

# Announcements



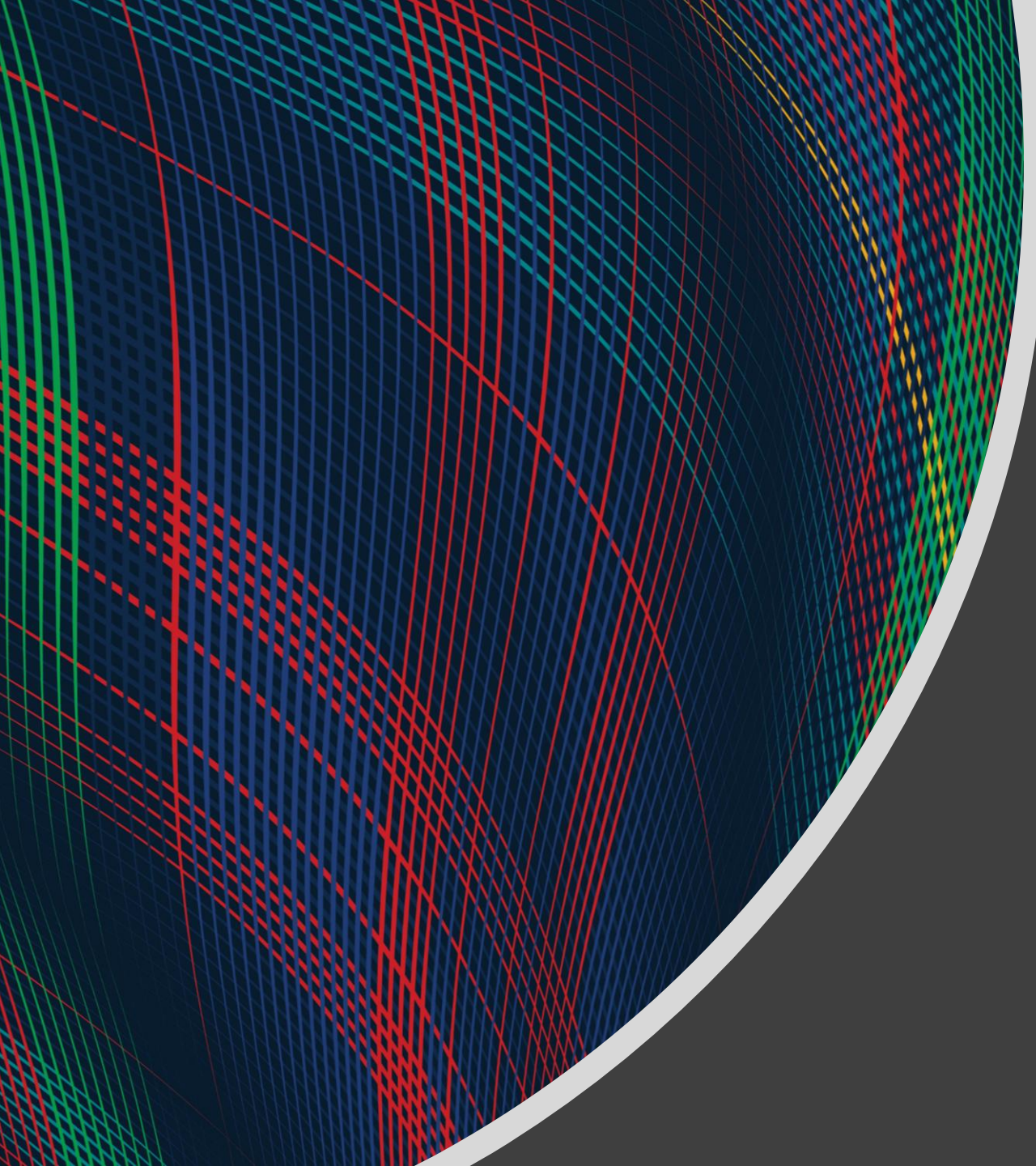
## Assignments

- HW4
  - Due Wed, 10/14, 11:59 pm
- HW5
  - Plan: Out tomorrow
  - Due Mon, 10/26, 11:59 pm

## Recitation

- No recitation the next two Fridays 😞
- We'll post a worksheet for neural nets and record a walk-through

## Survey



# Introduction to Machine Learning

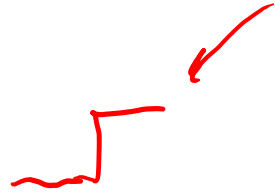
## Neural Networks

Instructor: Pat Virtue

# Plan

## Last Time

- Neural Networks
  - Perceptron
  - Multilayer perceptron



## Today

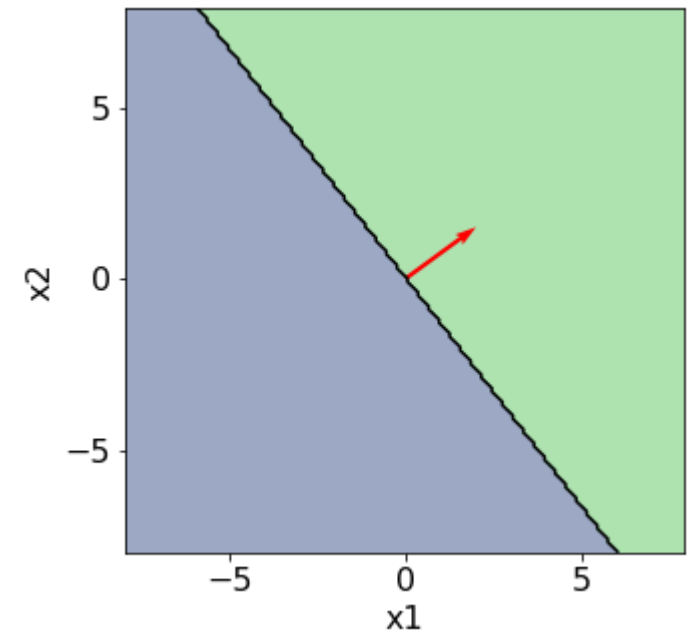
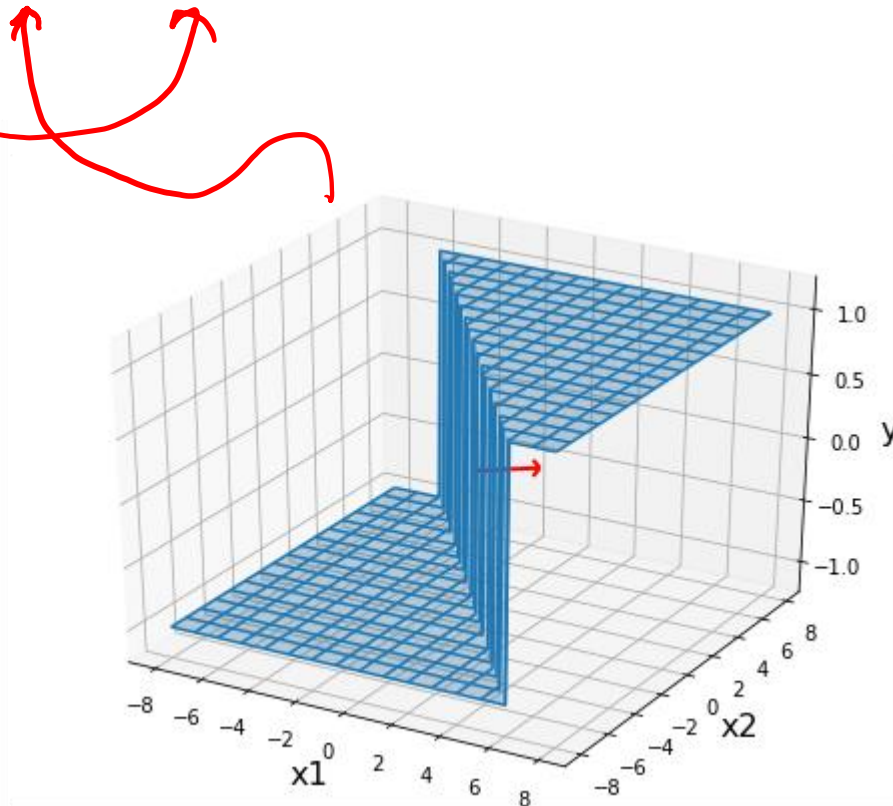
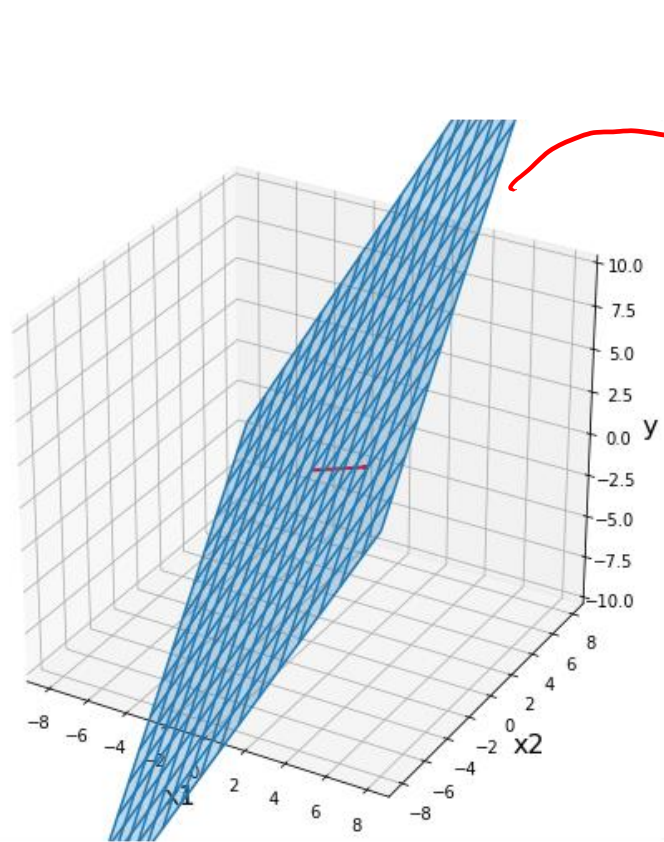
- Neural Networks
  - Building blocks
  - Optimization
    - Composite functions and chain rule
    - Forward-backward passes
    - Matrix calculus

# Perceptron

Classification: Hard threshold on linear model

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

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



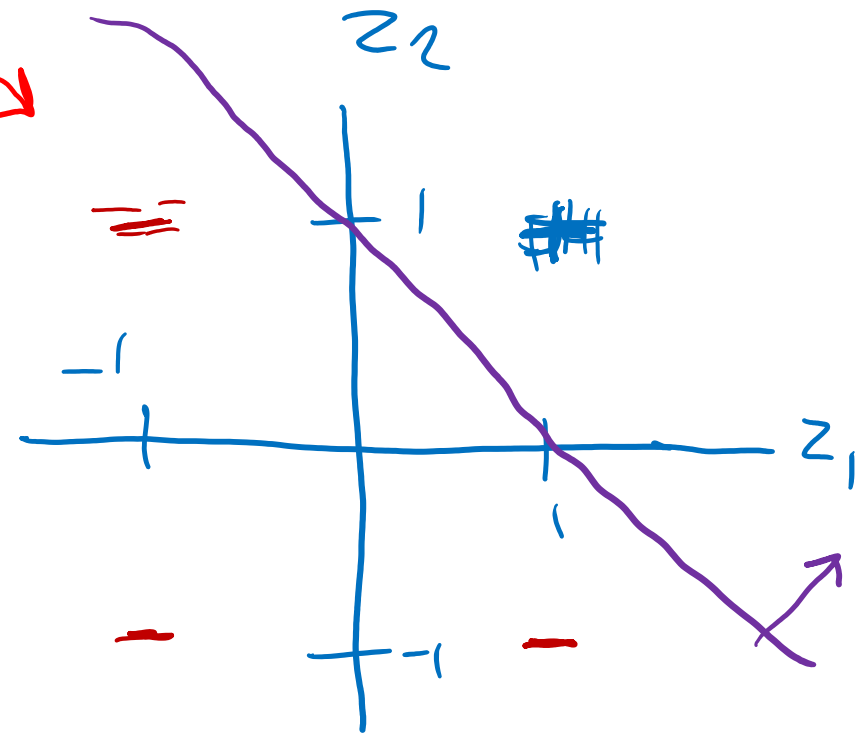
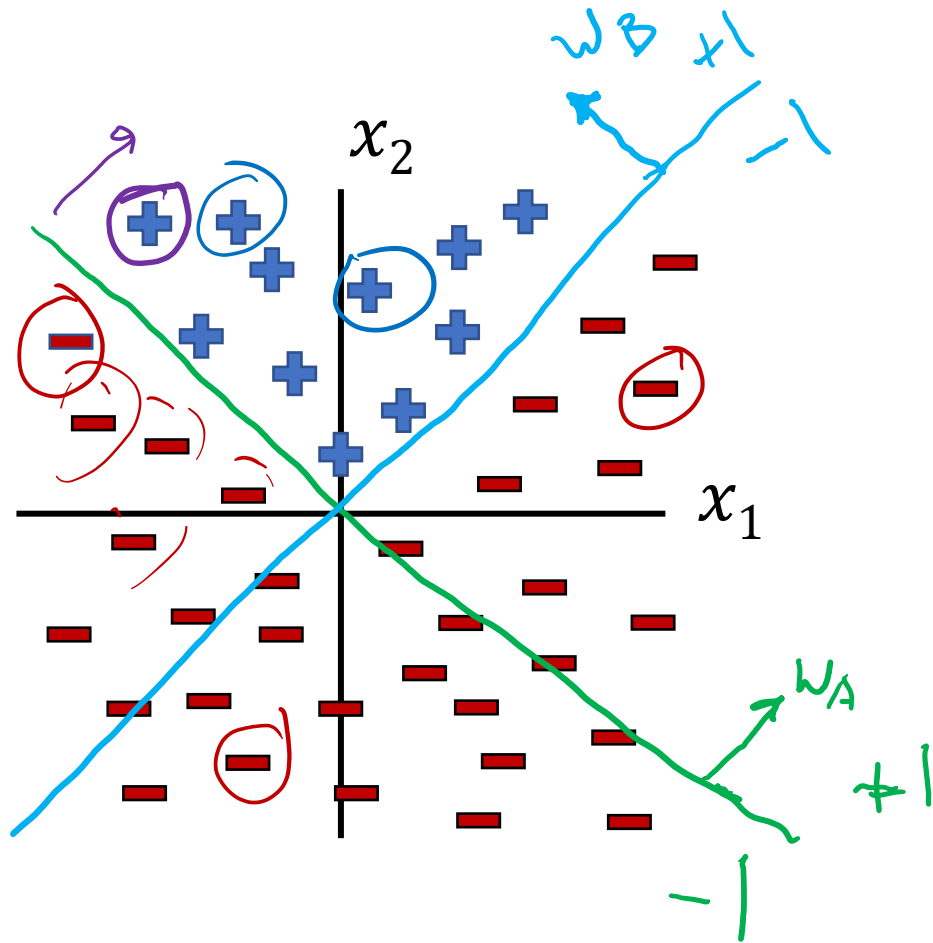
# Classification Design Challenge

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

$$z_1 = h_A(\mathbf{x}) = \text{sign}(\mathbf{w}_A^T \mathbf{x} + b_A)$$

$$z_2 = h_B(\mathbf{x}) = \text{sign}(\mathbf{w}_B^T \mathbf{x} + b_B)$$

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



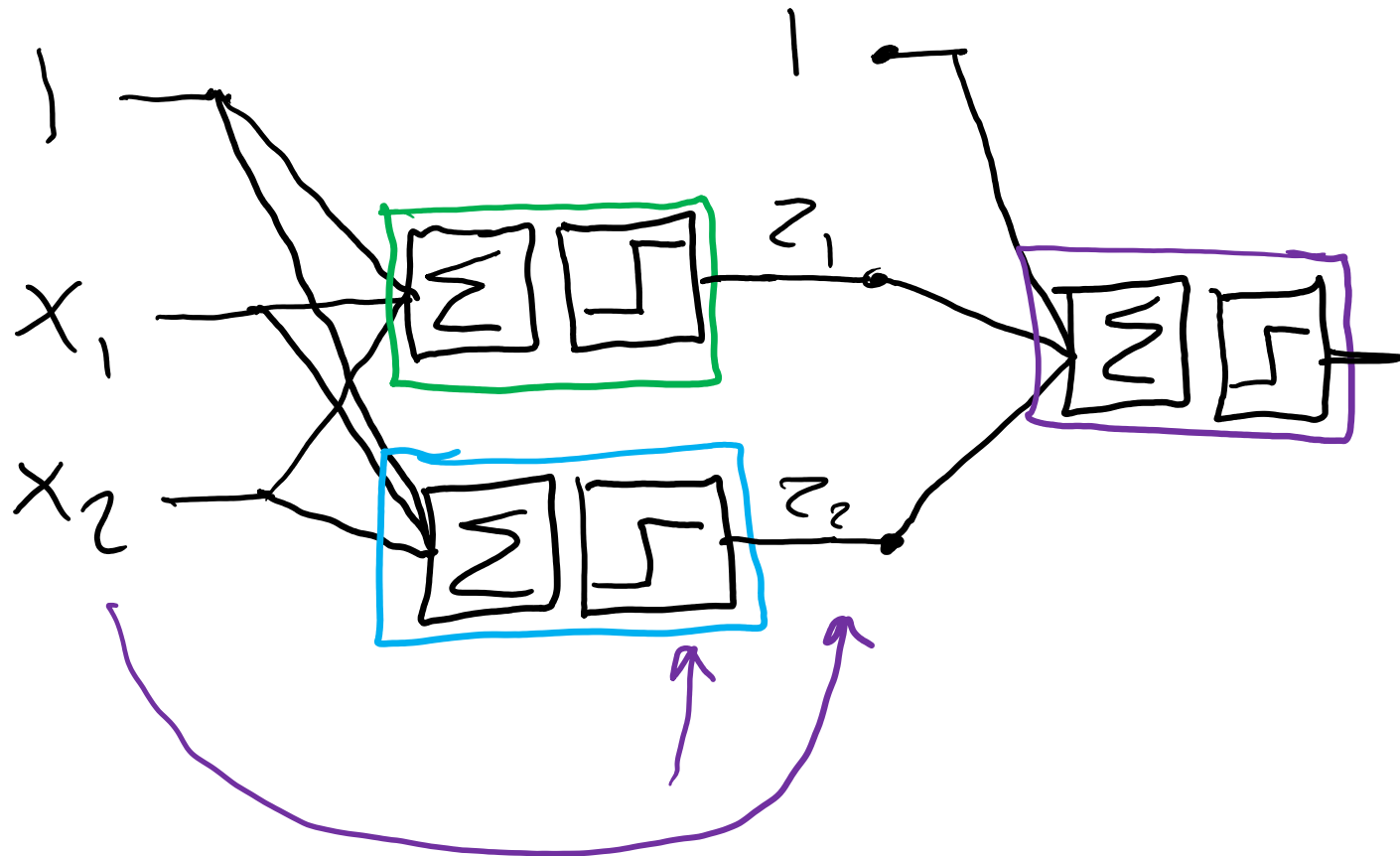
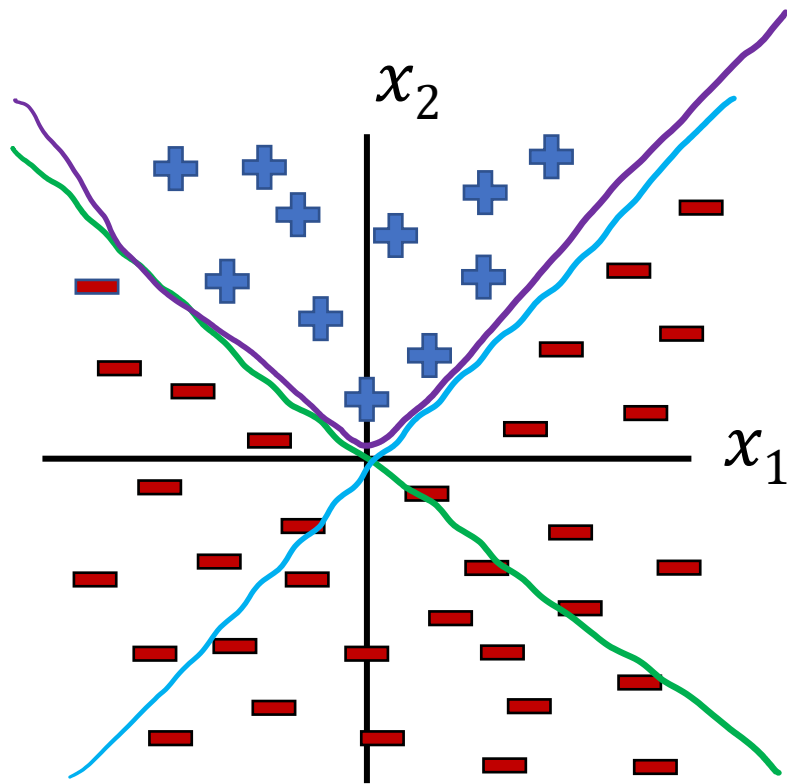
# Classification Design Challenge

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

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

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

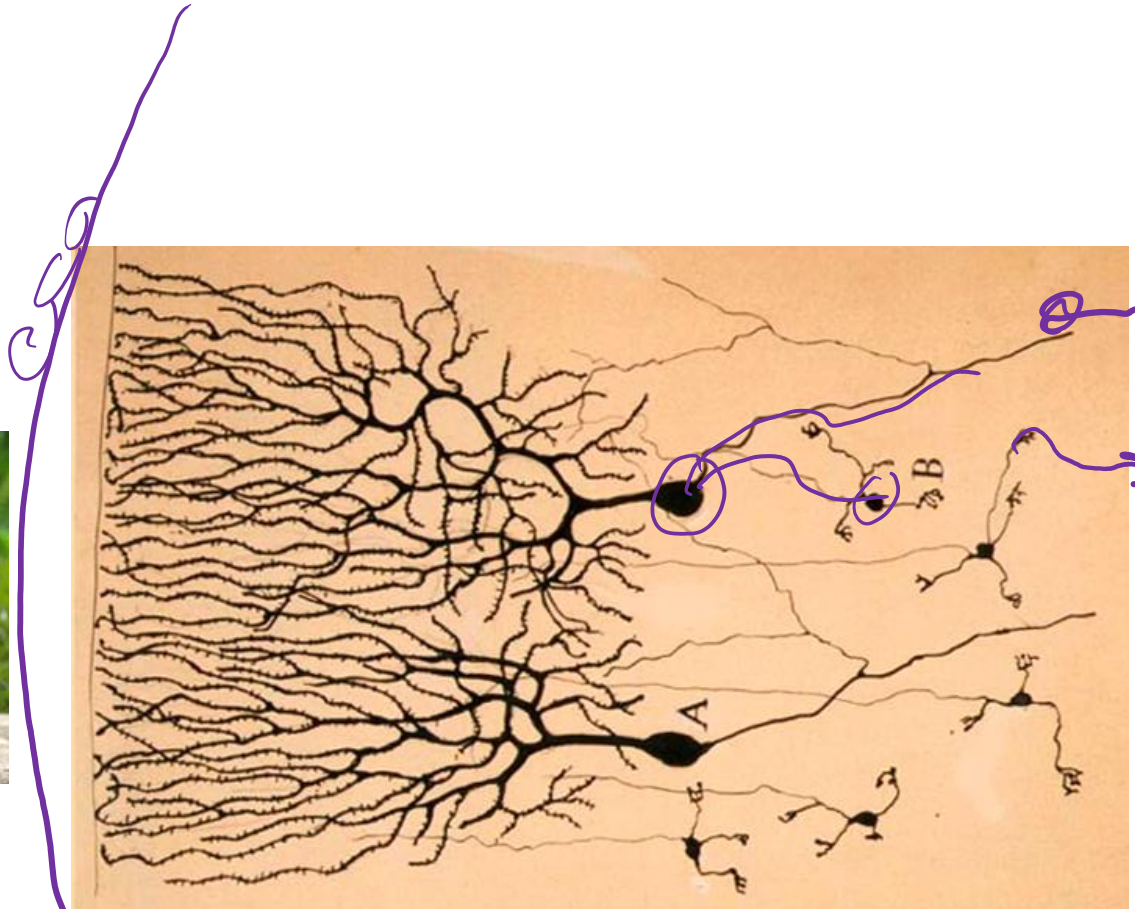
$$h_C(\mathbf{z}) = \text{sign}(\mathbf{w}_C^T \mathbf{z} + b_C)$$



# Neural Networks

Inspired by actual human brain

Input  
Signal



Output  
Signal

DOG



**CAT**



TREE



CAR

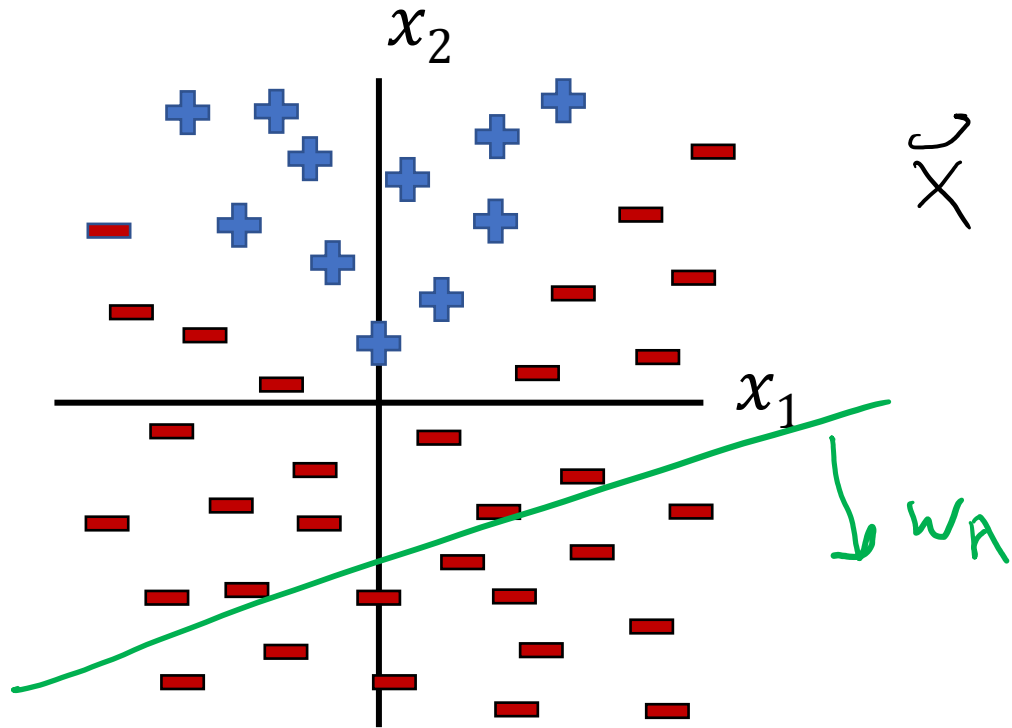


SKY

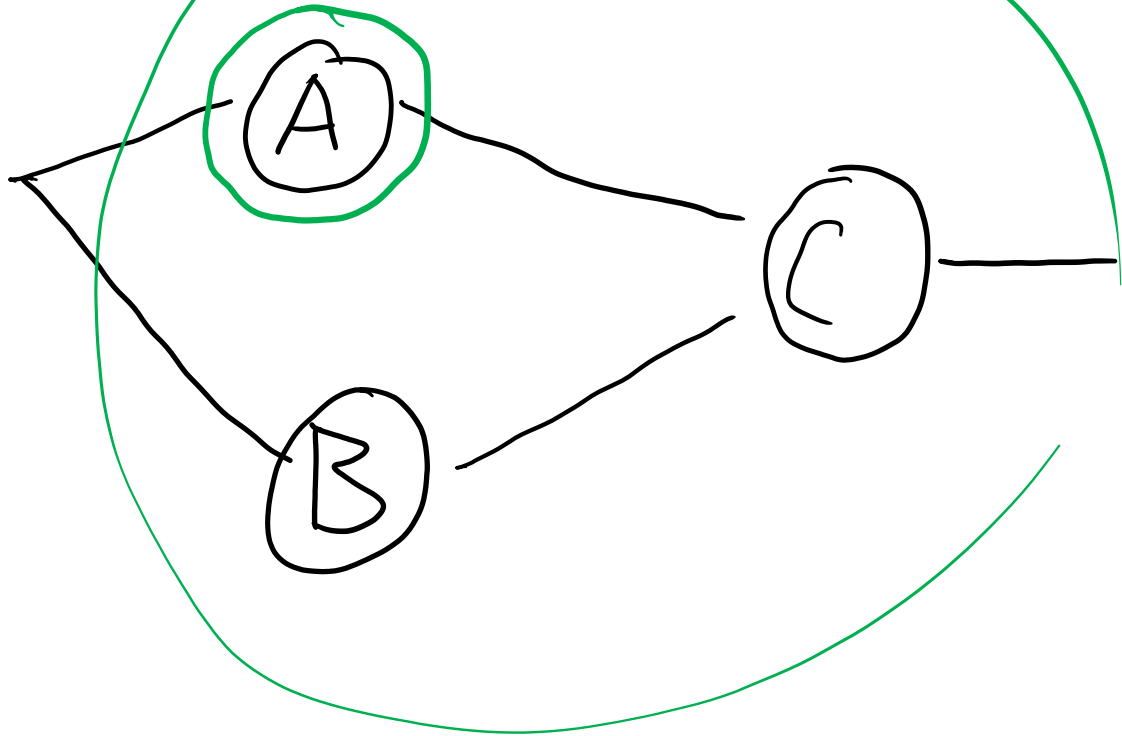
# Classification Design Challenge

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

$$\begin{aligned} h_A(\mathbf{x}) &= \text{sign}(w_A^T \mathbf{x} + b_A) \\ h_B(\mathbf{x}) &= \text{sign}(w_B^T \mathbf{x} + b_B) \\ h_C(\mathbf{x}) &= \text{sign}(w_C^T \mathbf{x} + b_C) \end{aligned}$$



$\mathbf{x}$

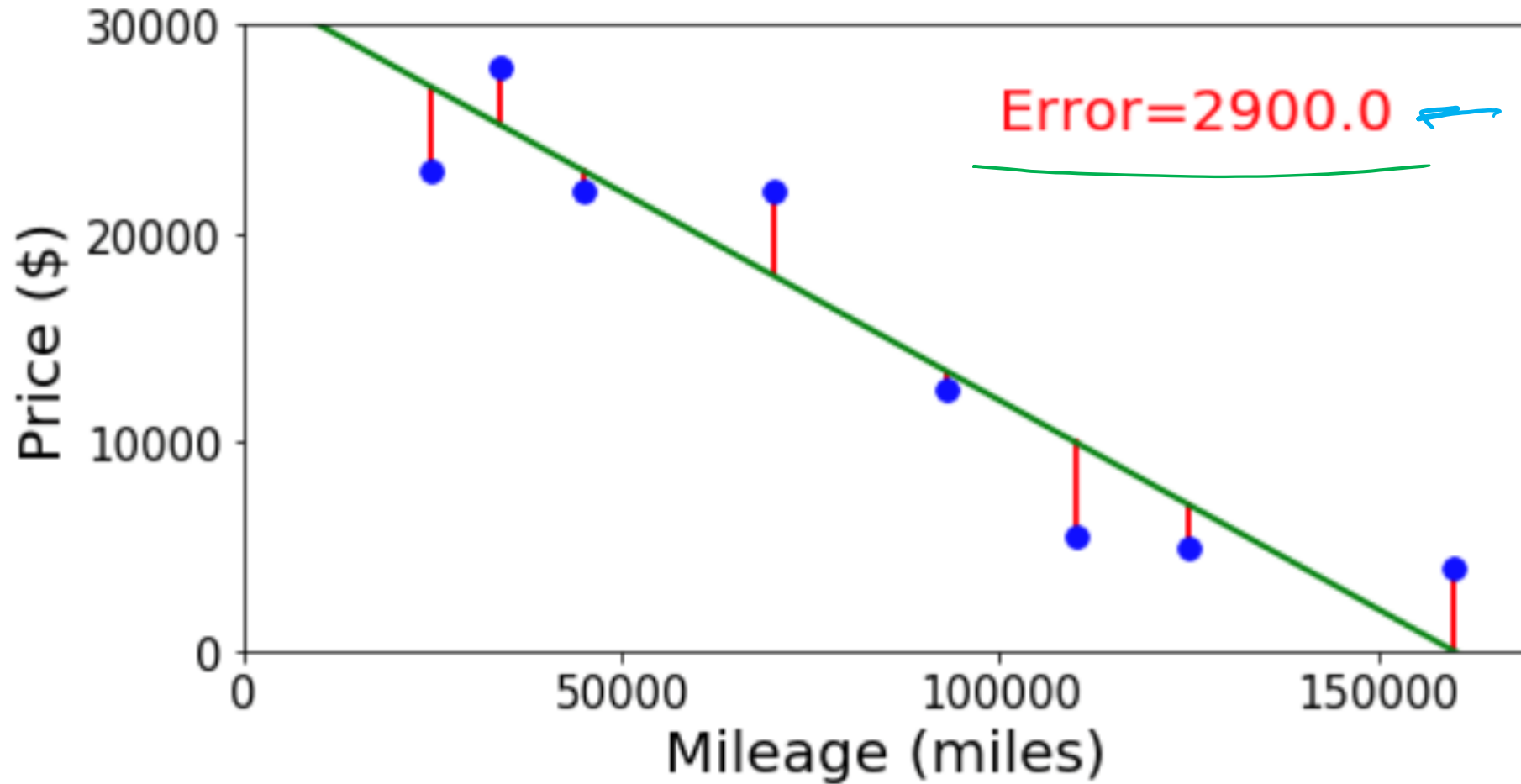
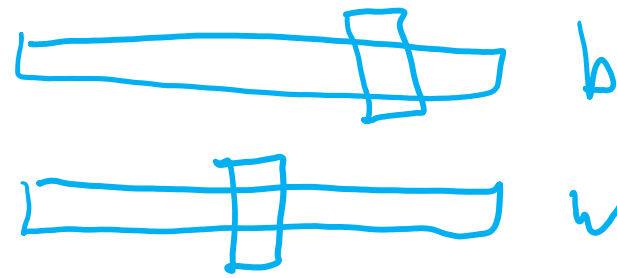




# Neural Networks

Simple single neuron example:

- Selling my car



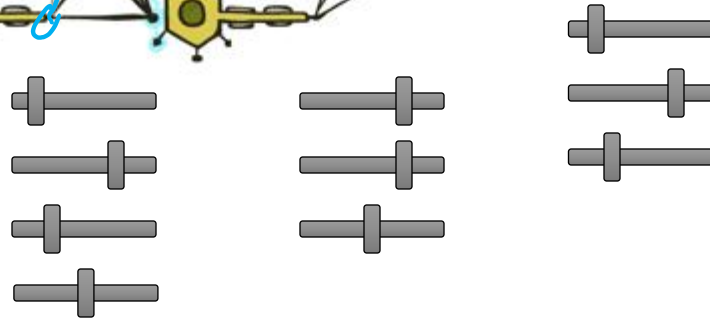
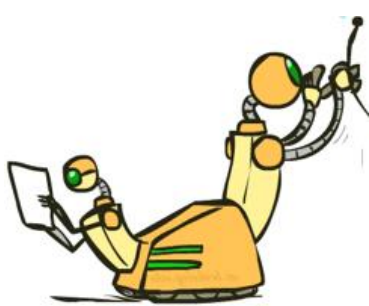
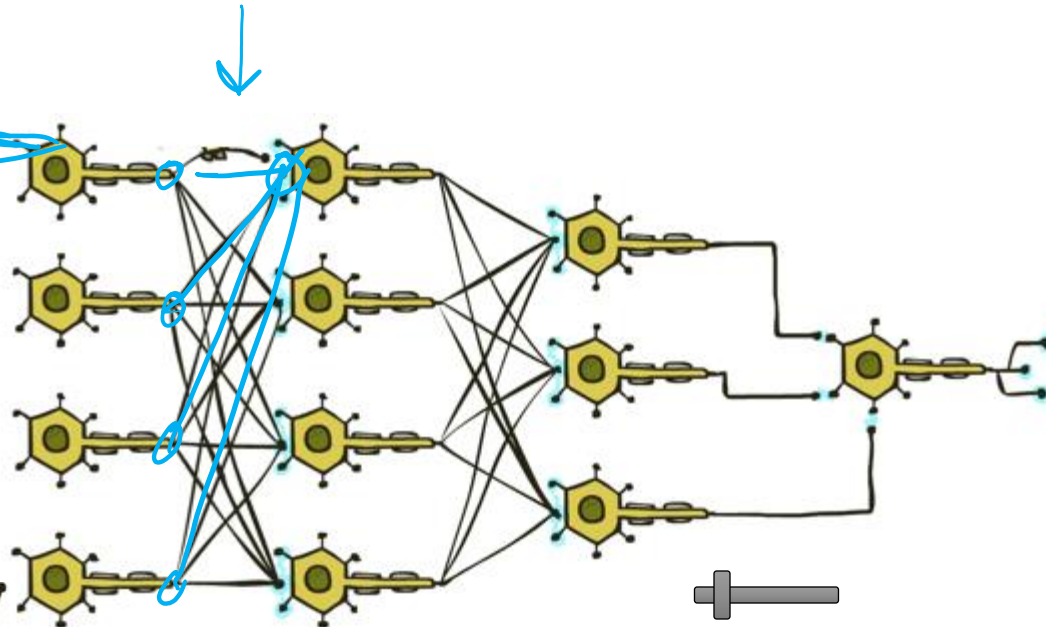
# Neural Networks

Many layers of neurons, millions of parameters

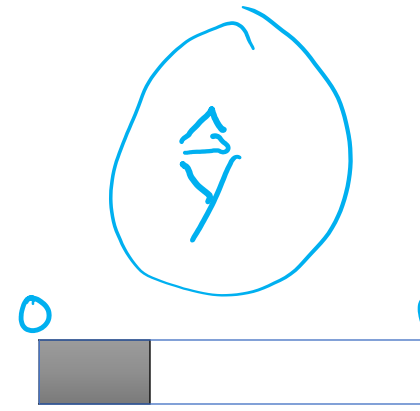
Input  
Signal



Handwritten blue scribbles and lines pointing from the kitten image to the input layer of the neural network.



one-hot  
vector



Output  
Signal

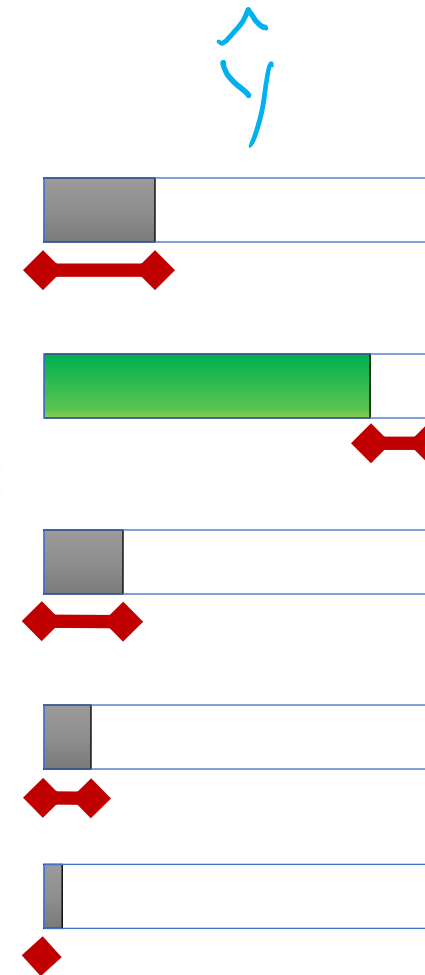
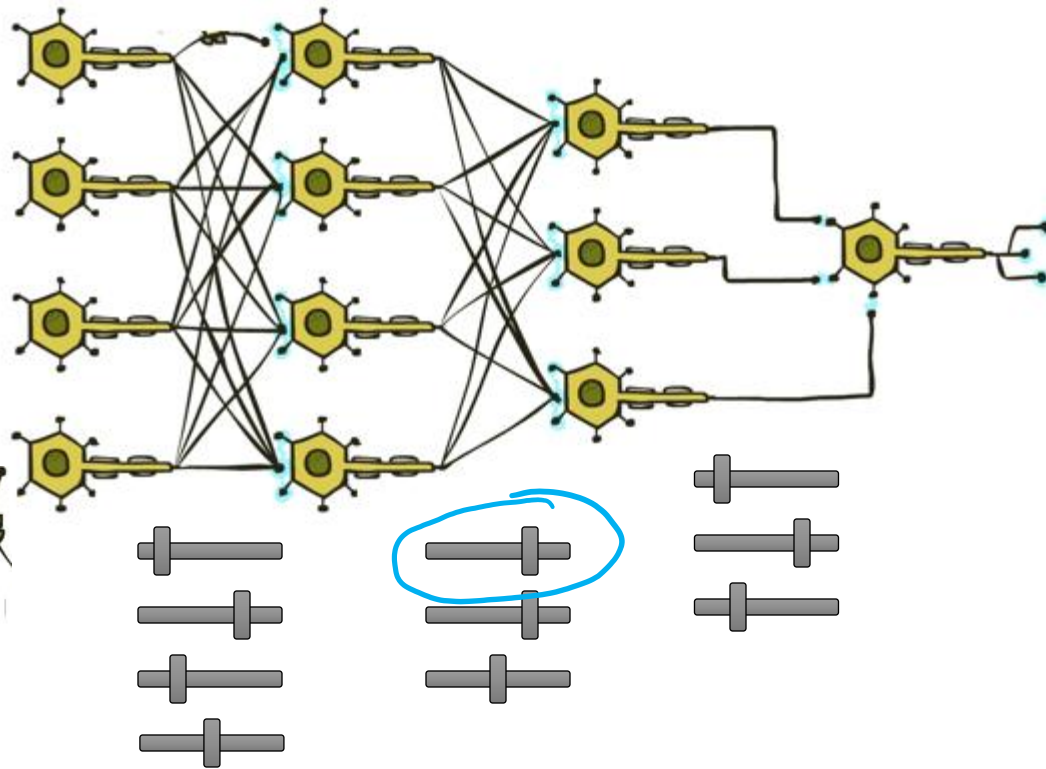
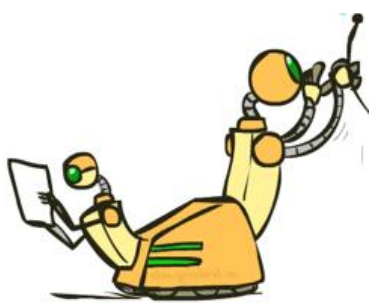


- DOG 0
- CAT** 1
- TREE 0
- CAR 0
- SKY 0

# Neural Networks

Many layers of neurons, millions of parameters

Input  
Signal



Output

Signal

DOG

CAT

TREE

CAR

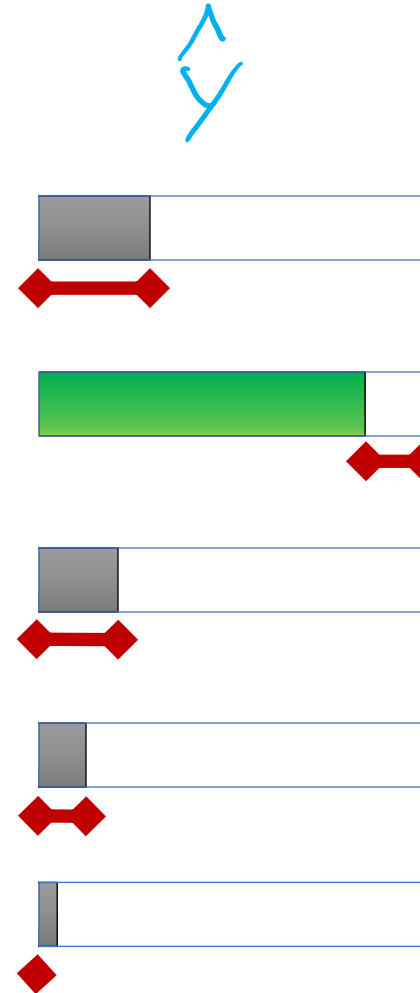
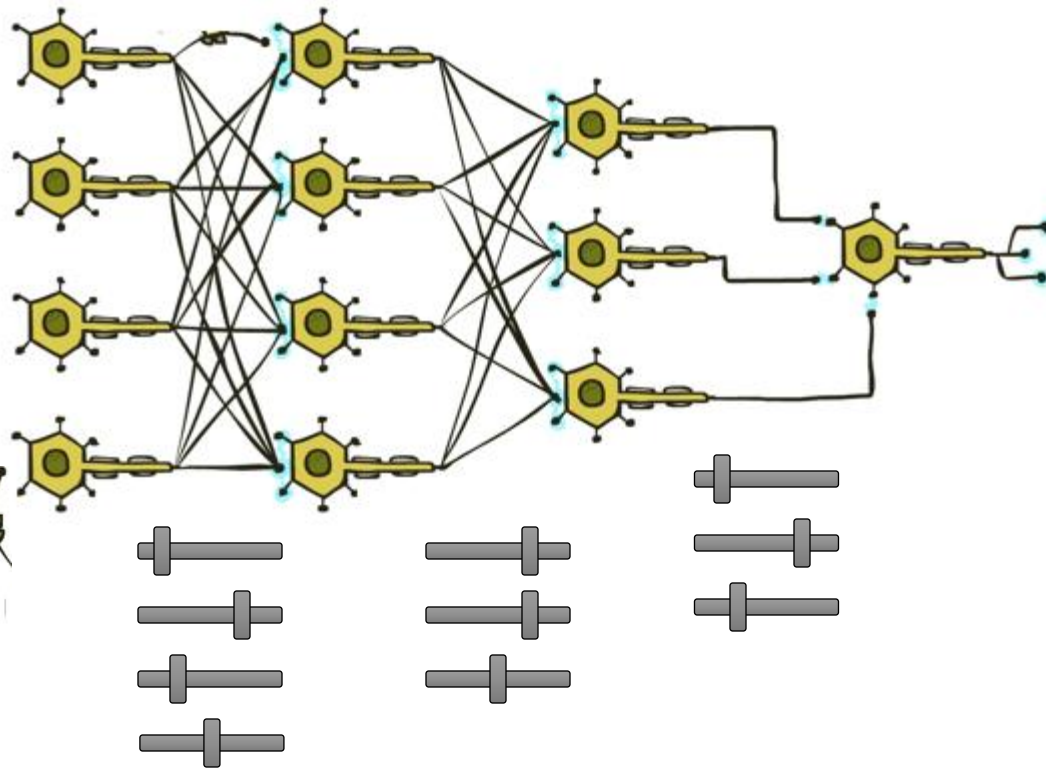
SKY

y

# Neural Networks

Many layers of neurons, millions of parameters

Input  
Signal



Output

Signal

LEFT

**RIGHT**

UP

DOWN

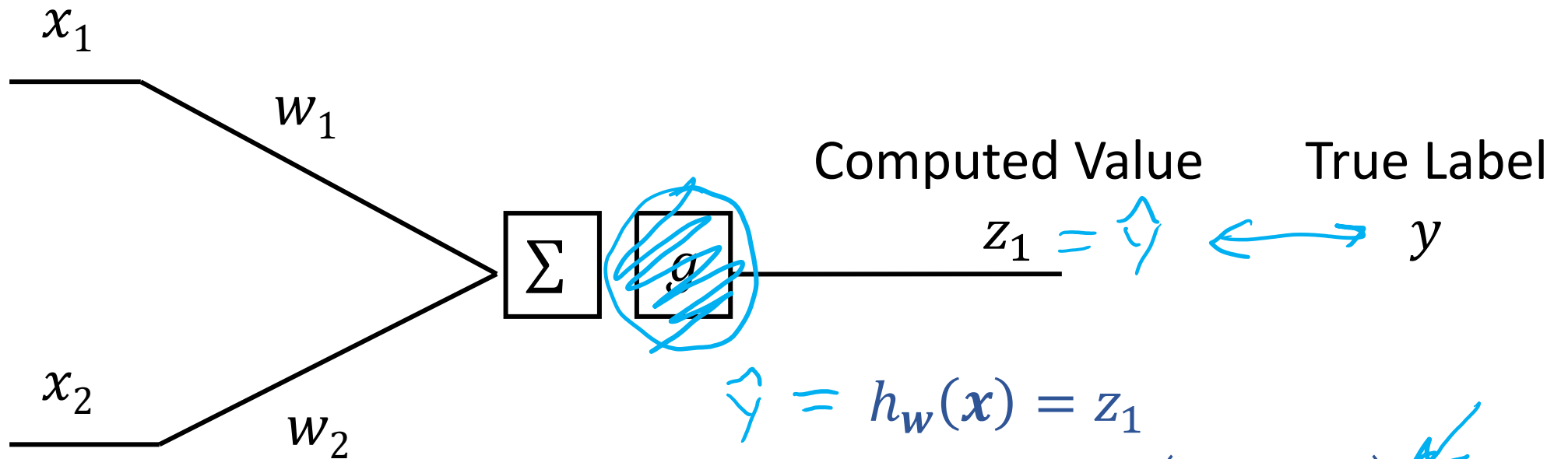
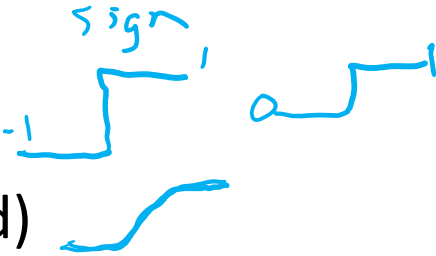
BUTTON



# Single Neuron

## Single neuron system

- Perceptron (if  $g$  is step function)
- Logistic regression (if  $g$  is sigmoid)
- Linear regression (if  $g$  is nothing)



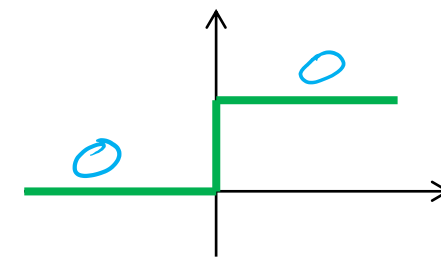
$$\hat{y} = h_w(\mathbf{x}) = z_1$$

$$= h_w(\mathbf{x}) = g\left(\sum_i w_i x_i\right)$$

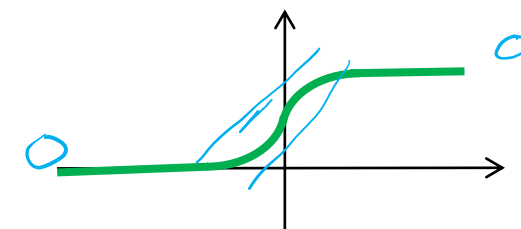
# Activation Functions

It would be really helpful to have a  $g(z)$  that was nicely differentiable

▪ Hard threshold:  $g(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$   $\frac{dg}{dz} = \begin{cases} 0 & z \geq 0 \\ 0 & z < 0 \end{cases}$

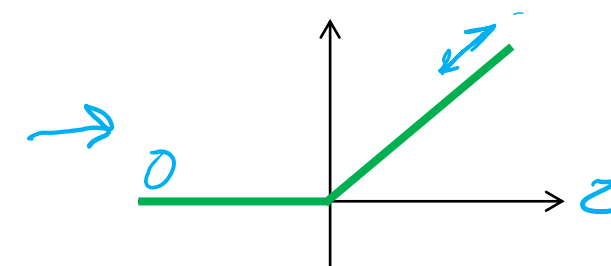


▪ Sigmoid:  $g(z) = \frac{1}{1+e^{-z}}$   $\frac{dg}{dz} = g(z)(1 - g(z))$



▪ (Softmax)

▪ ReLU:  $g(z) = \max(0, z)$   $\frac{dg}{dz} = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$



# Optimizing


How do we find the “best” set of weights?

$$l(y, \hat{y}) = (y - \hat{y})^2$$

$$J^{(i)}(\vec{w}) = l(y^{(i)}, h_w(x^{(i)}))$$

$$\nabla J(\vec{w})$$

$$\vec{w} = \vec{w} - \alpha \nabla_w J(\vec{w})$$

$$\hat{y} = h_w(\mathbf{x}) = g\left(\sum_j \underline{w_j} x_j\right)$$


# Loss Functions

## Regression

- Squared error:  $\ell(y, \hat{y}) = (y - \hat{y})^2$

$$J^{(i)}(\vec{w}) = (y^{(i)} - \hat{y}^{(i)})^2$$

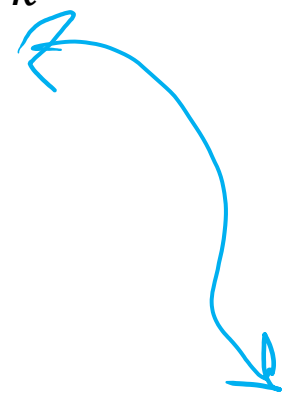
## Classification

- Cross entropy:  $\ell(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_k y_k \log \hat{y}_k$

$K=5$   
label<sup>(9)</sup> = 3

$$\mathbf{y}^{(9)} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\hat{\mathbf{y}}^{(9)} = \begin{bmatrix} .05 \\ .10 \\ .70 \\ 0 \\ .15 \end{bmatrix}$$



$$\prod_{k=1}^K \phi_k^{y_k}$$

*(Note: The original image has a red 'y\_k' above the exponent and a red line through the exponent in the original image, which has been corrected here.)*

$$-\log \prod_{k=1}^K \phi_k^{y_k}$$

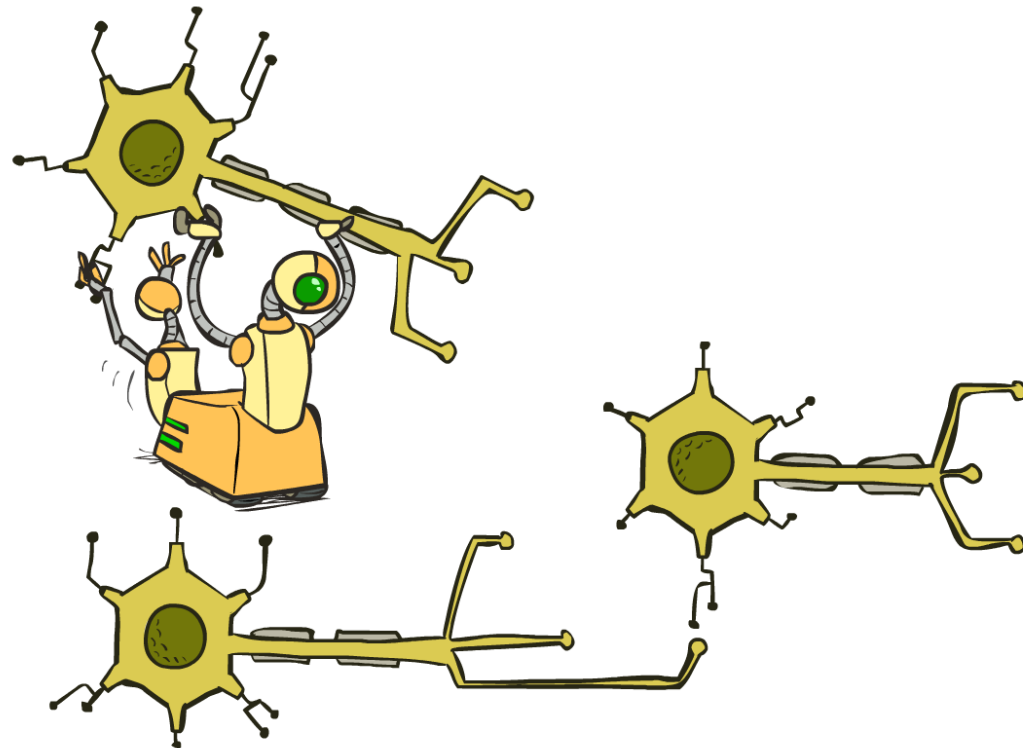
$$-\sum_{k=1}^K y_k \log \phi_k$$



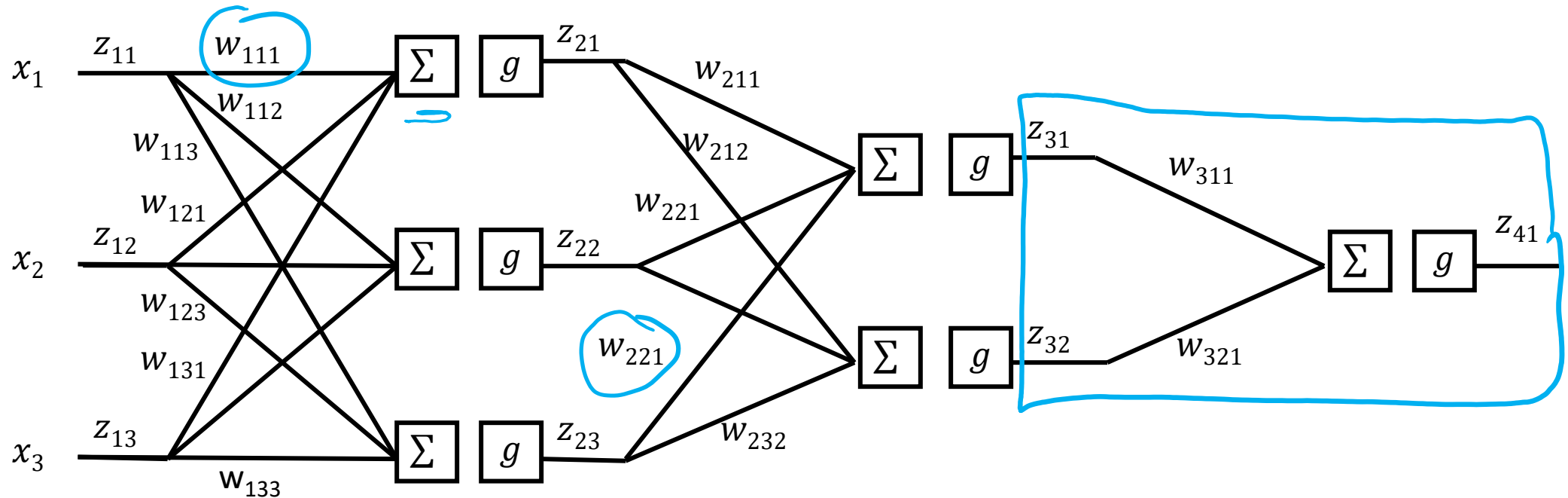
# Multilayer Perceptrons

A **multilayer perceptron** is a feedforward neural network with at least one **hidden layer** (nodes that are neither inputs nor outputs)

MLPs with enough hidden nodes can represent any function



# Neural Network Equations



$$h_w(\mathbf{x}) = z_{4,1}$$

$$z_{1,1} = x_1$$

$$\Rightarrow z_{4,1} = g\left(\sum_i w_{3,i,1} z_{3,i}\right)$$

$$z_{3,1} = g\left(\sum_i w_{2,i,1} z_{2,i}\right)$$

$$z_{d,1} = g\left(\sum_i w_{d-1,i,1} z_{d-1,i}\right)$$

$$h_w(x) = g\left(\sum_k w_{3,k,1} g\left(\sum_j w_{2,j,k} g\left(\sum_i w_{1,i,j} x_i\right)\right)\right)$$

# Optimizing

How do we find the “best” set of weights?

$$l(y, \hat{y}) = (y - \hat{y})^2$$

$$J^{(i)}(\vec{w}) = l(y^{(i)}, h_{\vec{w}}(\vec{x}^{(i)}))$$

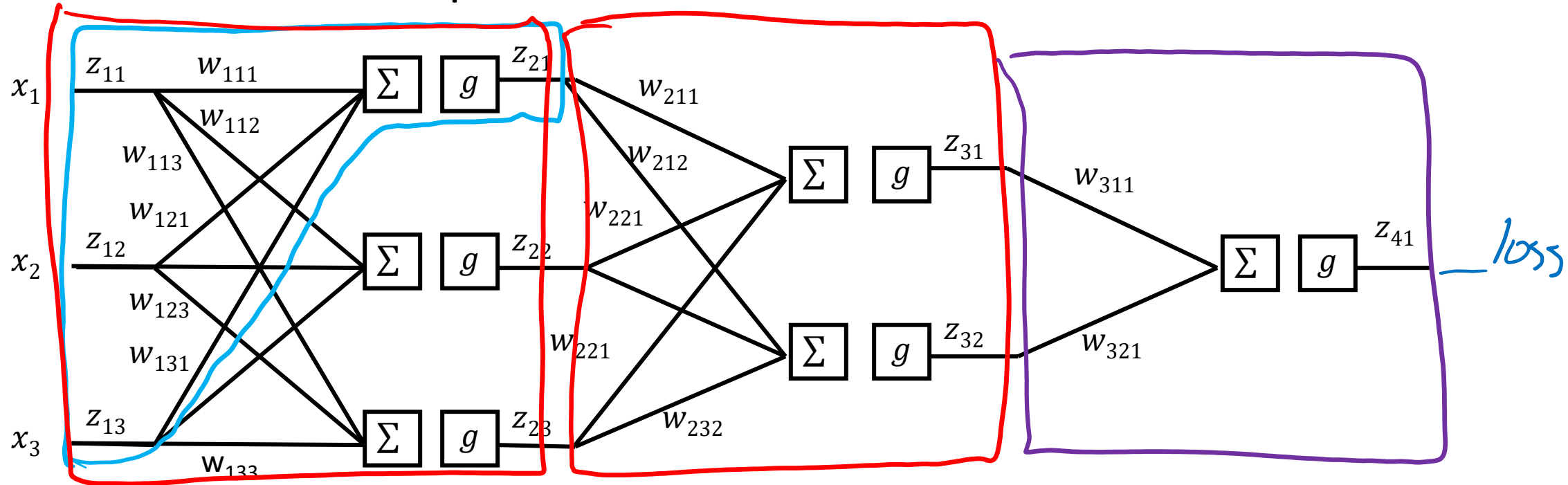
$$\nabla_{\vec{w}} J$$

$$\vec{w} = \vec{w} - \alpha \nabla_{\vec{w}} J(\vec{w})$$

$$h_{\vec{w}}(x) = g \left( \sum_k \underline{w_{3,k,1}} g \left( \sum_j \underline{w_{2,j,k}} g \left( \sum_i \underline{w_{1,i,j}} x_i \right) \right) \right)$$

$$\nabla J = \begin{bmatrix} \frac{\partial J}{\partial w_1} \\ \frac{\partial J}{\partial w_c} \\ \vdots \end{bmatrix}$$

# Neural Network Equations



How do we describe this network?

neuron (6) "node"

hidden layer

output layer

# Network Optimization Details

# Reminder: Calculus Chain Rule (scalar version)

$$y = f(z)$$

$$z = g(x)$$

$$y = f(g(x))$$

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

$$\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$$

# Network Optimization

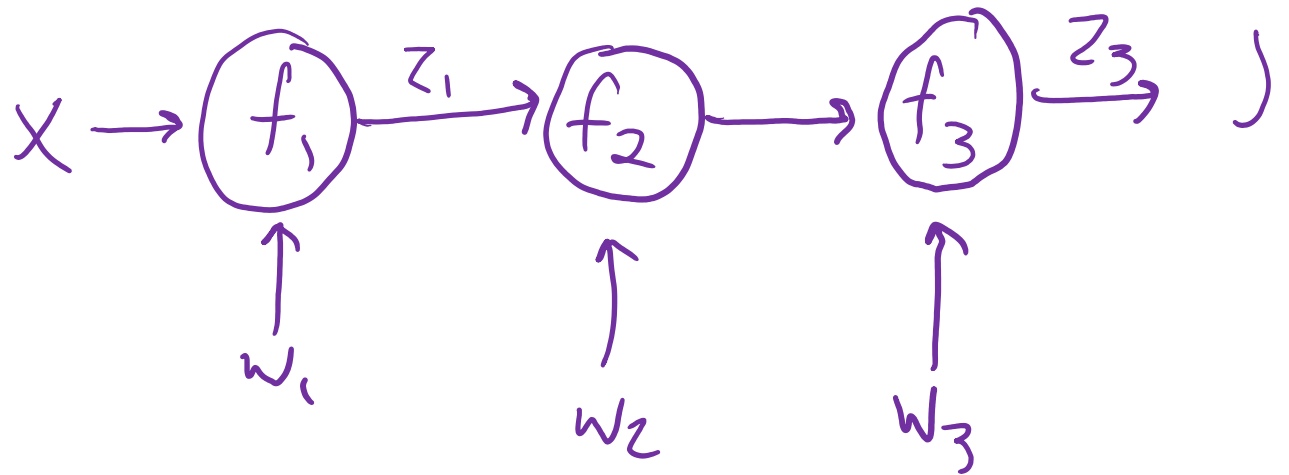
$$J(\mathbf{w}) = z_3$$

$$z_3 = f_3(w_3, z_2)$$

$$z_2 = f_2(w_2, z_1)$$

$$z_1 = f_1(w_1, x)$$

$$J(w_1, w_2, w_3) = f_3(w_3, f_2(w_2, f_1(w_1, x)))$$



$$\frac{dJ}{dw_3} = \frac{dJ}{dz_3} \frac{dz_3}{dw_3}$$

$$\frac{dJ}{dw_2} = \frac{dJ}{dz_3} \frac{dz_3}{dz_2} \frac{dz_2}{dw_2}$$

$$\frac{dJ}{dw_1} = \frac{dJ}{dz_3} \frac{dz_3}{dz_2} \frac{dz_2}{dz_1} \frac{dz_1}{dw_1}$$

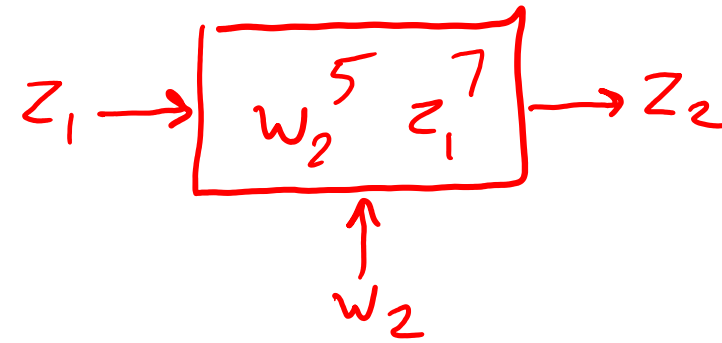
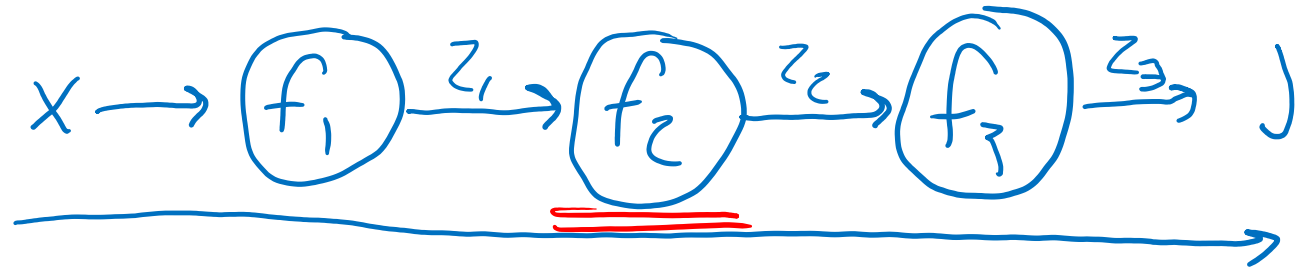
# Network Optimization: Forward then Backwards

$$J(\mathbf{w}) = z_3$$

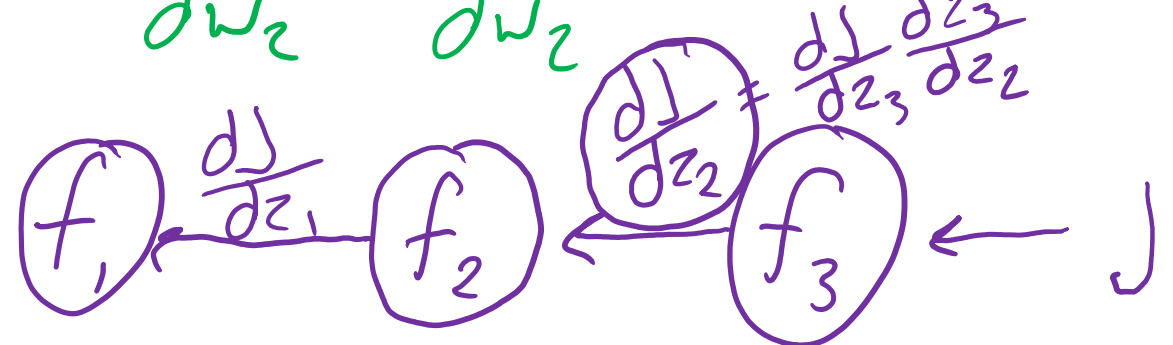
$$z_3 = f_3(w_3, z_2)$$

$$z_2 = f_2(w_2, z_1) = w_2^5 z_1^7$$

$$z_1 = f_1(w_1, x)$$



$$\frac{dz_2}{dw_2} = \frac{df_2}{dw_2} = 5w^4 z_1^7$$



$$\frac{\partial J}{\partial w_3} = \frac{\partial J}{\partial z_3} \frac{\partial z_3}{\partial w_3}$$

$$\frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial z_3} \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_2}$$

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial z_3} \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

Lots of repeated calculations



# Network Optimization: Layer Implementation

$$J(\mathbf{w}) = z_3$$

$$z_3 = f_3(w_3, z_2)$$

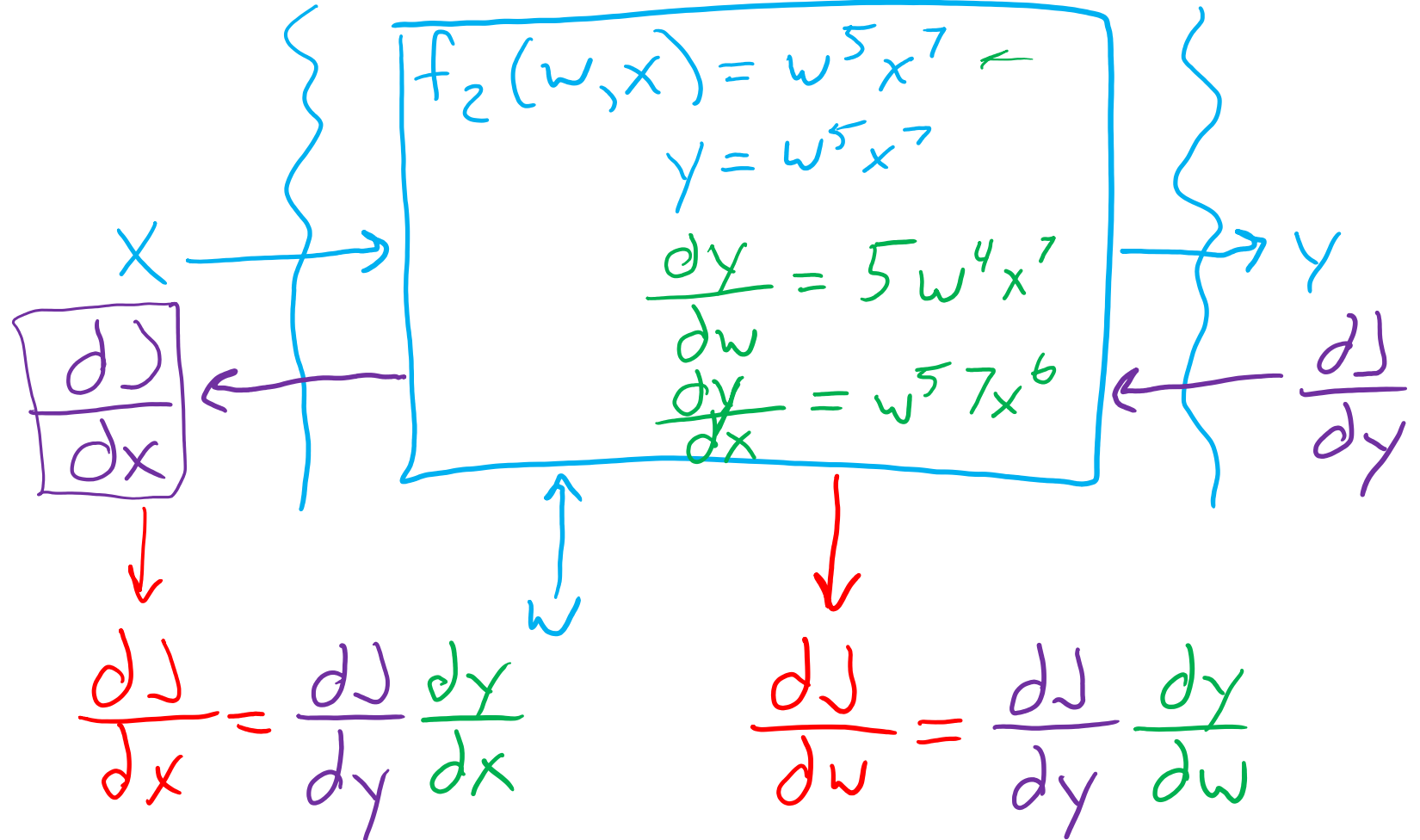
$$z_2 = f_2(w_2, z_1)$$

$$z_1 = f_1(w_1, x)$$

$$\frac{\partial J}{\partial w_3} = \frac{\partial J}{\partial z_3} \frac{\partial z_3}{\partial w_3}$$

$$\frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial z_3} \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial w_2}$$

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial z_3} \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$



Lots of repeated calculations

# Backpropagation (so-far)

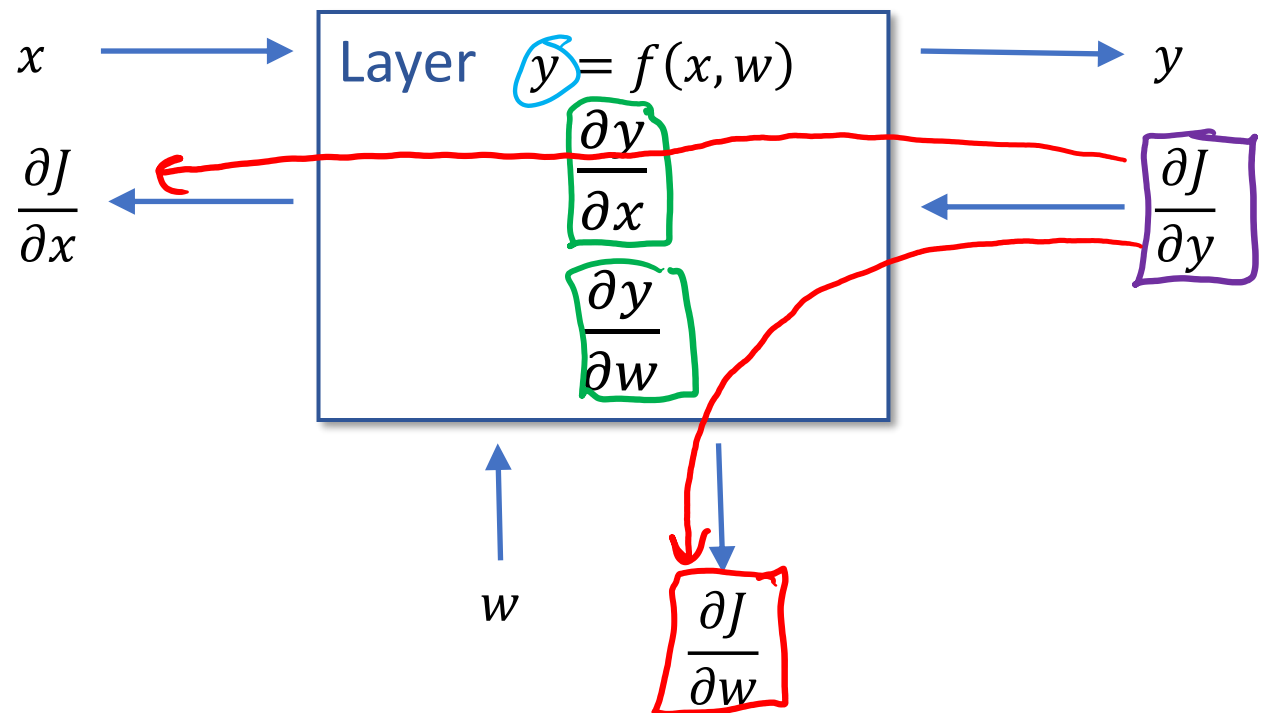
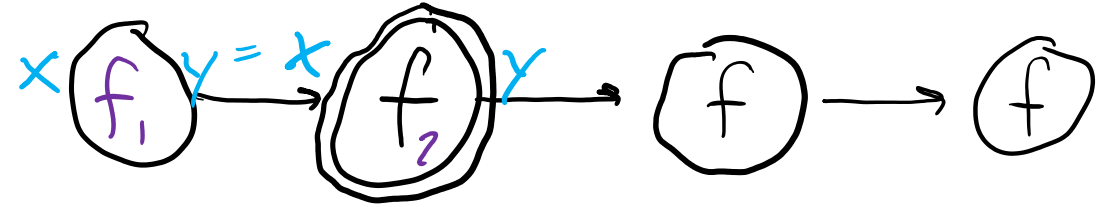
Compute derivatives per layer, utilizing previous derivatives

Objective:  $J(\mathbf{w})$

Arbitrary layer:  $y = f(x, w)$

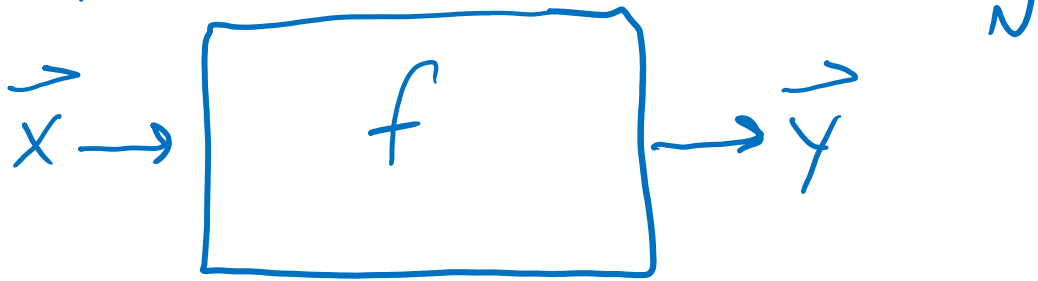
Need:

$$\begin{aligned} \frac{\partial J}{\partial x} &= \frac{\partial J}{\partial y} \frac{\partial y}{\partial x} \\ \frac{\partial J}{\partial w} &= \frac{\partial J}{\partial y} \frac{\partial y}{\partial w} \end{aligned}$$

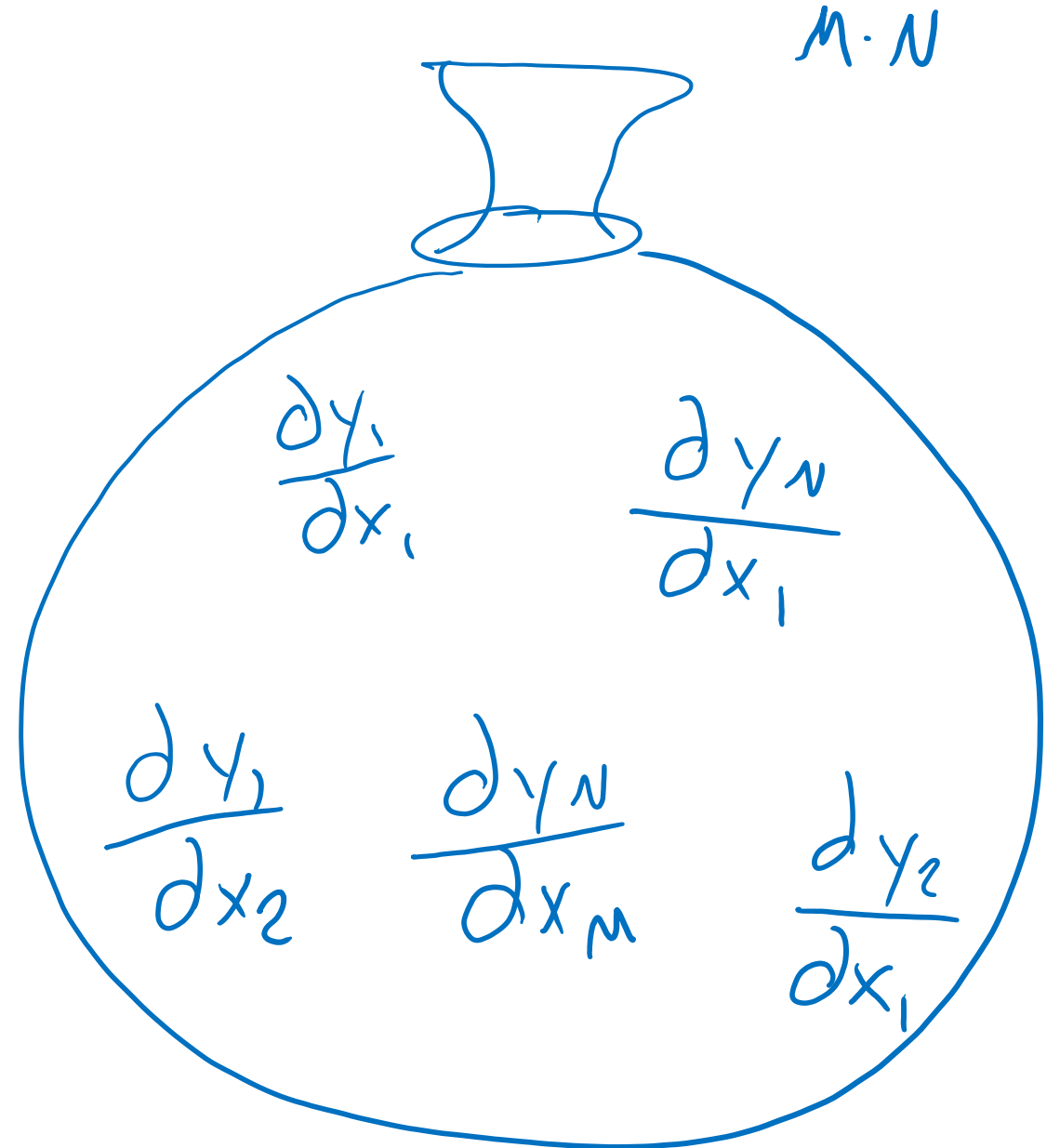


# Matrix Calculus

One way to think of it: Bag of Derivatives

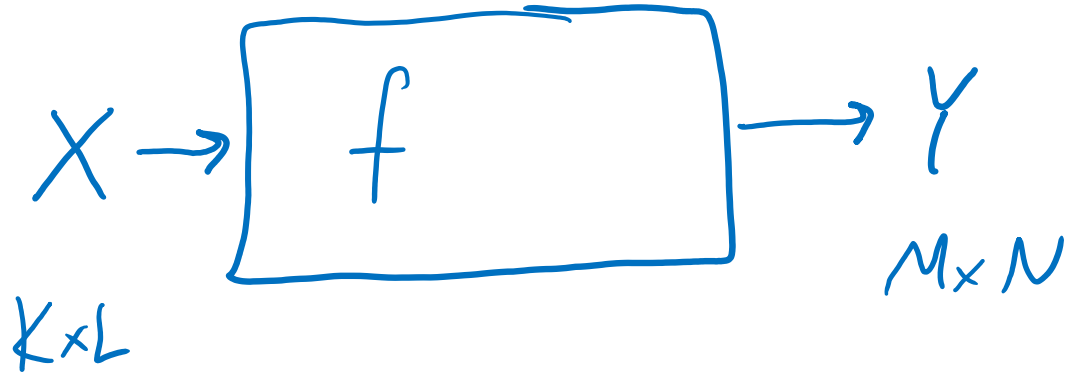


$$\frac{\partial y_i}{\partial x_j}$$

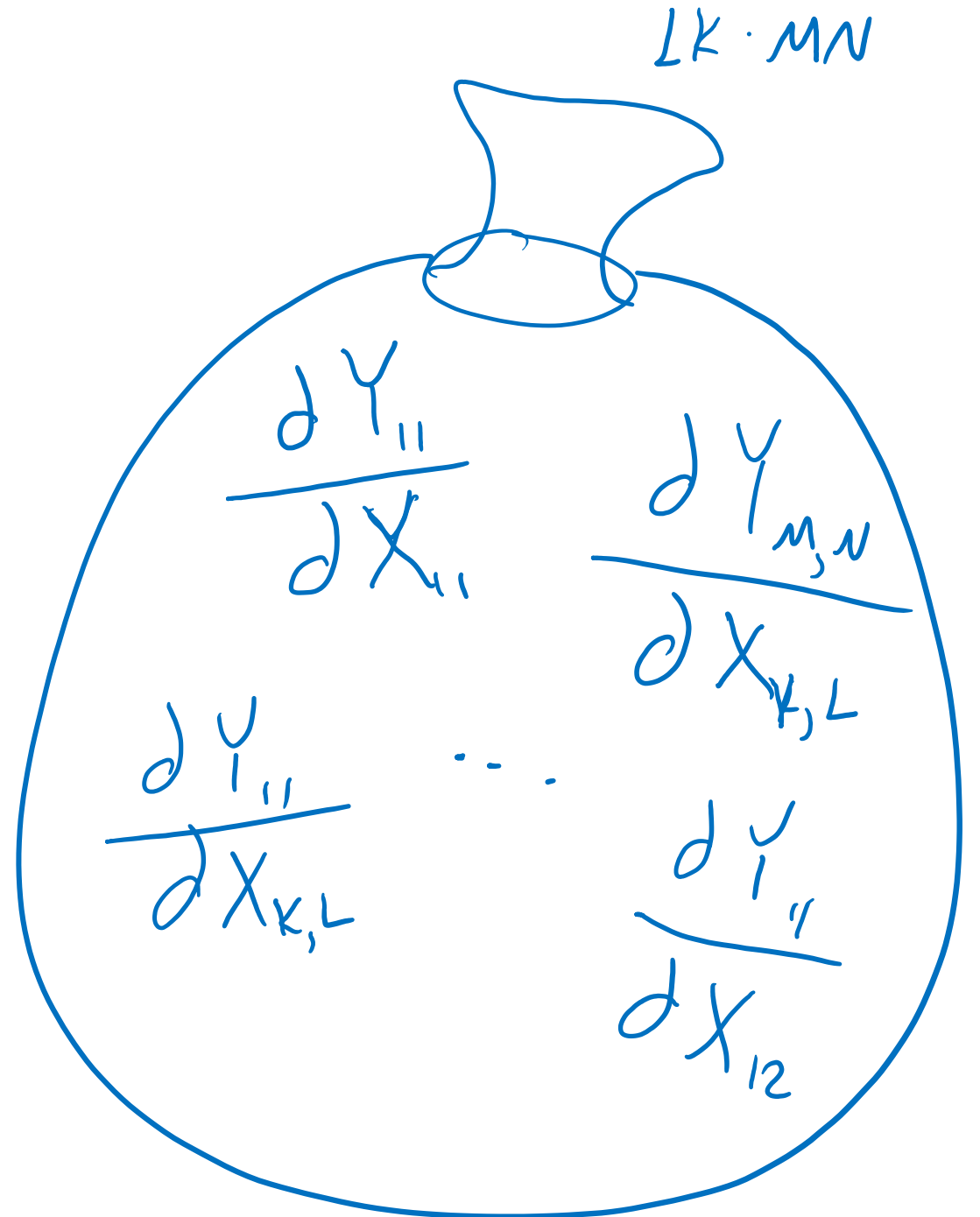


# Matrix Calculus

One way to think of it: Bag of Derivatives



$$\frac{dY_{i,j}}{dX_{k,l}}$$



# Matrix Calculus

Jacobian: Vector in, vector out

→ Numerator-layout

$$\mathbf{y} = f(\mathbf{x}) \quad \mathbf{y} \in \mathbb{R}^N, \quad \mathbf{x} \in \mathbb{R}^M, \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \in \mathbb{R}^{N \times M}$$

$$\frac{d\vec{y}}{d\vec{x}} = \begin{matrix} & \begin{matrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_M} \end{matrix} \\ \begin{matrix} N \\ \vdots \\ \frac{\partial y_N}{\partial x_1} \end{matrix} & \begin{matrix} \frac{\partial y_N}{\partial x_M} \end{matrix} \end{matrix}$$

# Matrix Calculus

Vector in, scalar out

Numerator-layout

$$y = f(\mathbf{x}) \quad y \in \mathbb{R}, \quad \mathbf{x} \in \mathbb{R}^M, \quad \frac{\partial y}{\partial \mathbf{x}} \in \mathbb{R}^{1 \times M}$$

$$\frac{\partial y}{\partial \vec{x}} = \left[ \frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_M} \right]$$

# Matrix Calculus

Scalar in, vector out

Numerator-layout

$$\mathbf{y} = f(x) \quad \mathbf{y} \in \mathbb{R}^N, \quad x \in \mathbb{R}, \quad \frac{\partial \mathbf{y}}{\partial x} \in \mathbb{R}^{N \times 1}$$

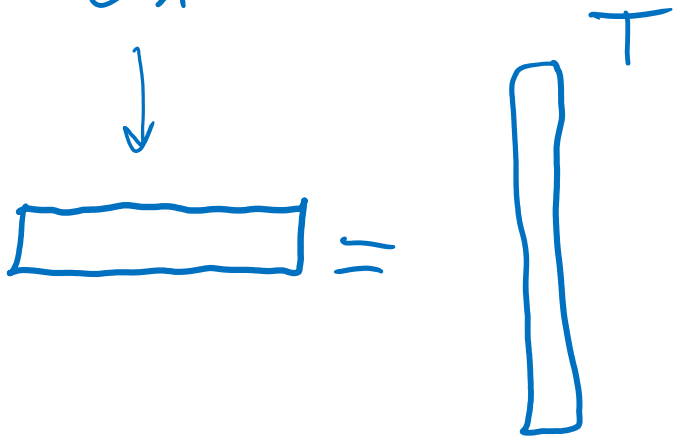
$$\frac{d\vec{y}}{dx} = \begin{bmatrix} \frac{dy_1}{dx} \\ \vdots \\ \frac{dy_N}{dx} \end{bmatrix}$$

# Matrix Calculus

Gradient: Vector in, scalar out

Transpose of numerator-layout

$$y = f(\mathbf{x}) \quad y \in \mathbb{R}, \quad \mathbf{x} \in \mathbb{R}^M, \quad \frac{\partial y}{\partial \mathbf{x}} \in \mathbb{R}^{1 \times M}, \quad \nabla_{\mathbf{x}} f \in \mathbb{R}^{M \times 1}$$

$$\frac{dy}{d\vec{x}} = \nabla_{\mathbf{x}} f^T$$




# Matrix Calculus

Matrix in, scalar out

Keep same dimensions as matrix

$$y = f(\mathbf{X}) \quad y \in \mathbb{R}, \quad \mathbf{X} \in \mathbb{R}^{N \times M}, \quad \frac{\partial y}{\partial \mathbf{X}} \in \mathbb{R}^{N \times M}$$

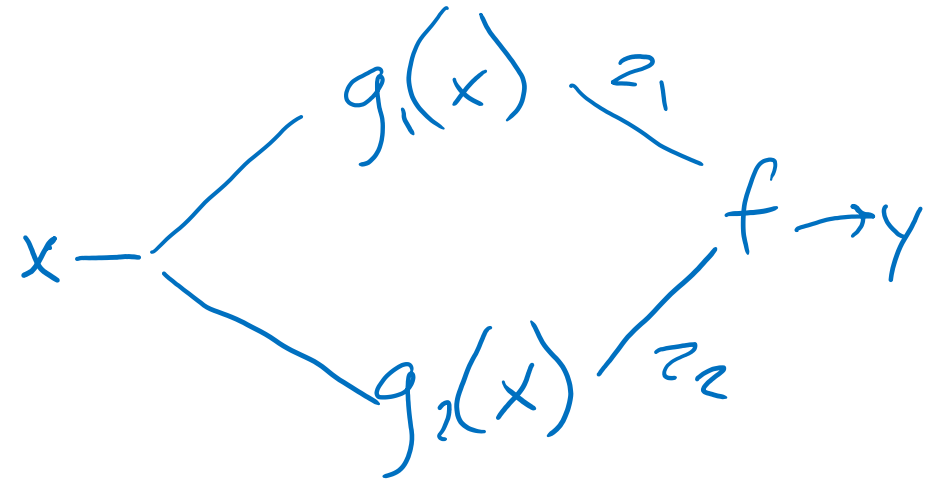
$$\frac{dy}{d\mathbf{X}} = \begin{bmatrix} \frac{dy}{dX_{1,1}} & \dots & \frac{dy}{dX_{1,M}} \\ \vdots & \ddots & \vdots \\ \frac{dy}{dX_{N,1}} & \dots & \frac{dy}{dX_{N,M}} \end{bmatrix}$$

# Multivariate Chain Rule

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

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

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



$$\frac{dy}{dx} = \left( \frac{dy}{dz_1} \right) \frac{dz_1}{dx} + \left( \frac{dy}{dz_2} \right) \frac{dz_2}{dx}$$

$$= (4z_1^3 e^{z_2} + 5 + 0) \cos(x) + (z_1^4 e^{z_2} + 0 + 7) 3x^2$$

# Multivariate Chain Rule

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

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

$$y = f(z_1, z_2) = \underline{z_1 z_2}$$

$$\begin{aligned} \frac{dy}{dx} &= \frac{dy}{dz_1} \frac{dz_1}{dx} + \frac{dy}{dz_2} \frac{dz_2}{dx} \\ &= z_2 \cos(x) + z_1 3x^2 \end{aligned}$$

# Calculus Chain Rule

Scalar:

$$y = f(z)$$

$$z = g(x)$$

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

Multivariate:

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

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

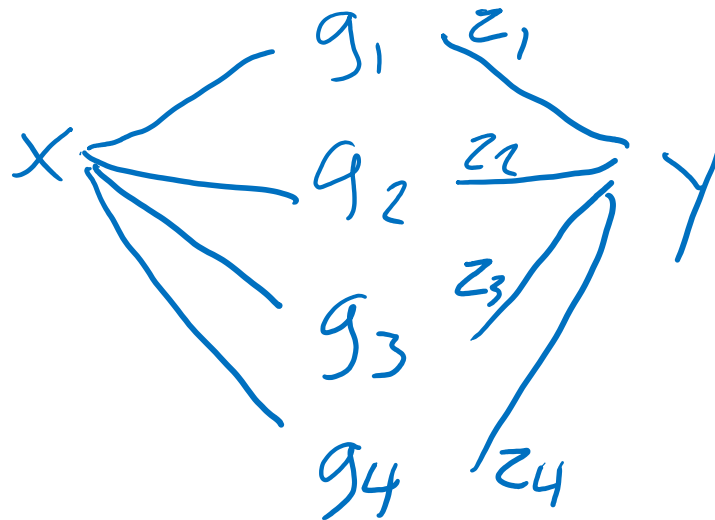
$$\frac{dy}{dx} = \sum_j \frac{\partial y}{\partial z_j} \frac{\partial z_j}{\partial x}$$

Multivariate:

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

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

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



# Network Optimization

$$J(\mathbf{w}) = z_4$$

$$z_4 = f_4(w_D, w_E, z_2, z_3)$$

$$z_3 = f_3(w_C, z_1)$$

$$z_2 = f_2(w_B, z_1)$$

$$z_1 = f_1(w_A, x)$$

Need multivariate chain rule!

# Network Optimization

$$J(\mathbf{w}) = z_4$$

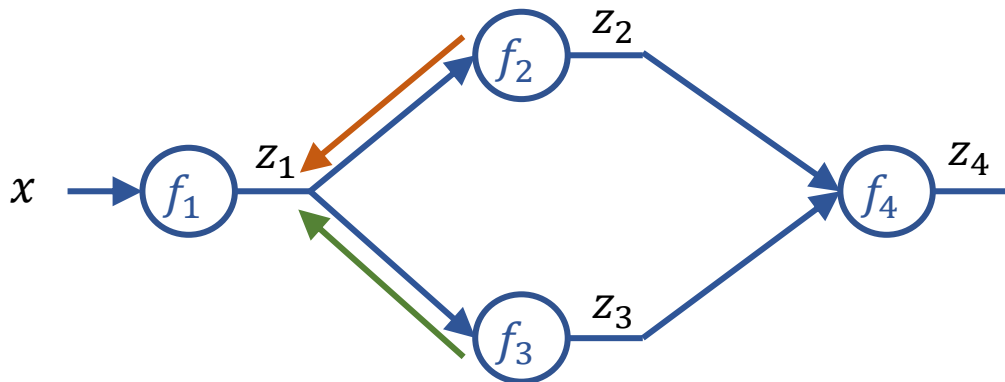
$$z_4 = f_4(w_D, w_E, z_2, z_3)$$

$$z_3 = f_3(w_C, z_1)$$

$$z_2 = f_2(w_B, z_1)$$

$$z_1 = f_1(w_A, x)$$

Need multivariate chain rule!



$$\frac{\partial J}{\partial w_E} = \frac{\partial J}{\partial z_4} \frac{\partial z_4}{\partial w_E}$$

$$\frac{\partial J}{\partial w_D} = \frac{\partial J}{\partial z_4} \frac{\partial z_4}{\partial w_D}$$

$$\frac{\partial J}{\partial z_3} = \frac{\partial J}{\partial z_4} \frac{\partial z_4}{\partial z_3}$$

$$\frac{\partial J}{\partial z_2} = \frac{\partial J}{\partial z_4} \frac{\partial z_4}{\partial z_2}$$

$$\frac{\partial J}{\partial w_C} = \frac{\partial J}{\partial z_3} \frac{\partial z_3}{\partial w_C}$$

$$\frac{\partial J}{\partial w_B} = \frac{\partial J}{\partial z_2} \frac{\partial z_2}{\partial w_B}$$

$$\frac{\partial J}{\partial z_1} = \frac{\partial J}{\partial z_2} \frac{\partial z_2}{\partial z_1} + \frac{\partial J}{\partial z_3} \frac{\partial z_3}{\partial z_1}$$

$$\frac{\partial J}{\partial w_A} = \frac{\partial J}{\partial z_1} \frac{\partial z_1}{\partial w_A}$$

# Backpropagation (updated)

Compute derivatives per layer, utilizing previous derivatives

Objective:  $J(\mathbf{w})$

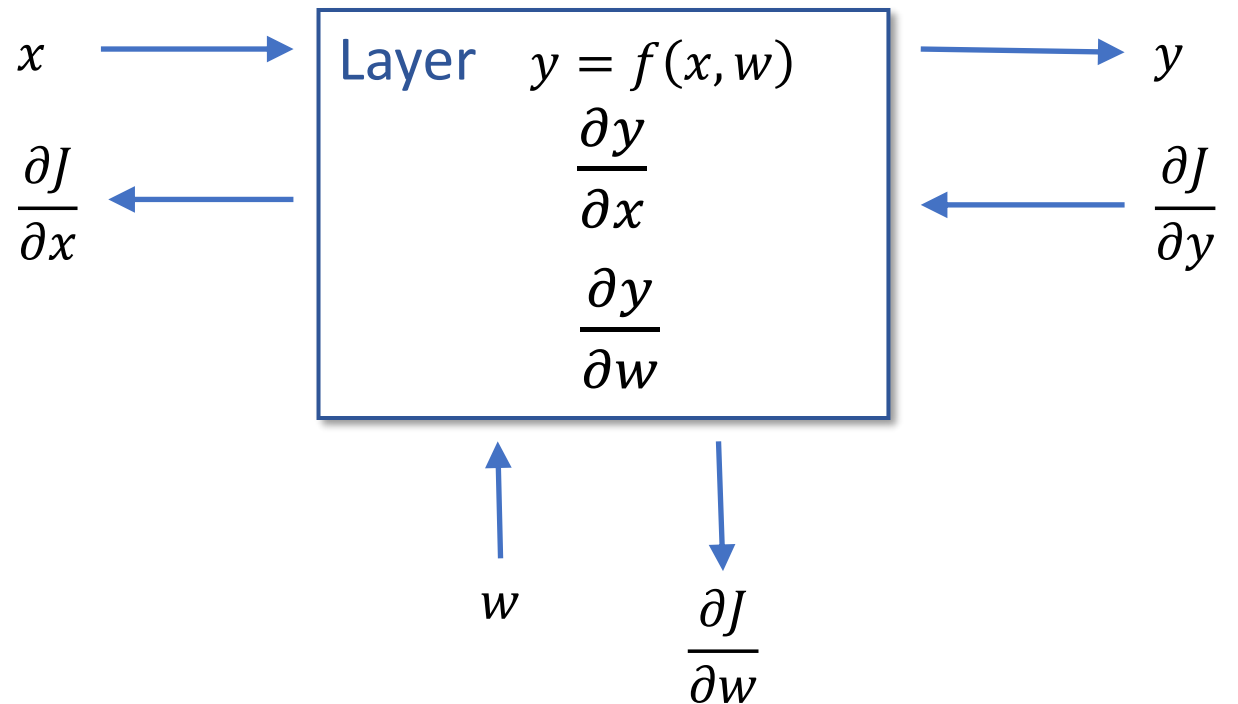
Arbitrary layer:  $y = f(x, w)$

Init:

- $\frac{\partial J}{\partial x} = 0$
- $\frac{\partial J}{\partial w} = 0$

Compute:

- $\frac{\partial J}{\partial x} + = \frac{\partial J}{\partial y} \frac{\partial y}{\partial x}$
- $\frac{\partial J}{\partial w} + = \frac{\partial J}{\partial y} \frac{\partial y}{\partial w}$



# Neural Networks Properties

## Practical considerations

- Large number of neurons
  - Danger for overfitting
- Modelling 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.