# 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

- Neural Networks
  - Perceptron
  - Multilayer perceptron

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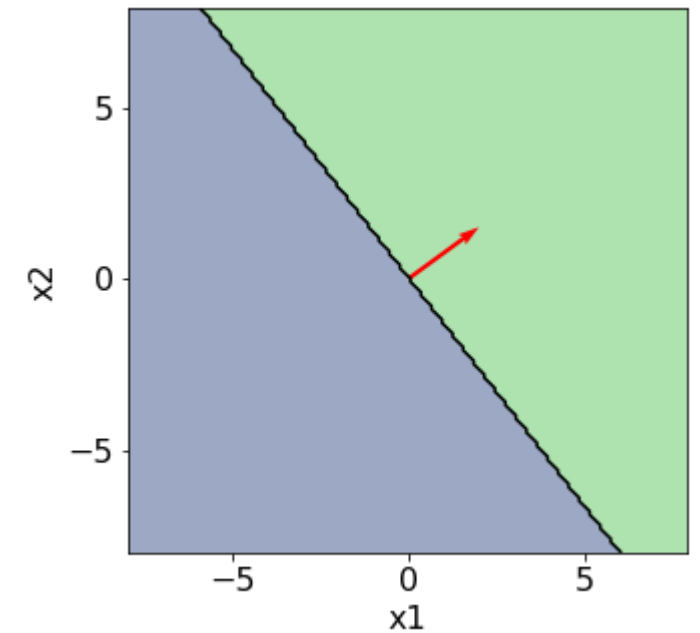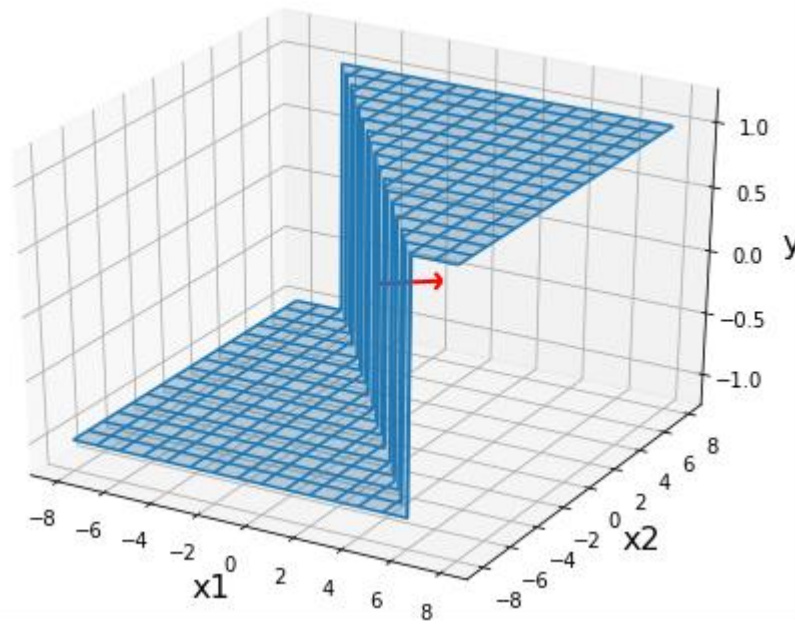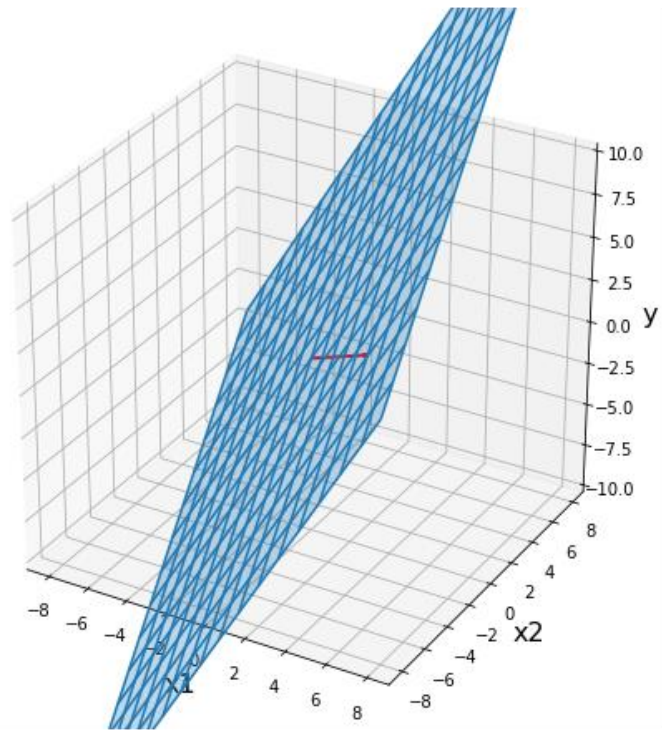- Neural Networks
  - Building blocks
  - Optimization
    - Composite functions and chain rule
    - Forward-backward passes
    - Matrix calculus

# Perceptron

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

Classification: Hard threshold on linear model

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

# Classification Design Challenge

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

$$z_1 = h_A(x) = sign(w_A^T x + b_A)$$
$$z_2 = h_B(x) = sign(w_B^T x + b_B)$$
$$h_C(\pmb{z}) = sign(w_C^T \pmb{z} + b_C)$$
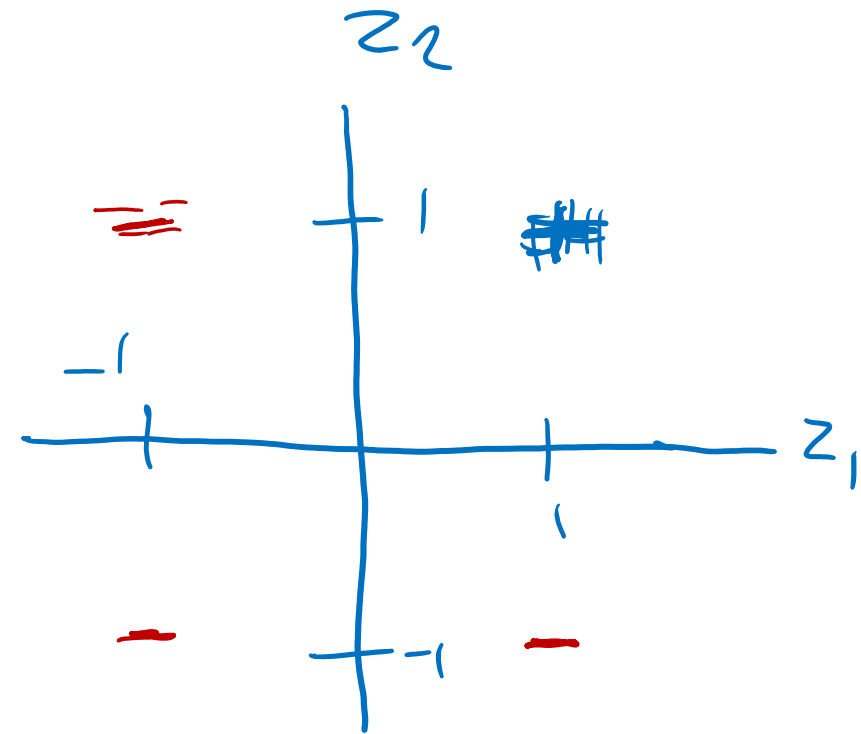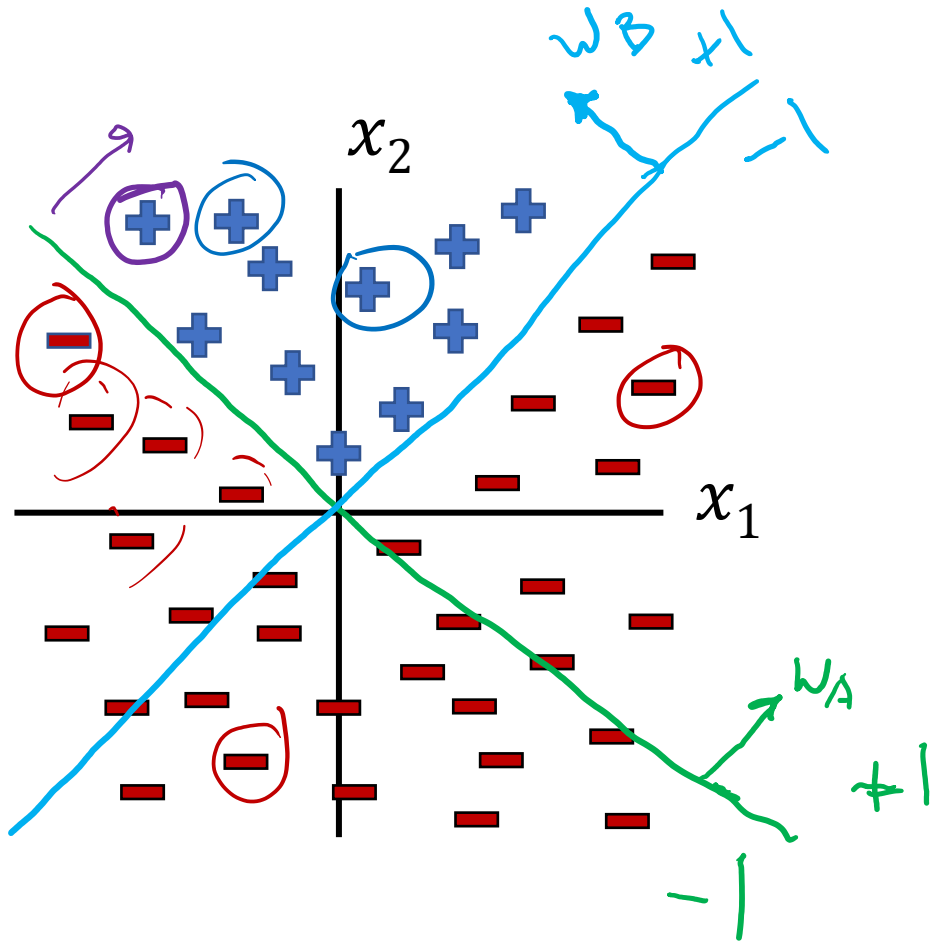
# Classification Design Challenge

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

$$h_A(\boldsymbol{x}) = sign(\boldsymbol{w}_A^T \boldsymbol{x} + b_A)$$
$$h_B(\boldsymbol{x}) = sign(\boldsymbol{w}_B^T \boldsymbol{x} + b_B)$$
$$h_C(\boldsymbol{z}) = sign(\boldsymbol{w}_C^T \boldsymbol{z} + b_C)$$

# Neural Networks
Inspired by actual human brain

Input
Signal

DOG

**CAT**

TREE

CAR

SKY

Image: https://en.wikipedia.org/wiki/Neuron
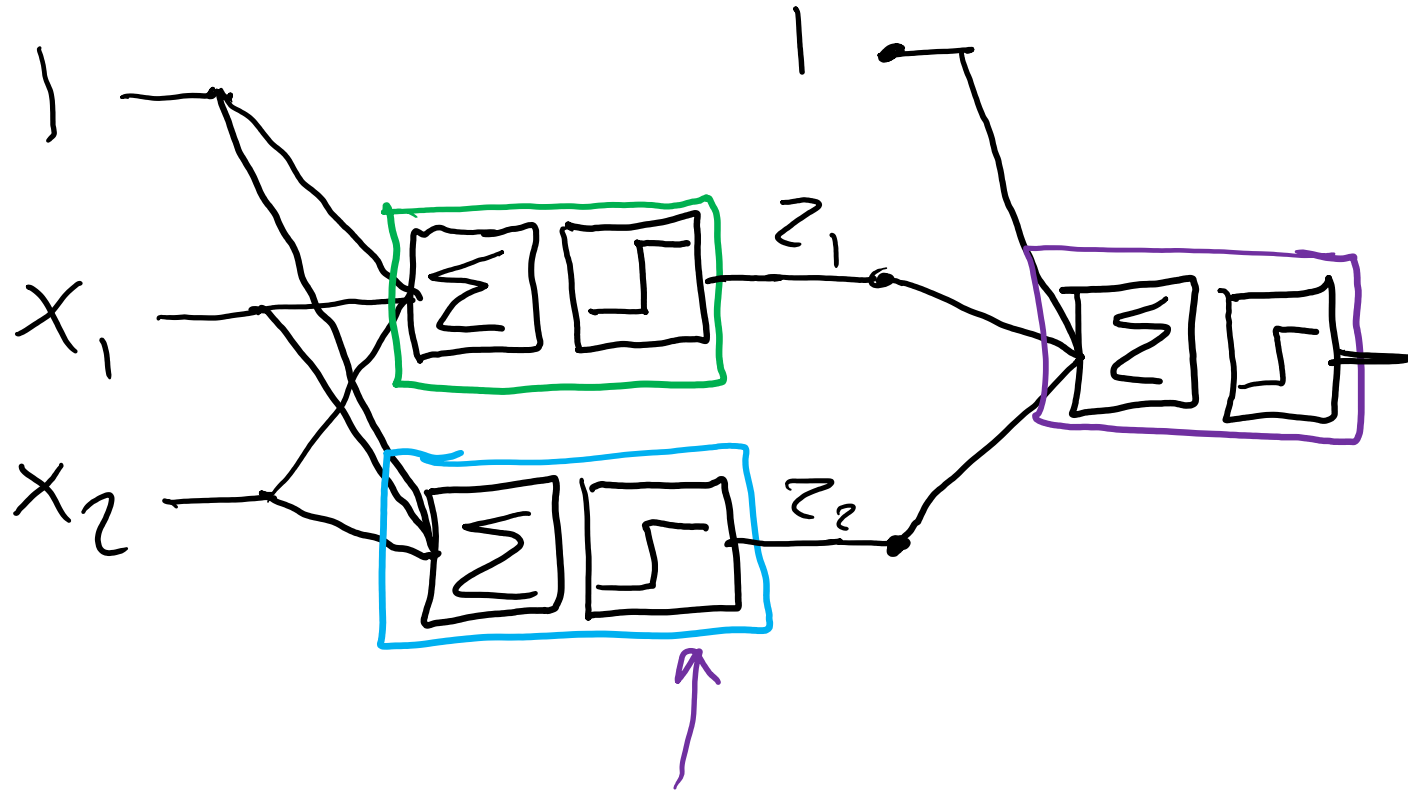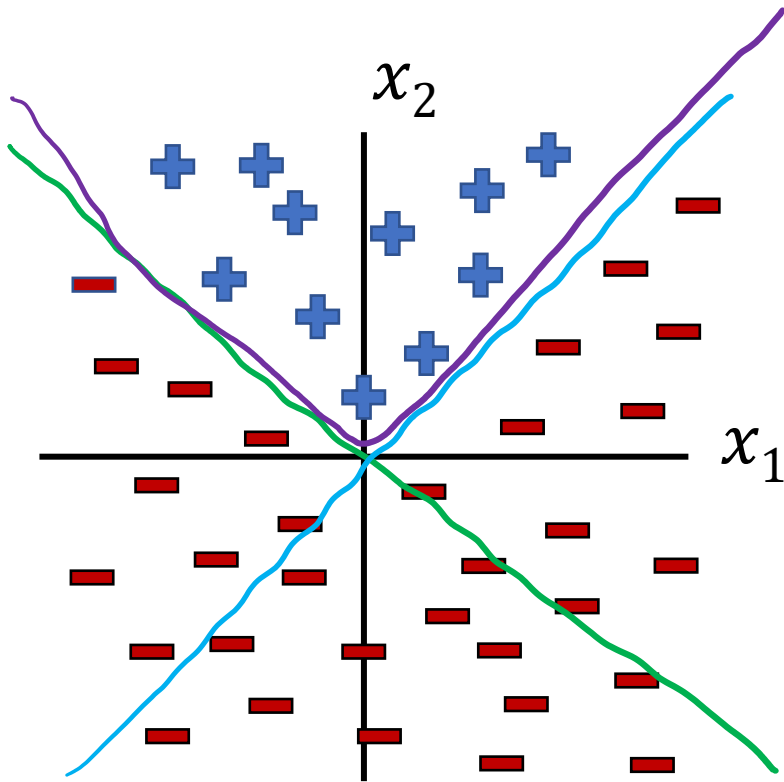
# Classification Design Challenge

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

$$h_A(\boldsymbol{x}) = sign(\boldsymbol{w}_A^T \boldsymbol{x} + b_A)$$
$$h_B(\boldsymbol{x}) = sign(\boldsymbol{w}_B^T \boldsymbol{x} + b_B)$$
$$h_C(\boldsymbol{x}) = sign(\boldsymbol{w}_C^T \boldsymbol{x} + b_C)$$

# Neural Networks

## Simple single neuron example:

- Selling my car

# Neural Networks
Many layers of neurons, millions of parameters

Output
Signal

Input
Signal



DOG

**CAT**

TREE

CAR

SKY

# Neural Networks
## Many layers of neurons, millions of parameters

Output Signal

Input Signal

DOG

**CAT**

TREE

CAR

SKY

# 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 regresion (if $g$ is nothing)



Computed Value

True Label

$z_1$

$y$

$h_{\boldsymbol{w}}(\boldsymbol{x}) = z_1$

$$h_{\boldsymbol{w}}(\boldsymbol{x}) = g\left(\sum_i w_i x_i\right)$$

# Activation Functions

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

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

- Sigmoid: $g(z) = \dfrac{1}{1+e^{-z}}$ $\qquad \dfrac{dg}{dz} = g(z)\big(1 - g(z)\big)$
- (Softmax)

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

# Optimizing

How do we find the "best" set of weights?

$$h_w(\boldsymbol{x}) = g\left(\sum_j w_j x_j\right)$$

# Loss Functions

## Regression

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

## Classification

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

# 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



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

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

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

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

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

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

# Optimizing

How do we find the "best" set of weights?

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

# Neural Network Equations



How do we describe this network?

# Network Optimization Details

# Reminder: Calculus Chain Rule (scalar version)

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

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

# Network Optimization

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

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

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

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

# Network Optimization: Forward then Backwards

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

$z_3 = f_3(w_3, z_2)$

$z_2 = f_2(w_2, z_1)$

$z_1 = f_1(w_1, x)$

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

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

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

Lots of repeated calculations

# 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: $J(\boldsymbol{w})$

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

Need:

- $\dfrac{\partial J}{\partial x} = \dfrac{\partial J}{\partial y}\dfrac{\partial y}{\partial x}$

- $\dfrac{\partial J}{\partial w} = \dfrac{\partial J}{\partial y}\dfrac{\partial y}{\partial w}$

$x \longrightarrow$ 
Layer $\quad y = f(x, w)$
$\dfrac{\partial y}{\partial x}$
$\dfrac{\partial y}{\partial w}$
$\longrightarrow y$

$\dfrac{\partial J}{\partial x} \longleftarrow$

$\longleftarrow \dfrac{\partial J}{\partial y}$

$w$

$\dfrac{\partial J}{\partial w}$

# Matrix Calculus

One way to think of it: Bag of Derivatives

# Matrix Calculus

One way to think of it: Bag of Derivatives

# Matrix Calculus

Jacobian: Vector in, vector out

Numerator-layout

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

# Matrix Calculus

Vector in, scalar out

Numerator-layout

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

# Matrix Calculus

Scalar in, vector out

Numerator-layout

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

# Matrix Calculus

Gradient: Vector in, scalar out

Transpose of numerator-layour

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

# Matrix Calculus

Matrix in, scalar out

Keep same dimensions as matrix

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

# Multivariate Chain Rule

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

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

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

# Multivariate Chain Rule

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

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

$y = f(z_1, z_2) = z_1 z_2$

# 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:

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

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

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

# Network Optimization

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

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

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

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

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

Need multivariate chain rule!

# Network Optimization
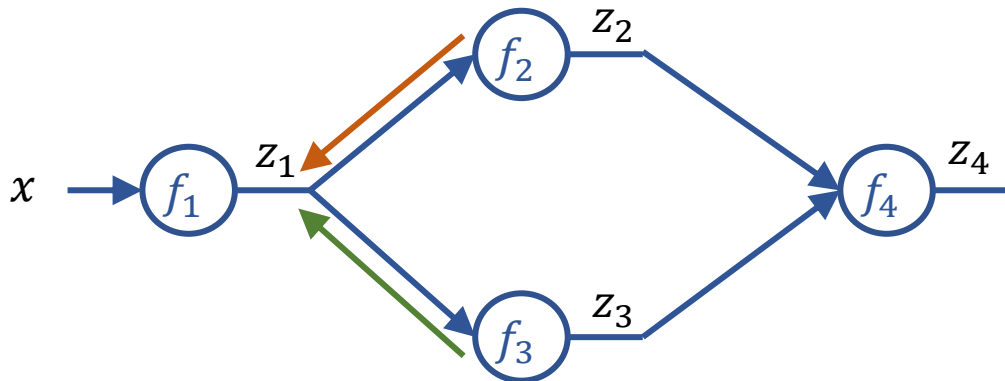
$$J(\boldsymbol{w}) = z_4$$

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

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

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

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

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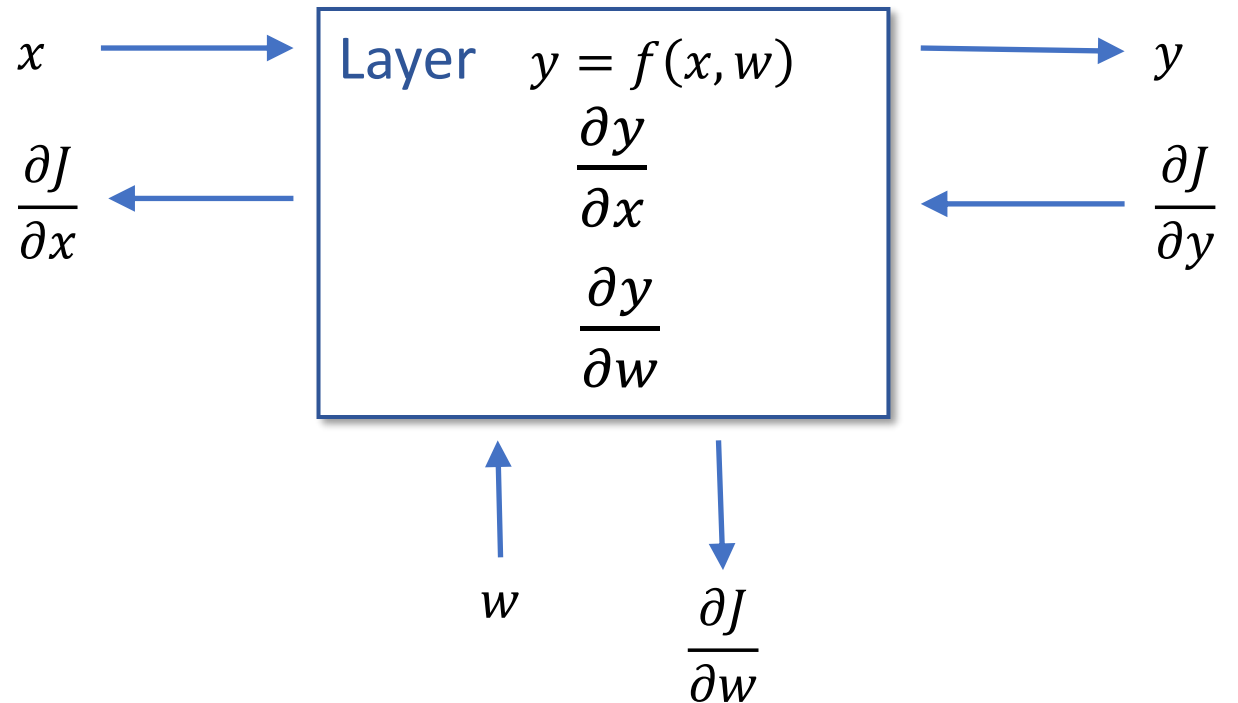
Objective: $J(\boldsymbol{w})$

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

Init:

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

Compute:

- $\dfrac{\partial J}{\partial x} \mathrel{+}= \dfrac{\partial J}{\partial y}\dfrac{\partial y}{\partial x}$

- $\dfrac{\partial J}{\partial w} \mathrel{+}= \dfrac{\partial J}{\partial y}\dfrac{\partial y}{\partial w}$

$x \longrightarrow$

Layer $\quad y = f(x, w)$

$\dfrac{\partial y}{\partial x}$

$\dfrac{\partial y}{\partial w}$

$\longrightarrow y$

$\dfrac{\partial J}{\partial x} \longleftarrow$

$\longleftarrow \dfrac{\partial J}{\partial y}$

$w$

$\dfrac{\partial J}{\partial w}$

# 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.