# Backprop

9

B

**CS 3244**
**Machine Learning**

NUS | Computing
National University
of Singapore

## Perceptron



Don't forget the **bias** term $w_0$

Activation function $g = \text{sgn}$

$\hat{y} = \text{sgn}\left(\sum_{r=0}^{n} w_r x_r\right)$
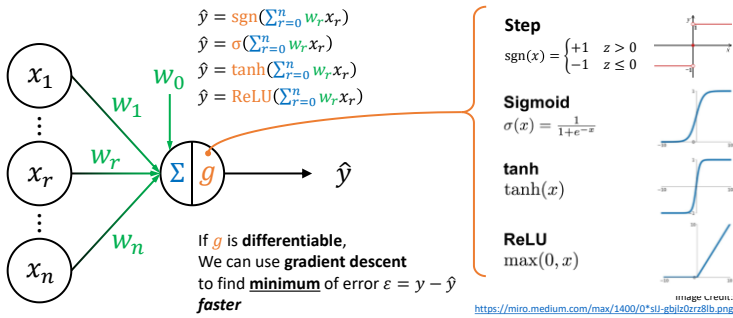
$w_r$ are weights

$x_r$ are inputs

## Perceptron Learning Algorithm

1. Initialize weights $\boldsymbol{w}$
   - Could be all zero, or random small values
2. For each instance $i$ with features $\boldsymbol{x}^{(i)}$
   - Classify $\hat{y}^{(i)} = \text{sgn}(\boldsymbol{w}^\top \boldsymbol{x}^{(i)})$
3. Select one misclassified instance
   - Update weights: $\boldsymbol{w} \leftarrow \boldsymbol{w} + \eta(y - \hat{y})\boldsymbol{x}$
4. Iterate steps 2 to 3 until
   - Convergence (classification error < threshold), or
   - Maximum number of iterations

$$\begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_r \\ \vdots \\ w_n \end{pmatrix} \leftarrow \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_r \\ \vdots \\ w_n \end{pmatrix} + \eta(y - \hat{y}) \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_r \\ \vdots \\ x_n \end{pmatrix}$$

$$w_r \leftarrow w_r + \eta(y - \hat{y})x_r$$

## Differentiable Activation Functions

$\hat{y} = \text{sgn}(\sum_{r=0}^{n} w_r x_r)$
$\hat{y} = \sigma(\sum_{r=0}^{n} w_r x_r)$
$\hat{y} = \tanh(\sum_{r=0}^{n} w_r x_r)$
$\hat{y} = \text{ReLU}(\sum_{r=0}^{n} w_r x_r)$



If $g$ is **differentiable**,
We can use **gradient descent**
to find **minimum** of error $\varepsilon = y - \hat{y}$
***faster***

**Step**
$\text{sgn}(x) = \begin{cases} +1 & z > 0 \\ -1 & z \le 0 \end{cases}$

**Sigmoid**
$\sigma(x) = \frac{1}{1+e^{-x}}$

**tanh**
$\tanh(x)$

**ReLU**
$\max(0, x)$

Image Credit:
https://miro.medium.com/max/1400/0*sIJ-gbjIz0zrz8lb.png

## Multi-Layer Perceptron (Neural Network)



Layer $l - 1$   Layer $l$   Layer $l + 1$

## Chain Rule

Consider composite function

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

$$g = g(f), f = f(x)$$

$$g'(x) = \frac{dg}{dx} = \frac{dg}{df}\frac{df}{dx}$$

**Intuition**
Rate of change of $g$ relative to $x$ is the product of
- rates of change of $g$ relative to $f$ and
- rates of change of $f$ relative to $x$

"If
- a car travels 2x fast as a bicycle and
- the bicycle is 4x as fast as a walking man,
then the car travels 2 × 4 = 8 times as fast as the man."
– George F. Simmons, Calculus with Analytic Geometry (1985)

# Week 09B: Lecture Outline

1. Perceptron
2. Perceptron Learning Algorithm (PLA)
3. Activation Functions
4. Gradient Descent
5. Neural Networks
   - Math Notation Primer
6. Backpropagation

# Math Primer

# Notation

- **Scalar**: not bolded, lower case

$$x$$

- **Vector**: bolded, lower case

$$\boldsymbol{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

- **Matrix**: bolded, upper case

$$\boldsymbol{X} = \begin{pmatrix} x_{11} & \cdots & x_{1m} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nm} \end{pmatrix}$$

# Functions with Vectors and Matrices

- Scalar-by-scalar:

  - $y(x) = wx$ for **scaling** input

- Scalar-by-vector:

  - $y(x) = \boldsymbol{w} \cdot \boldsymbol{x} = \boldsymbol{w}^\top \boldsymbol{x} = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = w_1 x_1 + w_2 x_2$ for **weighted sum**

- Vector-by-vector:

  - $\boldsymbol{y}(x) = w\boldsymbol{x} = w \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} wx_1 \\ wx_2 \end{pmatrix}$ for **scaled** outputs (same weight)

  - $\boldsymbol{\sigma}(\boldsymbol{z}) = \dfrac{e^{\boldsymbol{z}}}{\mathbf{1}^\top (\mathbf{1} + e^{\boldsymbol{z}})}$ for **softmax**

# Functions with Vectors and Matrices

- Matrix-by-matrix:

  - Using **Hadamard** product ∘ for element-wise multiplication
    - $y(X) = W \circ X = \begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix} = \begin{pmatrix} w_{11}x_{11} & w_{12}x_{12} \\ w_{21}x_{21} & w_{22}x_{22} \end{pmatrix}$
  - Using **Convolution** operator ∗ for element-wise multiplication then sum [W08b]
    - $Y(X) = W * X = \begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} \begin{pmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{pmatrix}$

    $= \begin{pmatrix} w_{11}x_{11} + w_{12}x_{12} + w_{21}x_{21} + w_{22}x_{22} & w_{11}x_{12} + w_{12}x_{13} + w_{21}x_{22} + w_{22}x_{23} \\ w_{11}x_{21} + w_{12}x_{22} + w_{21}x_{31} + w_{22}x_{32} & w_{11}x_{22} + w_{12}x_{23} + w_{21}x_{32} + w_{22}x_{33} \end{pmatrix}$

    - For computer vision filters (kernels)

> 1D Vectors, and 2D Matrices and ≥2D Tensors are for convenient notation of multiple similar calculations.

# Weighted Sum

**Summation Series** = Scalar

$$\sum_{r=0}^{n} w_n x_r$$

$$w_1 x_1 + \cdots + w_r x_r + \cdots + w_n x_n$$

**Transposed Vector Multiplication** = Scalar

$$\boldsymbol{w}^\top \boldsymbol{x} = \begin{pmatrix} w_1 & \cdots & w_r & \cdots & w_n \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_r \\ \vdots \\ x_n \end{pmatrix}$$

**Vector Dot Product** = Scalar

$$\boldsymbol{w} \cdot \boldsymbol{x} = \begin{pmatrix} w_1 \\ \vdots \\ w_r \\ \vdots \\ w_n \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ \vdots \\ x_r \\ \vdots \\ x_n \end{pmatrix}$$

**Transposed Matrix Multiplication** = Vector

$$\boldsymbol{W}^\top \boldsymbol{x} = \begin{pmatrix} w_{11} & \cdots & w_{1r} & \cdots & w_{1n} \\ \vdots & \ddots & \vdots & \iddots & \vdots \\ w_{r1} & \cdots & w_{rr} & \cdots & w_{rn} \\ \vdots & \iddots & \vdots & \ddots & \vdots \\ w_{n1} & \cdots & w_{nr} & \cdots & w_{nn} \end{pmatrix}^\top \begin{pmatrix} x_1 \\ \vdots \\ x_r \\ \vdots \\ x_n \end{pmatrix}$$

# Gradient

- **Derivative**: $d$
  - $\frac{dy}{dx}$ is the derivative of $y$ relative to $x$
- **Partial derivative**: $\partial$
  - $\frac{\partial y}{\partial x_1}$ is the derivative of $y$ relative to $x_1$
  - But $y$ also depends on other variables (e.g., $x_2$ so, we can also calculate $\frac{\partial y}{\partial x_2}$)
- **Gradient**: $\nabla$
  - To calculate the derivative relative to *all* $x_1$ and $x_2$ together
  - $\nabla y(\boldsymbol{x})$ is the gradient of $y$ relative to all variables $\boldsymbol{x} = (x_1, \dots, x_n)^\top$

$$\nabla y(\boldsymbol{x}) = \frac{dy}{d\boldsymbol{x}} = \begin{pmatrix} \partial y/\partial x_1 \\ \vdots \\ \partial y/\partial x_n \end{pmatrix} = \begin{pmatrix} \frac{\partial y}{\partial x_1} & \cdots & \frac{\partial y}{\partial x_n} \end{pmatrix}^\top$$

Vector in denominator means Derivative for each variable is put in separate, corresponding variable dimensions

- Assumes Cartesian coordinates (linear, orthogonal)

# Matrix Calculus

$n$ = Number of features in $\boldsymbol{x}$
$m$ = Number of instances in dataset
$N$ = Number of $y$ prediction tasks

Scalar-by-Vector (1D Vector)

$$\frac{dy}{d\boldsymbol{x}} = \begin{pmatrix} \dfrac{\partial y}{\partial x_1} \\ \vdots \\ \dfrac{\partial y}{\partial x_n} \end{pmatrix}$$

Scalar-by-Matrix (2D Matrix)

$$\frac{dy}{d\boldsymbol{X}} = \begin{pmatrix} \dfrac{\partial y}{\partial x_{11}} & \cdots & \dfrac{\partial y}{\partial x_{1m}} \\ \vdots & \ddots & \vdots \\ \dfrac{\partial y}{\partial x_{n1}} & \cdots & \dfrac{\partial y}{\partial x_{nm}} \end{pmatrix}$$

Vector-by-Vector (2D Matrix) – not in exam

$$\frac{d\boldsymbol{y}}{d\boldsymbol{x}} = \begin{pmatrix} \dfrac{\partial y_1}{\partial x_1} & \cdots & \dfrac{\partial y_N}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \dfrac{\partial y_1}{\partial x_n} & \cdots & \dfrac{\partial y_N}{\partial x_n} \end{pmatrix}$$

Vector-by-Matrix (3D Tensor) – not in exam

$$\frac{d\boldsymbol{y}}{d\boldsymbol{X}} = \begin{pmatrix} \dfrac{\partial y_1}{\partial x_{11}} & \cdots & \dfrac{\partial y_1}{\partial x_{1m}} \\ \vdots & \ddots & \vdots \\ \dfrac{\partial y_1}{\partial x_{n1}} & \cdots & \dfrac{\partial y_1}{\partial x_{nm}} \end{pmatrix} \ddots \begin{pmatrix} \dfrac{\partial y_N}{\partial x_{11}} & \cdots & \dfrac{\partial y_N}{\partial x_{1m}} \\ \vdots & \ddots & \vdots \\ \dfrac{\partial y_N}{\partial x_{n1}} & \cdots & \dfrac{\partial y_N}{\partial x_{nm}} \end{pmatrix}$$

Along 3rd dimension

This math informs what matrix **shapes** you need to implement
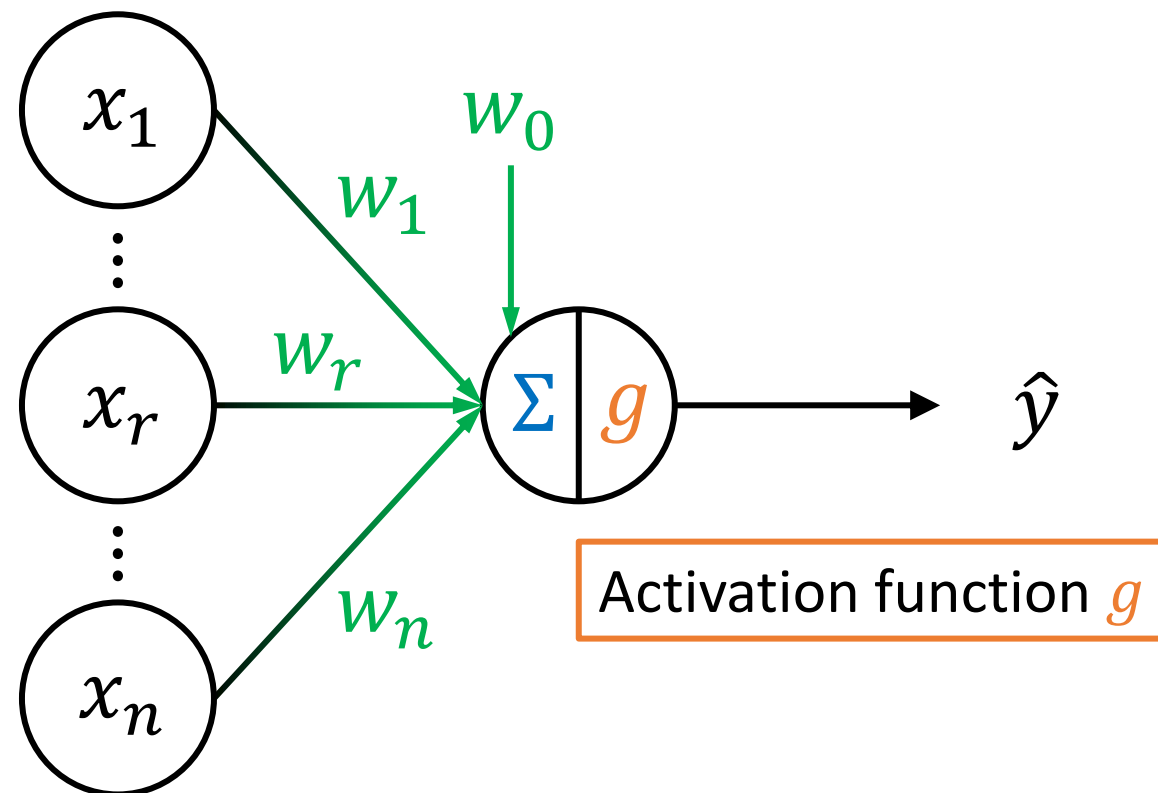
# Neural Network
## (recap)

# Single-Layer Perceptron



Activation function $g$

# Single-Layer Perceptron



**Activation**

$$a = g\big(f(\boldsymbol{x})\big), f(\boldsymbol{x}) = \boldsymbol{w}^\top \boldsymbol{x}$$

# Neural Network



$$W^{[l-1]} \qquad W^{[l]} \qquad W^{[l+1]} \qquad W^{[l+2]}$$

$$a^{[l-1]} \qquad a^{[l]} \qquad a^{[l+1]}$$

Layer $l-1$       Layer $l$       Layer $l+1$

# Layer Activation

$$a = g\big(f(x)\big), f(x) = \boldsymbol{w}^\top x$$

Single-Layer
Perceptron

Layer $l$
Activation
Function

Layer $l$
Weights

$$\boldsymbol{a}^{[l]} = g^{[l]}\left(\big(\boldsymbol{W}^{[l]}\big)^\top \boldsymbol{a}^{[l-1]}\right)$$

Layer $l$ in
Neural Network

Layer $l$
Activations

Layer $l-1$
Activations

# Gradient Descent Weight Update (Single Neuron)

Old weight

Direction of fastest error increase

Gradient of error

$$\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \nabla \varepsilon$$

$$\nabla \varepsilon = \frac{d\varepsilon}{d\boldsymbol{w}} = \begin{pmatrix} \partial \varepsilon / \partial w_1 \\ \vdots \\ \partial \varepsilon / \partial w_r \\ \vdots \\ \partial \varepsilon / \partial w_n \end{pmatrix}$$

New weight

Learning Rate

# Gradient Descent Weight Update (Neural Network)

Old weight

Direction of fastest error increase

Gradient of error

$$\boldsymbol{W} \leftarrow \boldsymbol{W} - \eta \nabla \varepsilon$$

New weight

Learning Rate

$$\nabla \varepsilon = \frac{d\varepsilon}{d\boldsymbol{W}} = \begin{pmatrix} \partial \varepsilon / \partial \boldsymbol{W}^{[1]} \\ \vdots \\ \partial \varepsilon / \partial \boldsymbol{W}^{[l]} \\ \vdots \\ \partial \varepsilon / \partial \boldsymbol{W}^{[L]} \end{pmatrix}$$

$$\frac{d\varepsilon}{d\boldsymbol{W}^{[l]}} = \frac{d\varepsilon}{d\hat{y}} \frac{d\hat{y}}{d\boldsymbol{W}^{[l]}}$$

$$\frac{d\varepsilon}{d\boldsymbol{W}^{[l]}} = \begin{pmatrix} \dfrac{\partial \varepsilon}{\partial w_{11}^{[l]}} & \cdots & \dfrac{\partial \varepsilon}{\partial w_{1m}^{[l]}} \\ \vdots & \ddots & \vdots \\ \dfrac{\partial \varepsilon}{\partial w_{n1}^{[l]}} & \cdots & \dfrac{\partial \varepsilon}{\partial w_{nm}^{[l]}} \end{pmatrix}$$

# Gradient Descent
# for Neural Networks
# Backpropagation

# Backpropagation

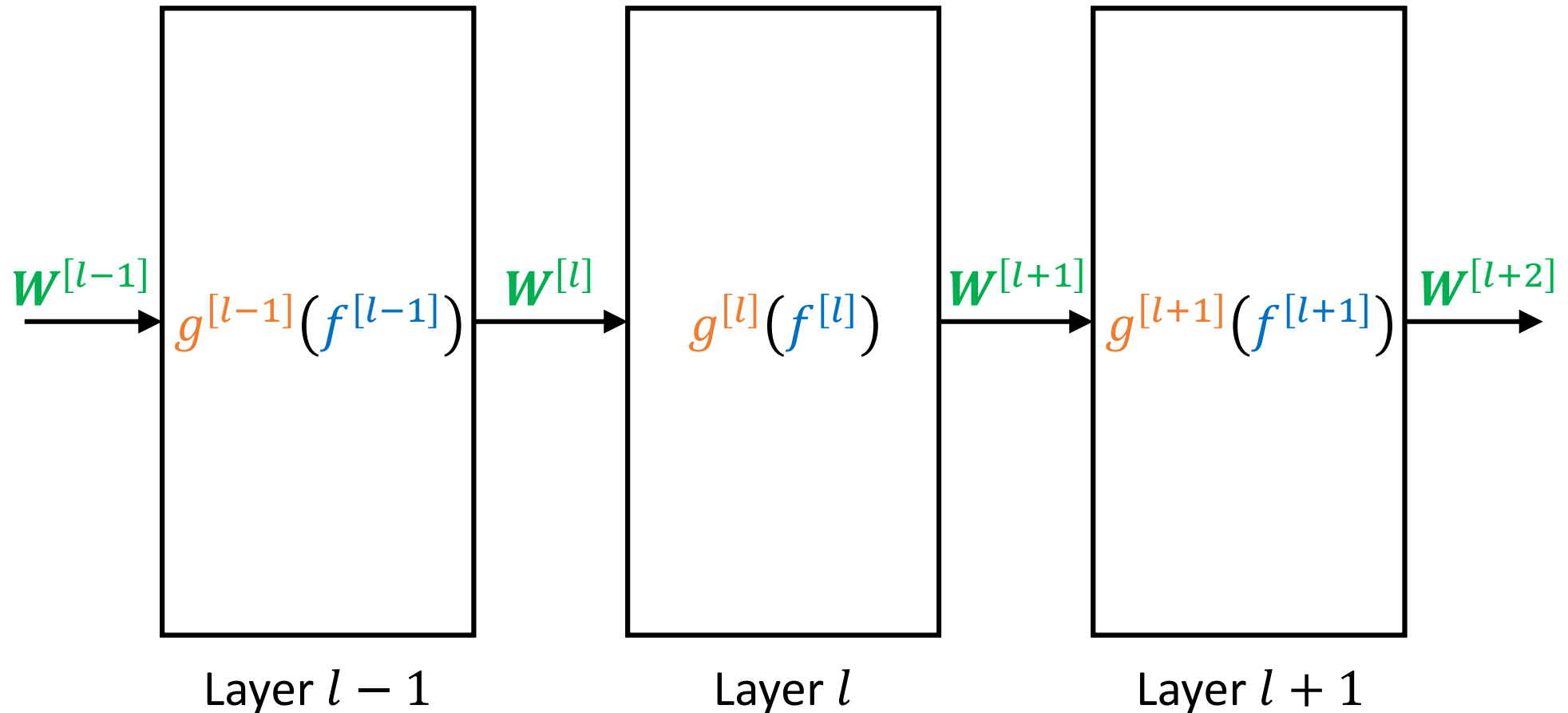Backpropagation **efficiently** computes the **<u>gradient</u>** by

- Avoiding duplicate calculations

- Not computing unnecessary intermediate values,

- Computing the **<u>gradient</u>** of *each* **layer**

Specifically, the gradient of the weighted input of each layer

is calculated from back $[l+1]$ to front $[l]$ :

$$\frac{d\hat{y}}{d\boldsymbol{W}^{[l]}} = \boldsymbol{a}^{[l-1]}\left(\boldsymbol{\delta}^{[l]}\right)^{\top} \quad \bigg| \quad \boldsymbol{\delta}^{[l]} = \left(\frac{dg^{[l]}}{df^{[l]}}\right)\left(\boldsymbol{W}^{[l+1]}\boldsymbol{\delta}^{[l+1]}\right)$$

Adapted from: https://en.wikipedia.org/wiki/Backpropagation

# Forward Propagation



$W^{[l-1]}$ → $g^{[l-1]}(f^{[l-1]})$ → $W^{[l]}$ → $g^{[l]}(f^{[l]})$ → $W^{[l+1]}$ → $g^{[l+1]}(f^{[l+1]})$ → $W^{[l+2]}$

Layer $l-1$      Layer $l$      Layer $l+1$

# Forward Propagation (Reverse Polish Notation)

$$(x^{[0]}, W^{[1]}) f^{[1]} g^{[1]} \cdots f^{[l]} g^{[l]} \cdots f^{[L-2]} g^{[L-2]} f^{[L-1]} g^{[L-1]} f^{[L]} g^{[L]} = \hat{y}$$

# Forward Propagation (Reverse Polish Notation)

$$\left(x^{[0]}, W^{[1]}\right) f^{[1]} g^{[1]} \cdots f^{[l]} g^{[l]} \cdots f^{[L-2]} g^{[L-2]} f^{[L-1]} g^{[L-1]} f^{[L]} g^{[L]} = \hat{y}$$

$$\left(a^{[L-1]}, W^{[L]}\right) f^{[L]} g^{[L]} = \hat{y}$$

$$\left(a^{[L-2]}, W^{[L-1]}\right) f^{[L-1]} g^{[L-1]} f^{[L]} g^{[L]} = \hat{y}$$

$$\left(a^{[L-3]}, W^{[L-2]}\right) f^{[L-2]} g^{[L-2]} f^{[L-1]} g^{[L-1]} f^{[L]} g^{[L]} = \hat{y}$$

# Gradients of Layer Weights (Backwards)

$$\left(x^{[0]}, W^{[1]}\right) f^{[1]} g^{[1]} \cdots f^{[l]} g^{[l]} \cdots f^{[L-2]} g^{[L-2]} f^{[L-1]} g^{[L-1]} f^{[L]} g^{[L]} = \hat{y}$$

$$\left(a^{[L-1]}, W^{[L]}\right) f^{[L]} g^{[L]} = \hat{y}$$

$$\frac{df^{[L]}}{dW^{[L]}} \frac{dg^{[L]}}{df^{[L]}} = \frac{\partial \hat{y}}{\partial W^{[L]}}$$

$$\left(a^{[L-2]}, W^{[L-1]}\right) f^{[L-1]} g^{[L-1]} f^{[L]} g^{[L]} = \hat{y}$$

$$\frac{df^{[L-1]}}{dW^{[L-1]}} \frac{dg^{[L-1]}}{df^{[L-1]}} \frac{df^{[L]}}{dg^{[L-1]}} \frac{dg^{[L]}}{df^{[L]}} = \frac{\partial \hat{y}}{\partial W^{[L-1]}}$$

$$\left(a^{[L-3]}, W^{[L-2]}\right) f^{[L-2]} g^{[L-2]} f^{[L-1]} g^{[L-1]} f^{[L]} g^{[L]} = \hat{y}$$

$$\frac{df^{[L-2]}}{dW^{[L-2]}} \frac{dg^{[L-2]}}{df^{[L-2]}} \frac{df^{[L-1]}}{dg^{[L-2]}} \frac{dg^{[L-1]}}{df^{[L-1]}} \frac{df^{[L]}}{dg^{[L-1]}} \frac{dg^{[L]}}{df^{[L]}} = \frac{\partial \hat{y}}{\partial W^{[L-2]}}$$

# Gradients of Layer Weights (Backwards)

$$\left(x^{[0]}, W^{[1]}\right) f^{[1]} g^{[1]} \cdots f^{[l]} g^{[l]} \cdots f^{[L-2]} g^{[L-2]} f^{[L-1]} g^{[L-1]} f^{[L]} g^{[L]} = \hat{y}$$

$$\left(a^{[L-1]}, W^{[L]}\right) f^{[L]} g^{[L]} = \hat{y}$$

$$\frac{df^{[L]}}{dW^{[L]}} \boxed{\delta^{[L]}} = \frac{\partial \hat{y}}{\partial W^{[L]}}$$

$$\left(a^{[L-2]}, W^{[L-1]}\right) f^{[L-1]} g^{[L-1]} f^{[L]} g^{[L]} = \hat{y}$$

$$\frac{df^{[L-1]}}{dW^{[L-1]}} \boxed{\delta^{[L-1]}} = \frac{\partial \hat{y}}{\partial W^{[L-1]}}$$

$$\left(a^{[L-3]}, W^{[L-2]}\right) f^{[L-2]} g^{[L-2]} f^{[L-1]} g^{[L-1]} f^{[L]} g^{[L]} = \hat{y}$$

$$\frac{df^{[L-2]}}{dW^{[L-2]}} \boxed{\delta^{[L-2]}} = \frac{\partial \hat{y}}{\partial W^{[L-2]}}$$

# **Recursive** Gradients of Layer Weights

$$\frac{\partial \hat{y}}{\partial W^{[l]}} = \frac{df^{[l]}}{dW^{[l]}} \frac{dg^{[l]}}{df^{[l]}} \frac{df^{[l+1]}}{dg^{[l]}} \frac{dg^{[l+1]}}{df^{[l+1]}} \cdots \frac{df^{[L]}}{dg^{[L-1]}} \frac{dg^{[L]}}{df^{[L]}}$$

$$\frac{\partial \hat{y}}{\partial W^{[l]}} = \frac{df^{[l]}}{dW^{[l]}} \frac{dg^{[l]}}{df^{[l]}} \frac{df^{[l+1]}}{dg^{[l]}} \delta^{[l+1]}$$

$$\frac{\partial \hat{y}}{\partial W^{[l]}} = \frac{df^{[l]}}{dW^{[l]}} \delta^{[l]}$$

$$\delta^{[l]} = \frac{dg^{[l]}}{df^{[l]}} \frac{df^{[l+1]}}{dg^{[l]}} \delta^{[l+1]}$$

**Reference**

$$a^{[l]} = g^{[l]}(f^{[l]})$$

$$f^{[l]} = (W^{[l]})^{\top} a^{[l-1]}$$

$$a^{[l-1]} = g^{[l-1]}$$

$$\frac{df^{[l]}}{dW^{[l]}} = \frac{d\left((W^{[l]})^{\top} a^{[l-1]}\right)}{dW^{[l]}} = a^{[l-1]}$$

$$\frac{df^{[l+1]}}{dg^{[l]}} = \frac{df^{[l+1]}}{da^{[l]}} = W^{[l+1]}$$

$$\frac{\partial \hat{y}}{\partial W^{[l]}} = a^{[l-1]} \delta^{[l]} \quad \bigg| \quad \delta^{[l]} = \frac{dg^{[l]}}{df^{[l]}} W^{[l+1]} \delta^{[l+1]}$$

**Recursive**

# Matrix multiplication to match shape (not in exam)

$$\frac{d\hat{y}}{d\boldsymbol{W}^{[l]}} = \begin{pmatrix} \dfrac{\partial \hat{y}}{\partial w_{11}^{[l]}} & \cdots & \dfrac{\partial \varepsilon}{\partial w_{1n^{[l]}}^{[l]}} \\ \vdots & \ddots & \vdots \\ \dfrac{\partial \varepsilon}{\partial w_{n^{[l-1]}1}^{[l]}} & \cdots & \dfrac{\partial \varepsilon}{\partial w_{n^{[l-1]}n^{[l]}}^{[l]}} \end{pmatrix}$$

$\boldsymbol{a}^{[l]}$

$\boldsymbol{a}^{[l-1]}$   $\boldsymbol{W}^{[l]}$

| Activation $\hat{y} = g(f)$ | | $\dfrac{dg}{df}$ |
|---|---|---|
| Sigmoid | $\dfrac{1}{1 + e^{-f}}$ | $(1 - g)g$ |
| tanh | $\tanh f$ | $1 - g^2$ |
| ReLU | $\max(0, f)$ | $[f > 0]$ |

$n_{\text{rows}} \times n_{\text{cols}}$   $n^{[l-1]} \times n^{[l]}$   $n^{[l-1]} \times 1\; 1 \times n^{[l]}$   $n^{[l]} \times 1$   $1 \times 1$   $n^{[l]} \times n^{[l+1]}\; n^{[l+1]} \times 1$

$$\frac{d\hat{y}}{d\boldsymbol{W}^{[l]}} = \boldsymbol{a}^{[l-1]}\left(\boldsymbol{\delta}^{[l]}\right)^{\top} \qquad \boldsymbol{\delta}^{[l]} = \left(\frac{dg^{[l]}}{df^{[l]}}\right)\left(\boldsymbol{W}^{[l+1]}\boldsymbol{\delta}^{[l+1]}\right)$$

# Backward Propagation



$$\frac{d\hat{y}}{dW^{[l]}} = a^{[l-1]}\left(\delta^{[l]}\right)^{\top} \qquad \delta^{[l]} = \left(\frac{dg^{[l]}}{df^{[l]}}\right)\left(W^{[l+1]}\delta^{[l+1]}\right)$$

# Backpropagation

Backpropagation **efficiently** computes the **<u>gradient</u>** by

- Avoiding duplicate calculations
- Not computing unnecessary intermediate values,
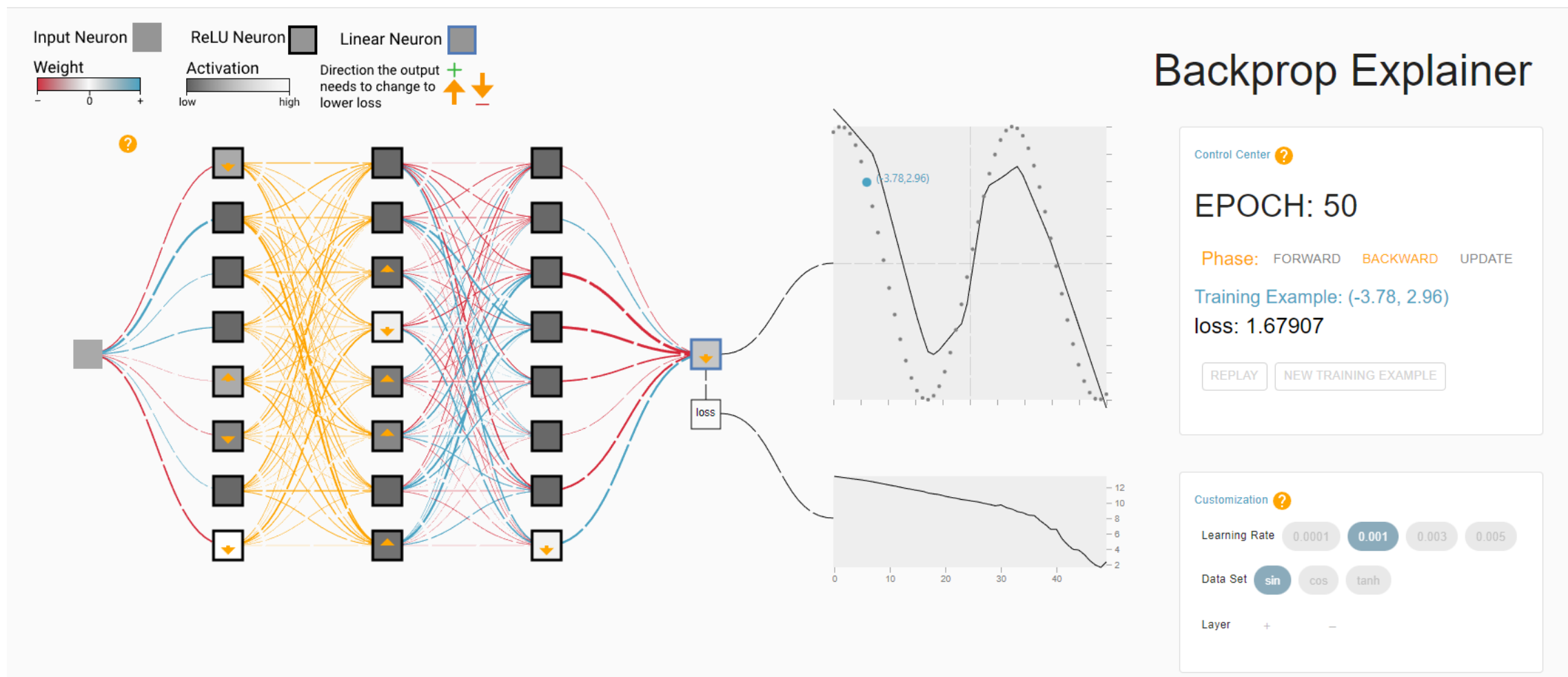- Computing the **<u>gradient</u>** of ***each* layer**

Specifically, the gradient of the weighted input of each layer

is calculated from back $[l+1]$ to front $[l]$ :

$$\frac{d\hat{y}}{d\boldsymbol{W}^{[l]}} = \boldsymbol{a}^{[l-1]}\left(\boldsymbol{\delta}^{[l]}\right)^{\top} \qquad \boldsymbol{\delta}^{[l]} = \left(\frac{d\boldsymbol{g}^{[l]}}{d\boldsymbol{f}^{[l]}}\right)\left(\boldsymbol{W}^{[l+1]}\boldsymbol{\delta}^{[l+1]}\right)$$

Adapted from: https://en.wikipedia.org/wiki/Backpropagation

https://xnought.github.io/backprop-explainer/

# Insert Web Page

This app allows you to insert secure web pages starting with https:// into the slide deck. Non-secure web pages are not supported for security reasons.

Please enter the URL below.

| https:// | xnought.github.io/backprop-explainer/ |
|----------|---------------------------------------|

Note: Many popular websites allow secure access. Please click on the preview button to ensure the web page is accessible.

# Practice Backprop during tutorial

CS3244, Solution to Tutorial 07—Perceptrons and Neural Networks                    1

National University of Singapore
School of Computing
CS3244: Machine Learning
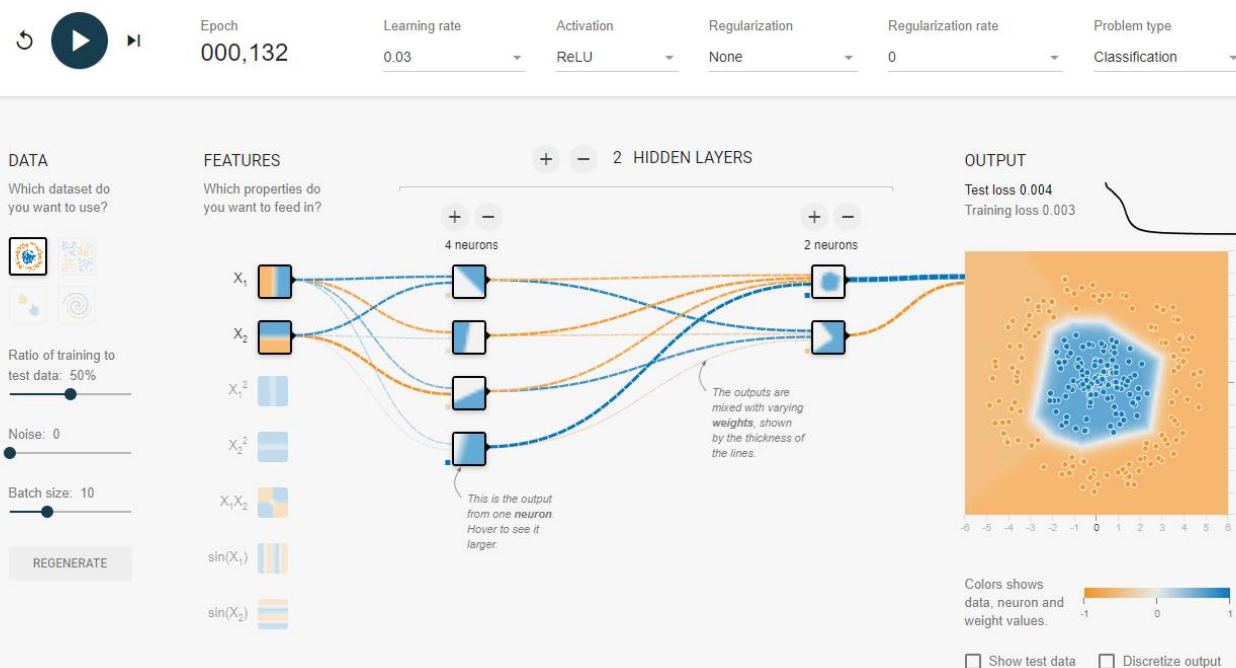Solution to Tutorial 07

**Perceptrons and Neural Networks**

**Colab Notebook :**   Perceptrons and Neural Networks

1. **Backpropagation algorithm.** In this question, we're going to use a neural network with **a**

# Resources for self-study

- [What is backpropagation really doing?](), [Backpropagation calculus]() - [3Blue1Brown]()
- [A worked example of backpropagation]() - Alexander Schiendorfer
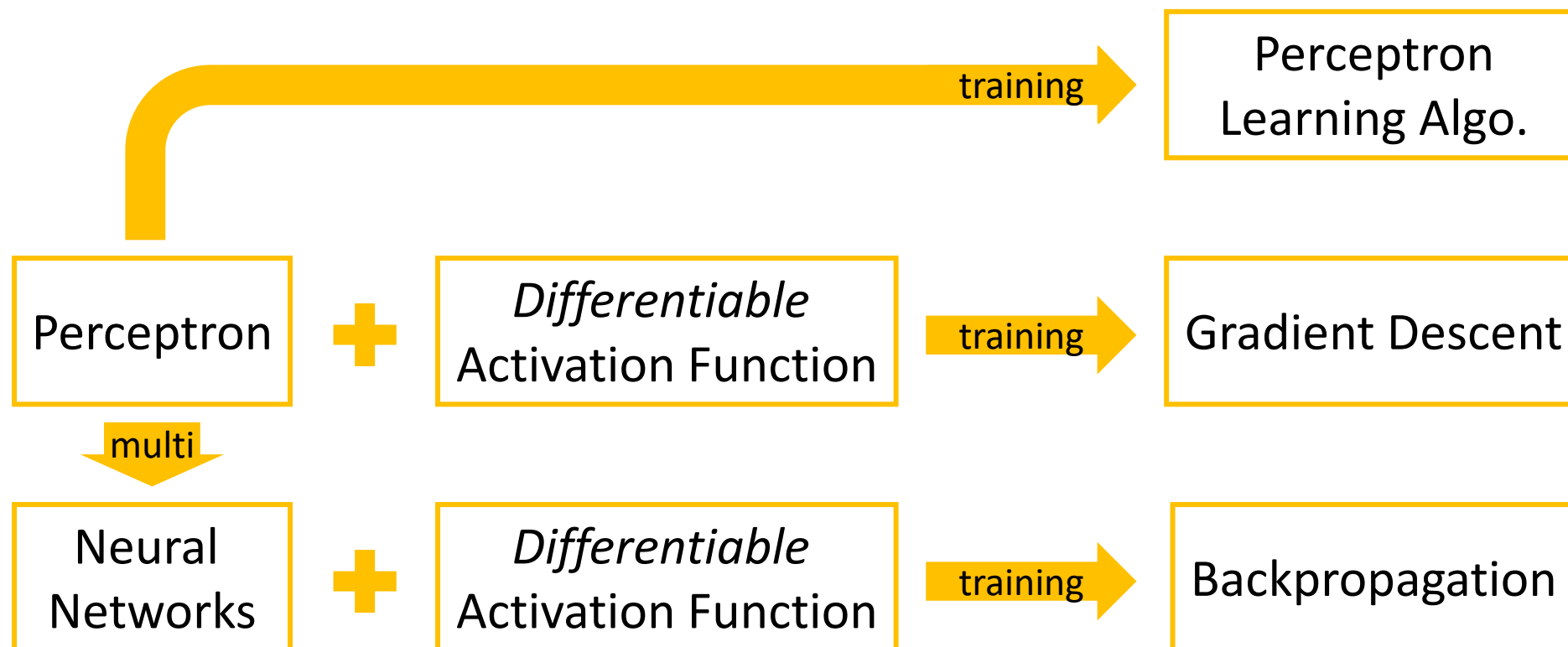- [TensorFlow Playground]()

# Auto Differentiation for Backprop

- Even with backprop, implementing the gradients is tedious
- Deep learning APIs have automated differentiation.
    - Tensor Flow autodiff
    - PyTorch autograd
    - Implement derivatives of many common functions
    - You just need to implement your layers and neurons; API will handle gradients
- Caution
    - If you want to implement **custom functions/layers** (not simple weighted sum)
    - They need to be **differentiable** to be able to calculate their **gradients**
    - Otherwise, backprop cannot update weights accurately

# Wrapping Up

# AI
## HISTORY

# What did we learn?

Next week (W10A):
Deep Learning (CNN)

Image credit: https://rigu.co.uk/kaleidoscope-filter-77mm

NUS CS3244: Machine Learning

Next week (W10B):
Deep Learning (RNN)

Image credit: https://flaunt.com/content/anthony-james-portal-series

NUS CS3244: Machine Learning

# W10 Pre-Lecture Task (due before next Mon)

**Watch**

- [Who Invented A.I.? - The Pioneers of Our Future](#) by [ColdFusion](#)

**Play**

- [https://distill.pub/2018/building-blocks/](#)
  - Don't worry about reading the whole article

**Discuss**

1. Identify what is strange, funny, or erroneous in the deep learning model in Building-Blocks
2. Take a screenshot of the issue and share with your tutorial mates
3. Try to explain why the model was behaving as identified
3. Post a 2–3 sentence description to the topic in your tutorial group: #tg-xx