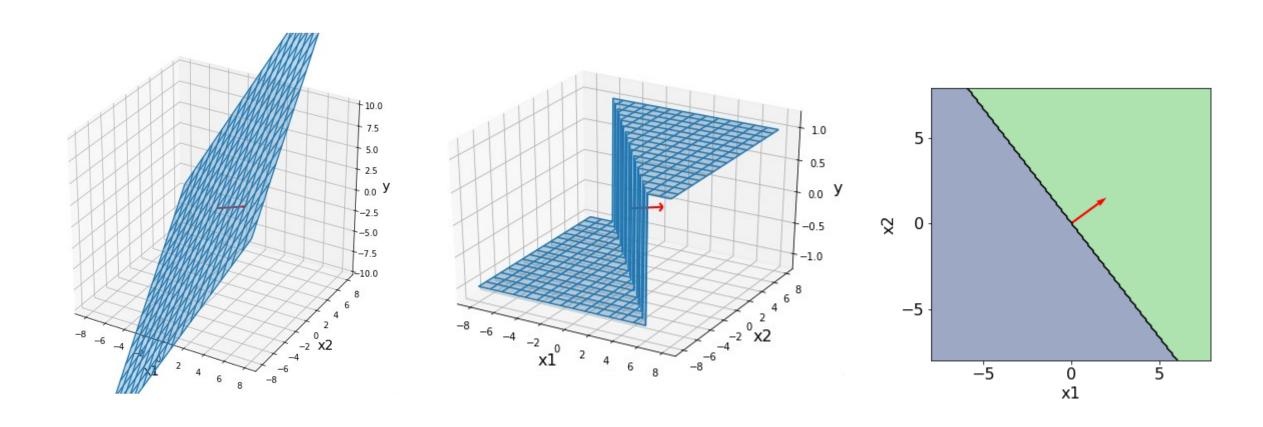
$$h_B(\mathbf{x}) = sign(\mathbf{w}_B^{\mathsf{T}}\mathbf{x} + b_B)$$

$$h_C(\mathbf{z}) = sign(\mathbf{w}_C^{\mathsf{T}}\mathbf{z} + b_C)$$

#### Perceptron

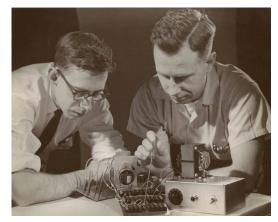
$$sign(\mathbf{z}) = \begin{cases} 1, & if \ z \ge 0 \\ -1, & if \ z < 0 \end{cases}$$

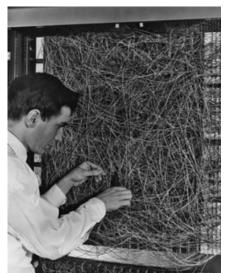
# Classification: Hard threshold on linear model $h(\mathbf{x}) = sign(\mathbf{w}^{\mathsf{T}}\mathbf{x} + b)$



#### Perceptron History

#### Frank Rosenblatt, 1957





The New Yorker, December 6, 1958 P. 44

Talk story about the perceptron, a new electronic brain which hasn't been built, but which has been successfully simulated on the I.B.M. 704. Talk with Dr. Frank Rosenblatt, of the Cornell Aeronautical Laboratory, who is one of the two men who developed the prodigy; the other man is Dr. Marshall C. Yovits, of the Office of Naval Research, in Washington. Dr. Rosenblatt defined the perceptron as the first non-biological object which will achieve an organization o its external environment in a meaningful way. It interacts with its environment, forming concepts that have not been made ready for it by a human agent. If a triangle is held up, the perceptron's eye picks up the image & conveys it along a random succession of lines to the response units, where the image is registered. It can tell the difference betw. a cat and a dog, although it wouldn't be able to tell whether the dog was to theleft or right of the cat. Right now it is of no practical use, Dr. Rosenblatt conceded, but he said that one day it might be useful to send one into outer space to take in impressions for us.

#### Exercise

# Which of the following perceptron parameters will perfectly classify this data?

$$h(\mathbf{x}) = sign(\mathbf{w}^{\mathsf{T}}\mathbf{x} + b)$$

$$sign(\mathbf{z}) = \begin{cases} 1, & if \ z \ge 0 \\ -1, & if \ z < 0 \end{cases}$$

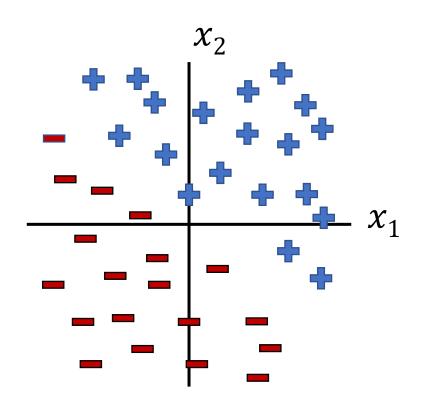
**A.** 
$$\mathbf{w} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$
,  $b = 0$ 

**B.** 
$$\mathbf{w} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}, b = 0$$

C. 
$$\mathbf{w} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$
,  $b = 0$ 

**D.** 
$$\mathbf{w} = \begin{bmatrix} -1 \\ -1 \end{bmatrix}, b = 0$$

E. None of the above



#### Poll

# Which of the following perceptron parameters will perfectly classify this data?

$$h(\mathbf{x}) = sign(\mathbf{w}^{\mathsf{T}}\mathbf{x} + b)$$

$$sign(\mathbf{z}) = \begin{cases} 1, & if \ z \ge 0 \\ -1, & if \ z < 0 \end{cases}$$

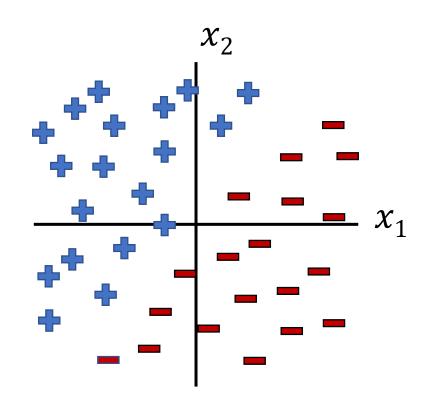
**A.** 
$$\mathbf{w} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$
,  $b = 0$ 

**B.** 
$$\mathbf{w} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}, b = 0$$

C. 
$$\mathbf{w} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$
,  $b = 0$ 

**D.** 
$$\mathbf{w} = \begin{bmatrix} -1 \\ -1 \end{bmatrix}, b = 0$$

E. None of the above



#### Poll

# Which of the following parameters of $h_C(\mathbf{z})$ will perfectly classify this data?

$$h_C(\mathbf{z}) = sign(\mathbf{w}_C^{\mathsf{T}}\mathbf{z} + b_C)$$

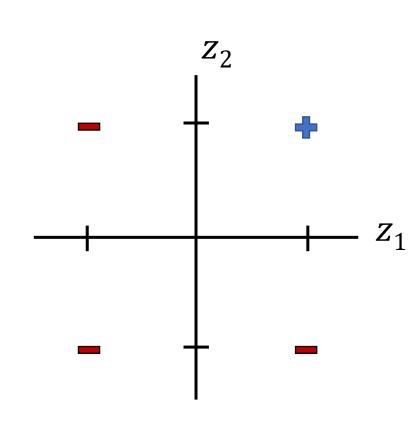
$$sign(\mathbf{x}) = \begin{cases} 1, & if \ x \ge 0 \\ -1, & if \ x < 0 \end{cases}$$

$$A. \mathbf{w}_C = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, b_C = 0$$

B. 
$$\mathbf{w}_C = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$
,  $b_C = 1$ 

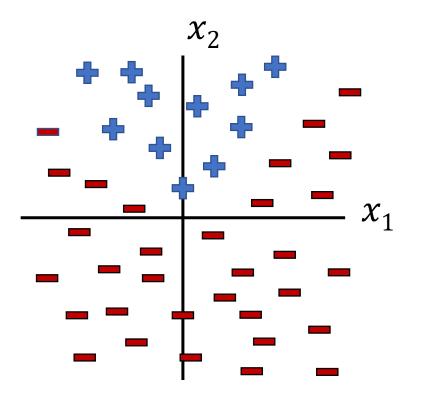
C. 
$$\mathbf{w}_C = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$
,  $b_C = -1$ 

D. None of the above



### Classification Design Challenge

How could you configure three specific perceptrons to classify this data?



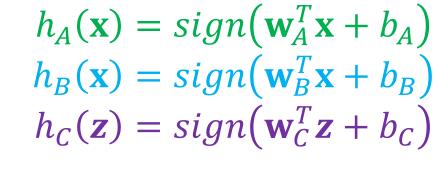
$$h_A(\mathbf{x}) = sign(\mathbf{w}_A^{\mathsf{T}}\mathbf{x} + b_A)$$

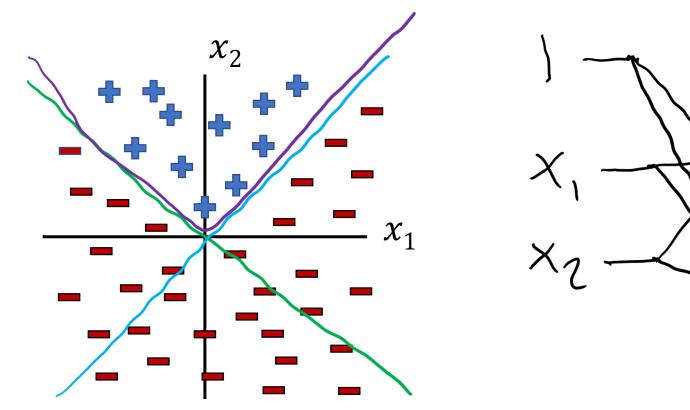
$$h_B(\mathbf{x}) = sign(\mathbf{w}_B^{\mathsf{T}}\mathbf{x} + b_B)$$

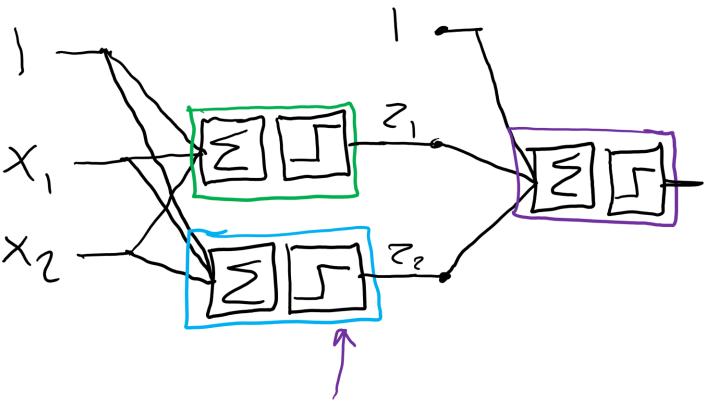
$$h_C(\mathbf{z}) = sign(\mathbf{w}_C^{\mathsf{T}}\mathbf{z} + b_C)$$

## Classification Design Challenge

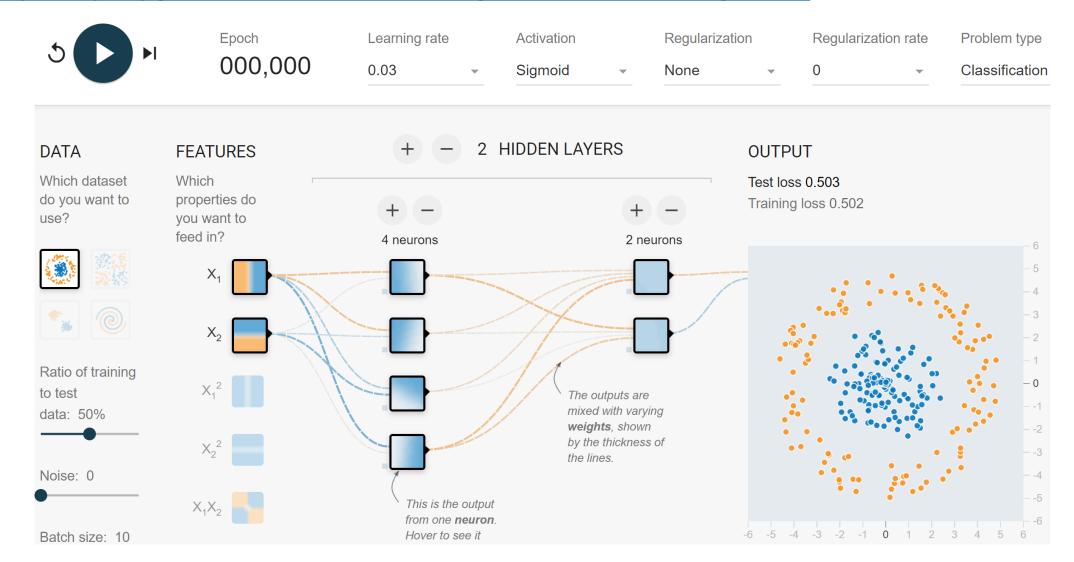
How could you configure three specific perceptrons to classify this data?



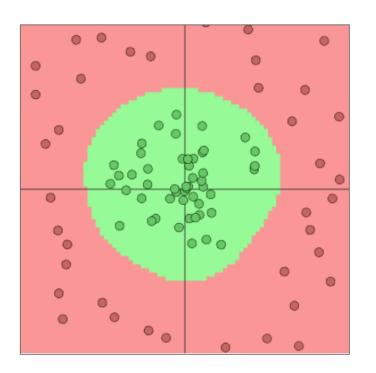




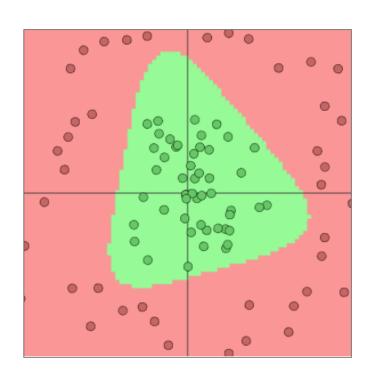
https://playground.tensorflow.org/#activation=sigmoid

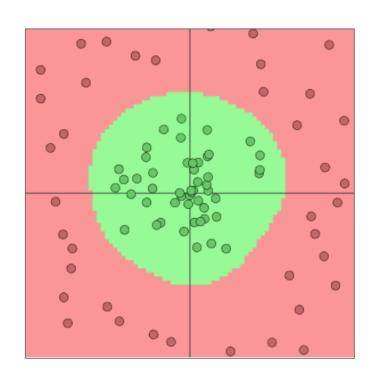


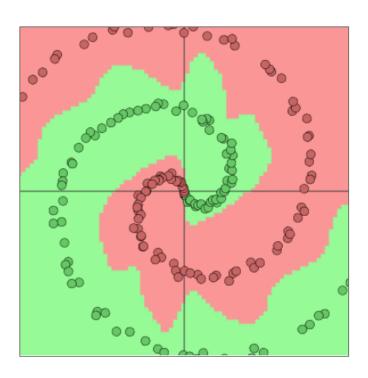
Approximate arbitrary decision boundary



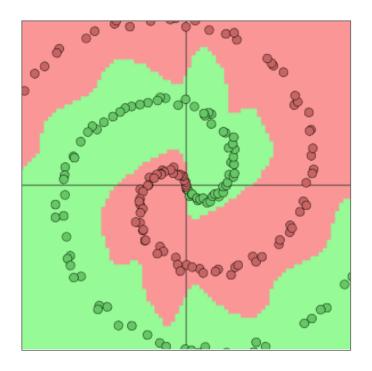
Approximate arbitrary decision boundary







Approximate arbitrary decision boundary



## Neural Network Properties

(Over?) Fitting any function

#### Neural Networks Properties

#### **Universal Approximation Theorem**

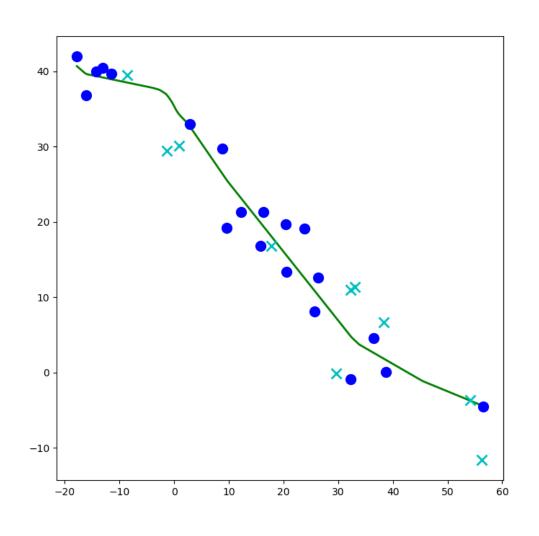
 A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy

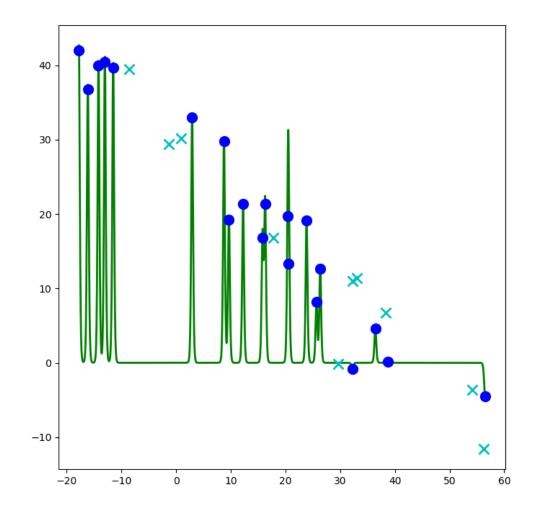
#### Reminder

#### Just because it can fit any function will it?

- Objective function is non-convex
  - → it will get stuck in local minima
- Stochastic gradient decent take pseudorandom steps
  - → Helps pop out of local minima
- Is getting stuck at a local minima necessarily a bad thing??

## Network to Approximate a 1-D Function





### Debugging Overfitting and Underfitting

#### Underfitting (check this first!)

- Evidence: poor training loss (and poor validation loss)
  - Note: Compare with human performance as a baseline, i.e., make sure the task isn't impossible for a human (with unlimited resources)
- Try: Adding more capacity to network (wider or deeper)
  - But how do we choose a new network structure?

#### Overfitting

- Evidence: good training loss, but poor validation loss
- Evidence: really large parameter values
- Try: Regularization (we'll learn about this soon!)
- Try: Adding more data

#### Summary: Neural Networks Properties

#### **Practical considerations**

- Large number of neurons
  - Danger for overfitting
- Modeling assumptions vs data assumptions trade-off
- Gradient descent can easily get stuck local optima

#### What if there are no non-linear activations?

 A deep neural network with only linear layers can be reduced to an exactly equivalent single linear layer

#### **Universal Approximation Theorem:**

 A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.