# Perceptron & Neural Networks

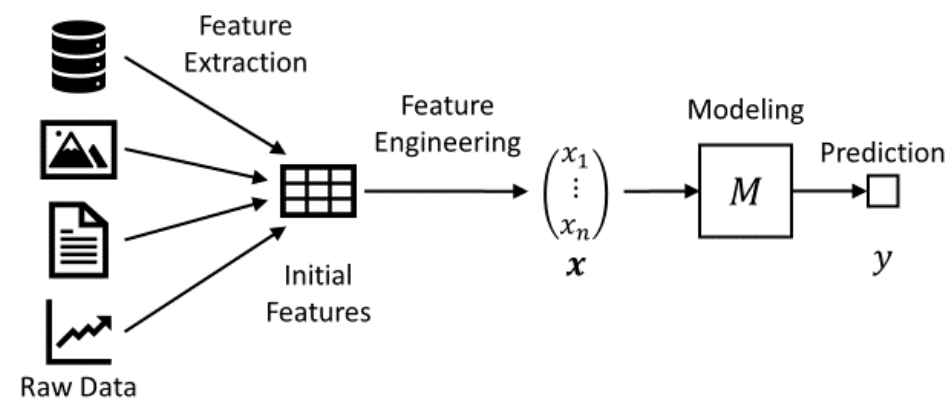**9**

**A**

CS 3244
Machine Learning

NUS | Computing
National University of Singapore

Feature Extraction/Engineering → Modeling
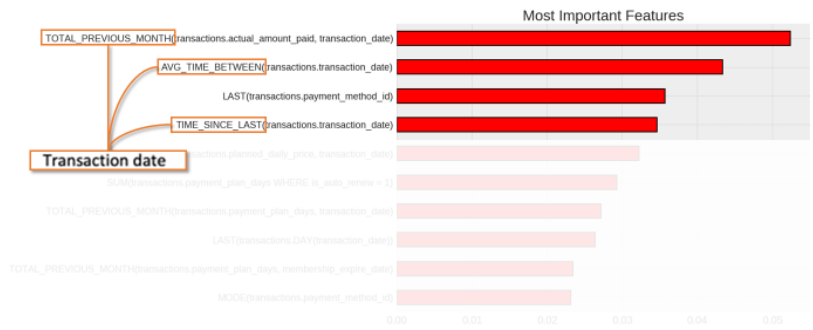
[W08b] Student Learning Outcomes

What did we learn?

1. Describe **issues** when extracting features for various data types
2. Describe **techniques** of feature extraction/engineering for different data types

| Tabular | Temporal | Image | Text |
|---------|----------|-------|------|
| • Domain-specific custom equations<br>• Features from counting, aggregation, difference, min, max | • Features from previous values, aggregate statistics, linear regression<br>• Wave analysis features | • RGB image as 3D tensor<br>• Color features from RGB histogram<br>• Shape features from edge detection<br>• Edge detection via Convolution | • Tokenization<br>• Stemming, Lemmatization<br>• Stop words<br>• Bag-of-Words encoding |

Tabular Feature Engineering:
## Counting, Aggregation, Difference, Min, Max

Most Important Features

TOTAL_PREVIOUS_MONTH(transactions.actual_amount_paid, transaction_date)
AVG_TIME_BETWEEN(transactions.transaction_date)
LAST(transactions.payment_method_id)
TIME_SINCE_LAST(transactions.transaction_date)

Transaction date

Source: https://github.com/Featuretools/predict-customer-churn

NUS CS3244: Machine Learning        4

## Sliding Time Window

- Prediction Task: **Price Prediction**
- Features
  - Moving Average
  - Moving Standard Deviation
  - Moving Range (Min, Max)
  - Moving Trend (Slope of linear fit)

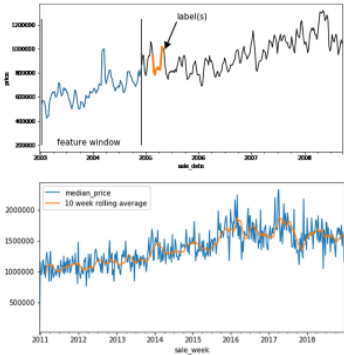Image credit: https://cloud.google.com/blog/products/ai-machine-learning/how-to-quickly-solve-machine-learning-forecasting-problems-using-pandas-and-bigquery
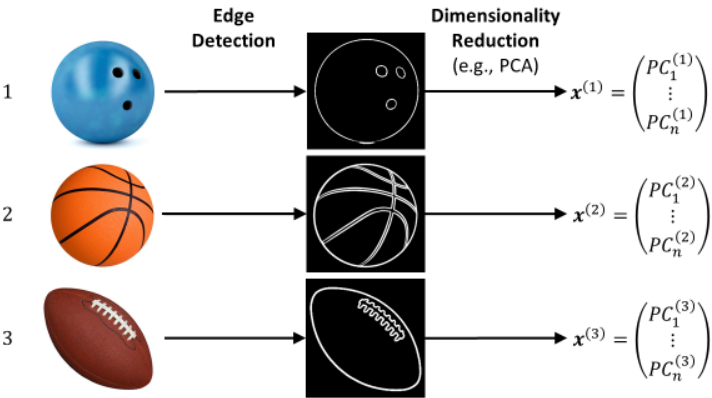
NUS CS3244: Machine Learning        19

## Feature: Edge Detection Kernels (2D)

$$\frac{\partial c}{\partial p_x} \approx \frac{c_{(x+1,y)} - c_{(x-1,y)}}{p_{(x+1,y)} - p_{(x-1,y)}}$$

$$\frac{\partial c}{\partial p_y} \approx \frac{c_{(x,y+1)} - c_{(x,y-1)}}{p_{(x,y+1)} - p_{(x,y-1)}}$$

$$I_{p_x} = \begin{pmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{pmatrix}$$

$$I_{p_y} = \begin{pmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{pmatrix}$$

$$\nabla c = \frac{\partial c}{\partial p_x} + \frac{\partial c}{\partial p_y} \sim \left( I_{p_x} + I_{p_y} \right) * \boldsymbol{x}$$

NUS CS3244: Machine Learning        46

## Feature: Shape Feature Vector

Edge Detection     Dimensionality Reduction (e.g., PCA)

1    $x^{(1)} = \begin{pmatrix} PC_1^{(1)} \\ \vdots \\ PC_n^{(1)} \end{pmatrix}$

2    $x^{(2)} = \begin{pmatrix} PC_1^{(2)} \\ \vdots \\ PC_n^{(2)} \end{pmatrix}$

3    $x^{(3)} = \begin{pmatrix} PC_1^{(3)} \\ \vdots \\ PC_n^{(3)} \end{pmatrix}$

NUS CS3244: Machine Learning

## Bag-of-Words (BOW) Encoding

1. Preprocess string $s$ to array of words $\boldsymbol{w}$
2. Array of words → One-hot vector (fixed length)
3. BOW($\boldsymbol{w}$) → $\boldsymbol{x}$
4. Problem: high dimensions if many words

$$\boldsymbol{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_{13} \end{pmatrix}$$

| # | Original Text $s$ | Pre-Processed Words $\boldsymbol{w}$ | chicken | wings | amazing | honestly | but | way | too | long | wait | not | worth | salty | expensive |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | "Chicken wings were amazing honestly" | ['chicken', 'wings', 'amazing', 'honestly'] | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | "Amazing wings, but waaaay to long to wait." | ['amazing', 'wings', 'but', 'way', 'too', 'long', 'wait'] | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 |
| 3 | "Not worth it! Too salty chicken and expensive!" | ['not', 'worth', 'too', 'salty', 'chicken', 'expensive'] | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 1 |

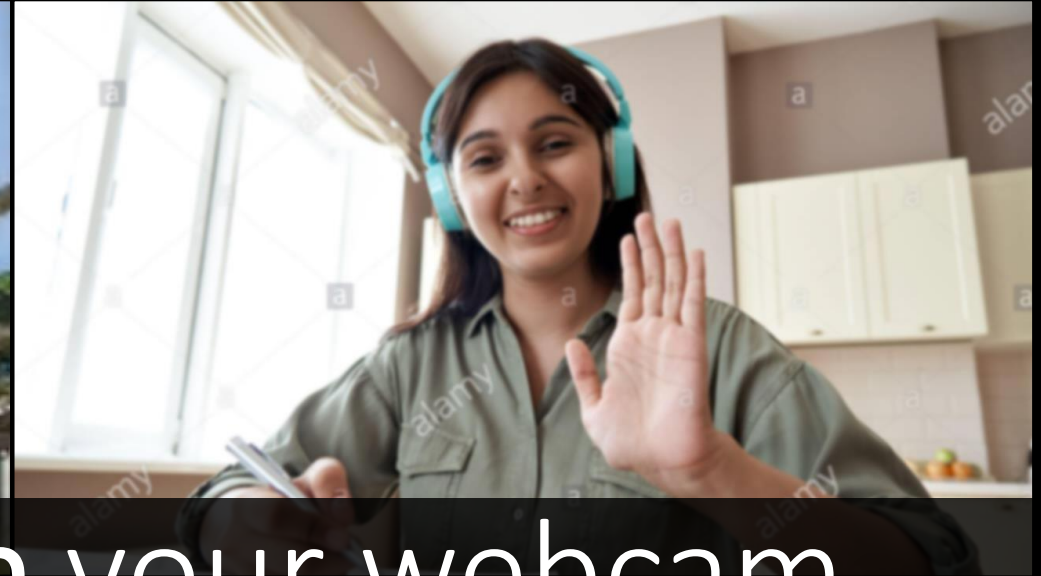The word "too" could predict negative sentiment        46

# Week 09A: Learning Outcomes

1. Describe the **structure** of Perceptrons and how it performs classification

2. Understand how Perceptrons are **trained** with the PLA

3. Understand how to **compose** multiple Perceptrons into a Neural Network

4. Describe how Neural Networks are **trained** with gradient descent and backpropagation

# Week 09A: Lecture Outline

1. Perceptron

2. Activation Functions

3. Gradient Descent

4. Neural Networks

5. Backpropagation

# Please <u>turn on</u> your webcam

Mystery Student

AI
HISTORY

# Perceptron

# Linear Classifiers

- Logistic Regression [W04A]
- Linear SVM [W04B]
- Perceptron

# Perceptron

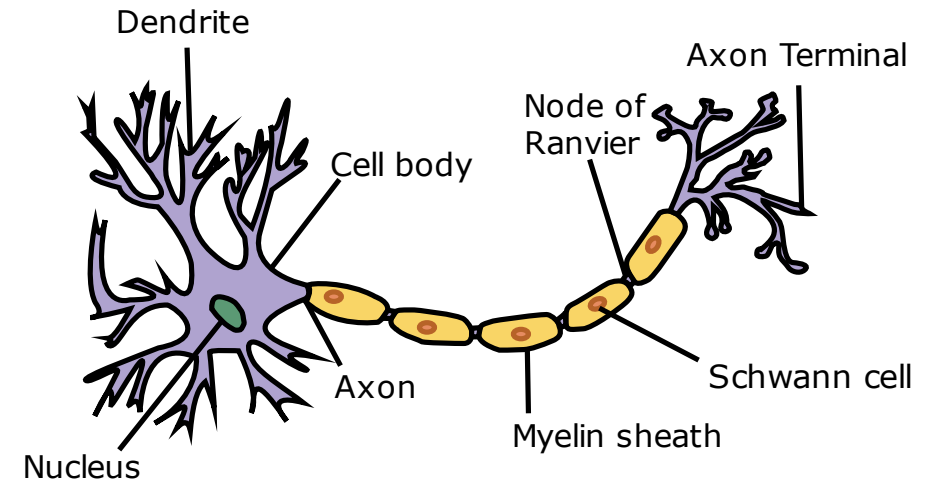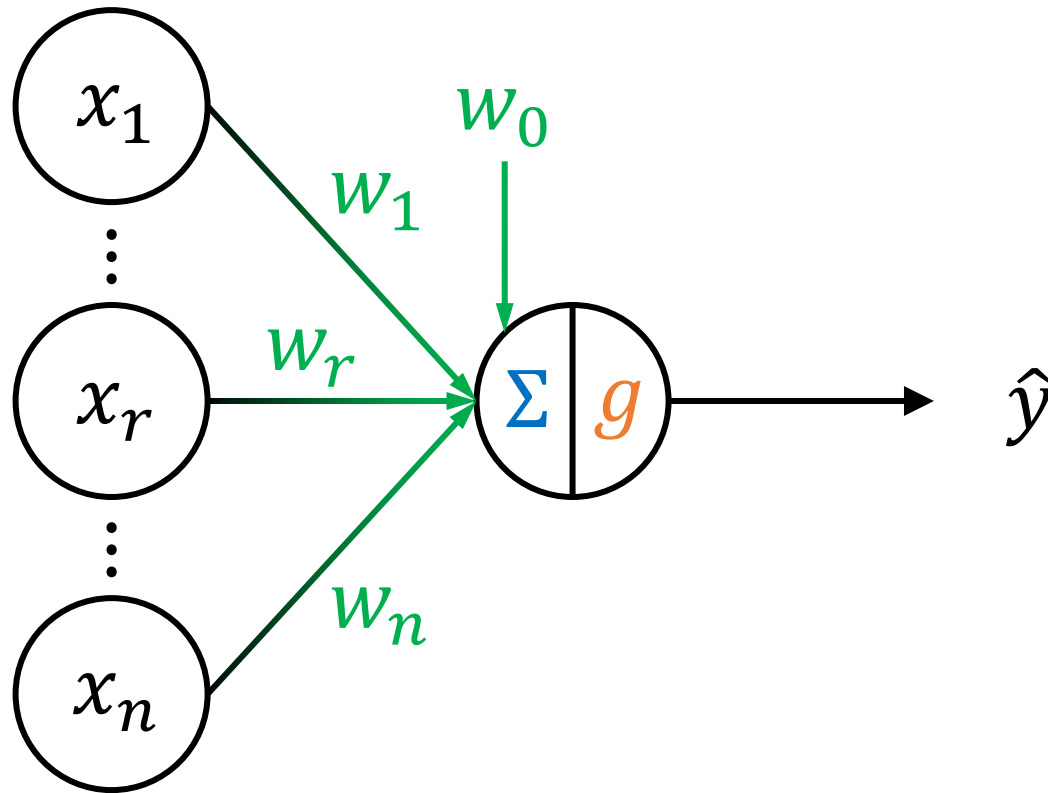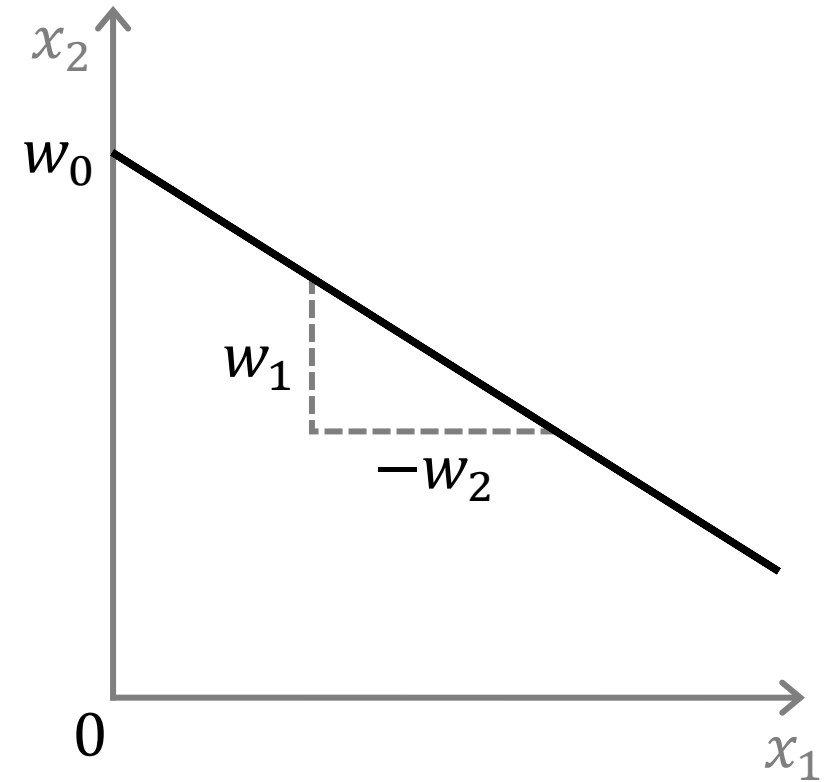- What is a perceptron?
- How to train it?

# Perceptron

$x_1$

$x_r$

$x_n$

$w_0$

$w_1$

$w_r$

$w_n$

$\Sigma \mid g$

$\hat{y}$

Dendrite

Cell body

Axon Terminal

Node of Ranvier

Axon

Schwann cell

Myelin sheath

Nucleus

Diagram credits: Dhp1080 - Own work, CC BY-SA 3.0 via Wikimedia Commons.

# Line Equation

$$x_2 = mx_1 + c$$

$$w_2 x_2 + w_1 x_1 + w_0 = 0$$

$$\sum_{r=0}^{n} w_r x_r = 0 \quad x_0 = 1$$

Image credit: https://en.wikipedia.org/wiki/Sign_function

# Linear Classification

$$x_2 = mx_1 + c$$

$$w_2 x_2 + w_1 x_1 + w_0 = 0$$

$$\sum_{r=0}^{n} w_r x_r = 0 \quad x_0 = 1$$

$$\sum_{r=0}^{n} w_r x_r > 0 \qquad \sum_{r=0}^{n} w_r x_r \le 0$$

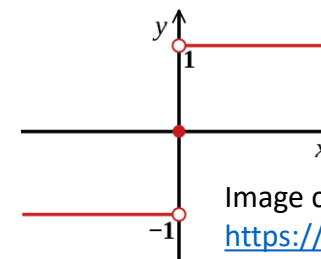$$\hat{y} = \text{sgn}\left(\sum_{r=0}^{n} w_r x_r\right), \text{sgn}(z) = \begin{cases} +1 & z > 0 \\ -1 & z \le 0 \end{cases}$$
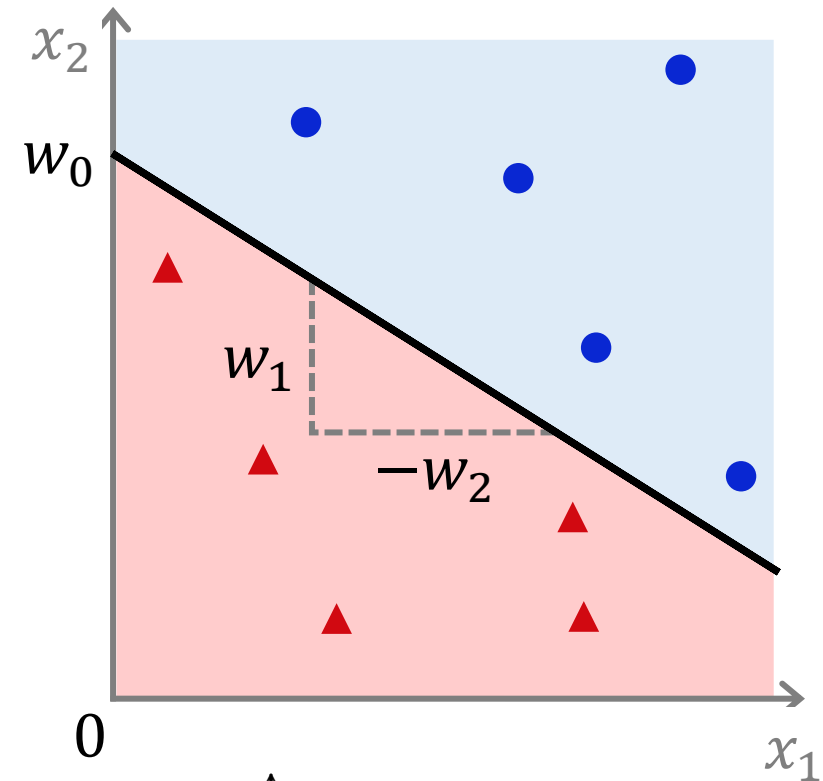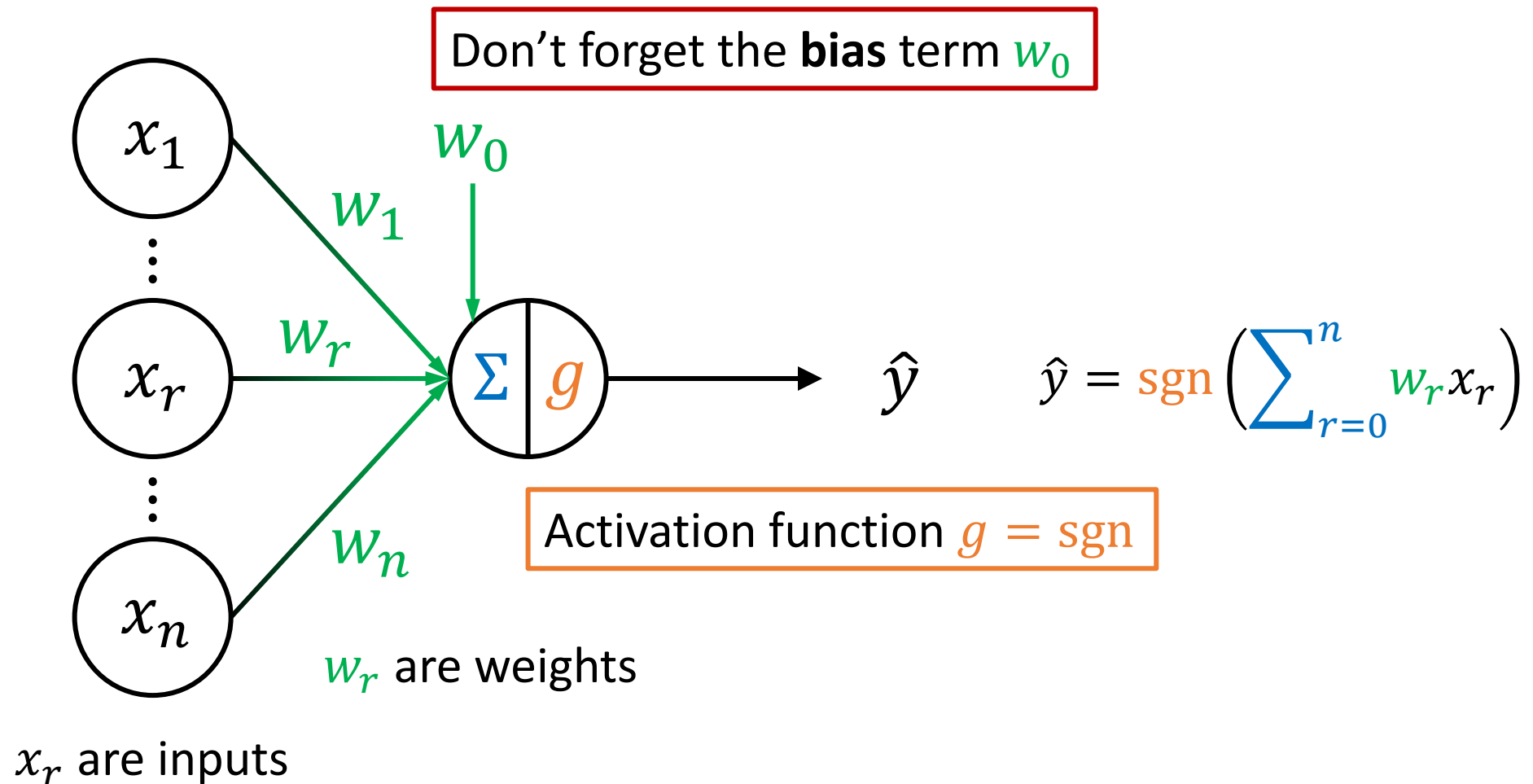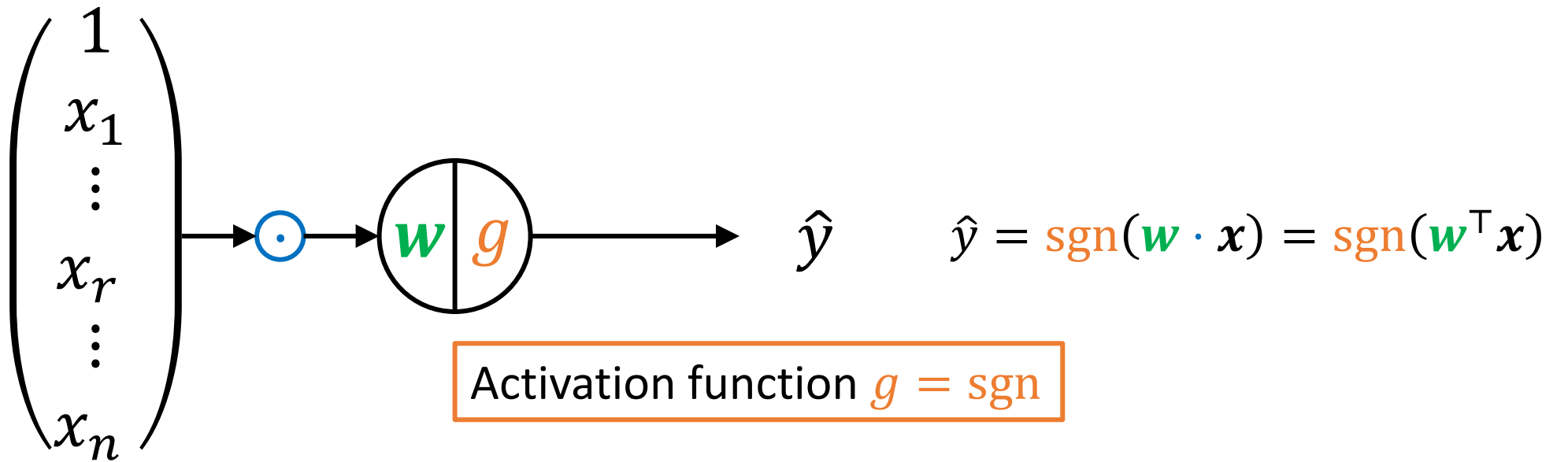


Image credit:
https://en.wikipedia.org/wiki/Sign_function

# Perceptron



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

Activation function $g = \text{sgn}$

$\hat{y}$

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

$w_r$ are weights

$x_r$ are inputs

# Perceptron Classification



$$\hat{y} = \text{sgn}(\boldsymbol{w} \cdot \boldsymbol{x}) = \text{sgn}(\boldsymbol{w}^\top \boldsymbol{x})$$

Activation function $g = \text{sgn}$

# Perceptron Learning Algorithm (PLA)

1. Initialize weights $\boldsymbol{w}$
   - Could be all zero, or random small values

2. For each instance $i$ with features $\boldsymbol{x}^{(i)}$
   - Classify $\hat{y}^{(i)} = \text{sgn}\left(\boldsymbol{w}^\top \boldsymbol{x}^{(i)}\right)$

3. Select one <span style="color:red">mis</span>classified instance
   - Update weights: $\boldsymbol{w} \leftarrow \boldsymbol{w} + \boxed{\Delta \boldsymbol{w}}$

4. Iterate steps 2 to 3 until
   - Convergence (classification error < threshold), or
   - Maximum number of iterations

How to calculate?
- What direction?
- What magnitude?

# Perceptron Weight Update

$$\boldsymbol{w} \leftarrow \boldsymbol{w} + \eta(y - \hat{y})\boldsymbol{x}$$

Old weight

Learning Error

New weight

Learning Rate

# Vector Distances and Similarity



**Cosine Similarity**

$$s = \cos(\theta) = \frac{\hat{\boldsymbol{y}} \cdot \boldsymbol{y}}{\|\hat{\boldsymbol{y}}\|\|\boldsymbol{y}\|}$$

**Angular Distance**

$$\theta = \cos^{-1}(s)$$
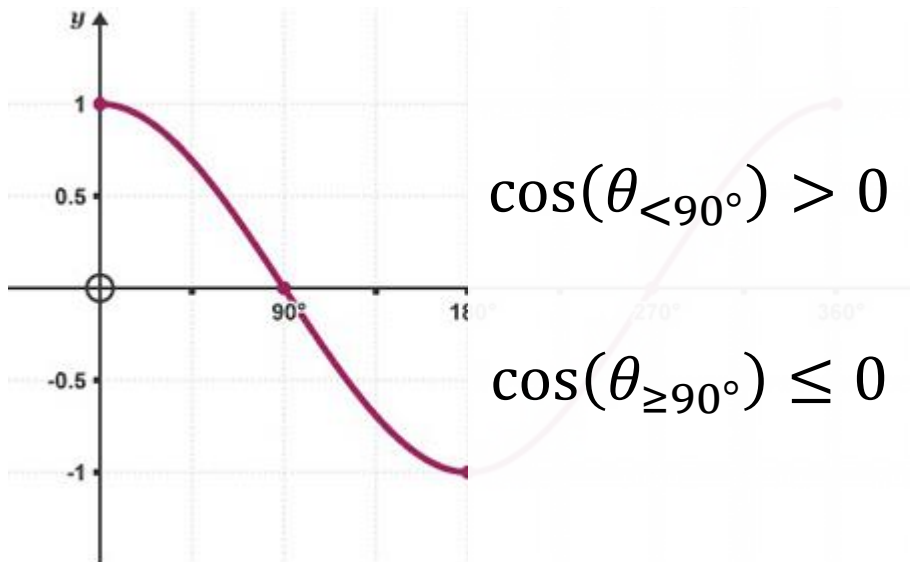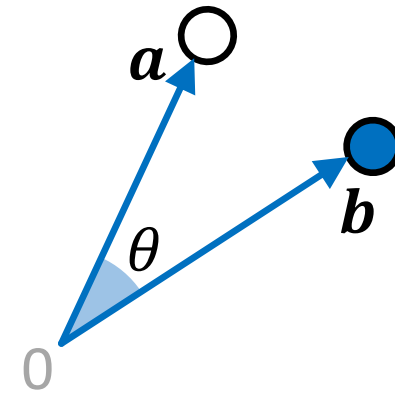
# Vector Distances and Similarity

**Cosine Curve**



$$\cos(\theta_{<90°}) > 0$$

$$\cos(\theta_{\geq 90°}) \leq 0$$

Image credit:

**Cosine Similarity**

$$s = \cos(\theta) = \frac{\boldsymbol{a} \cdot \boldsymbol{b}}{\|\boldsymbol{a}\|\|\boldsymbol{b}\|}$$

**Angular Distance**

$$\theta = \cos^{-1}(s)$$

# Perceptron Weight Update



$$\cos(\theta_{<90°}) > 0$$
$$\cos(\theta_{\geq 90°}) \leq 0$$

**Left column:**

Consider this misclassification:

$$y = +1, \qquad \hat{y} = \text{sgn}(w \cdot x) = -1$$

- $\hat{y} = -1 \Rightarrow w \cdot x \leq 0 \Rightarrow \theta = \cos^{-1}\left(\frac{w}{|w|} \cdot \frac{x}{|x|}\right) \geq 90°$

But we want

- $\hat{y} = +1 \Rightarrow w \cdot x > 0 \Rightarrow \theta < 90°$

- i.e., $w$ to point in a <u>more similar direction</u> as $x$

<u>Adding $x$ to $w$</u> will make a more positive result, i.e.,

$$w' = w + x$$

**Right column:**

Consider this misclassification:

$$y = -1, \qquad \hat{y} = \text{sgn}(w \cdot x) = +1$$

- $\hat{y} = +1 \Rightarrow w \cdot x > 0 \Rightarrow \theta = \cos^{-1}\left(\frac{w}{|w|} \cdot \frac{x}{|x|}\right) < 90°$

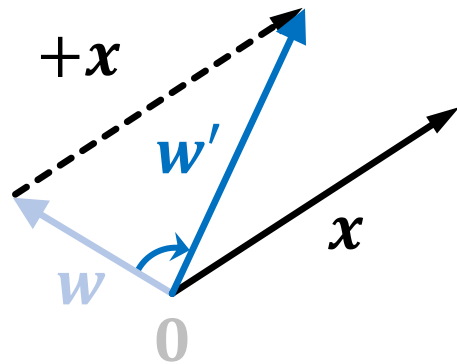But we want

- $\hat{y} = -1 \Rightarrow w \cdot x \leq 0 \Rightarrow \theta > 90°$

- i.e., $w$ to point in a <u>less similar direction</u> as $x$

<u>Negating $x$ from $w$</u> will make a less positive result, i.e.,

$$w' = w - x$$

# Perceptron Weight Update



$$y = +1$$
$$\hat{y} = -1$$

$$y - \hat{y} = +2$$
$$\Delta \boldsymbol{w} = +2\eta \boldsymbol{x}$$

Old weight

Learning Error

$$\boldsymbol{w} \leftarrow \boldsymbol{w} + \eta(y - \hat{y})\boldsymbol{x}$$

New weight

Learning Rate

$$y = -1$$
$$\hat{y} = +1$$

$$y - \hat{y} = -2$$
$$\Delta \boldsymbol{w} = -2\eta \boldsymbol{x}$$

# Perceptron Learning Algorithm
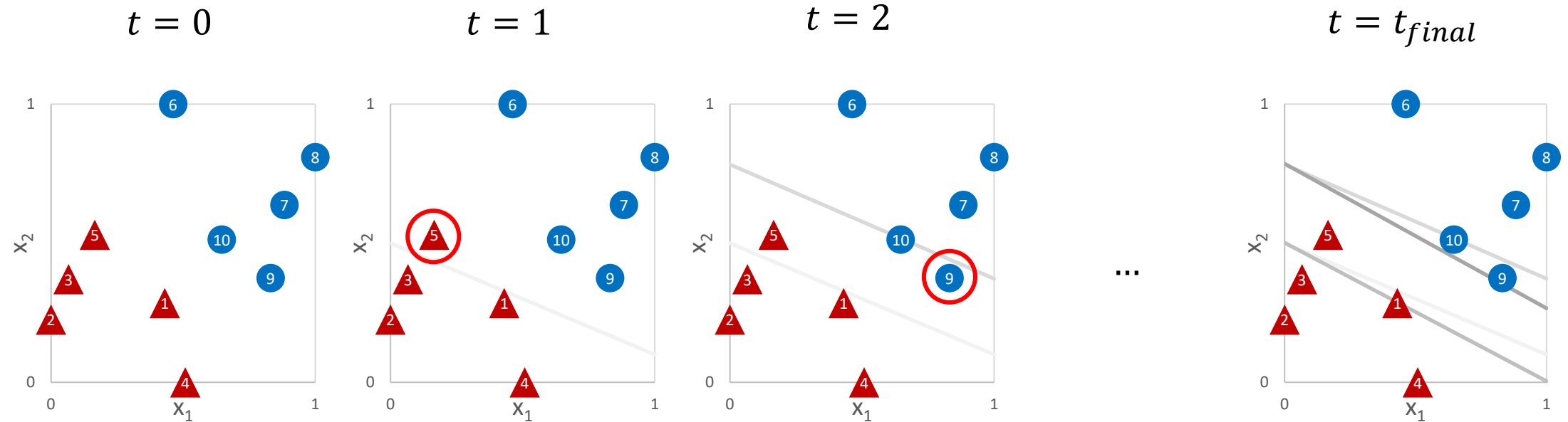
1. Initialize weights $\boldsymbol{w}$
   - Could be all zero, or random small values

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 <span style="color:red">mis</span>classified instance
   - Update weights: $\boldsymbol{w} \leftarrow \boldsymbol{w} + \eta(y - \hat{y})\boldsymbol{x}$

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

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

4. Iterate steps 2 to 3 until
   - Convergence (classification error < threshold), or
   - Maximum number of iterations

# Perceptron Learning Algorithm in Action

# 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:// | tamas.xyz/perceptron-demo/app/ |
|----------|--------------------------------|

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

# What are the differences?
# Perceptron vs. Linear SVM

In Slack #general
1. Write to thread to suggest feature
2. Emote (👍 :+1:) to vote for feature

**Perceptron**                    **Linear Support Vector Machine (SVM)**

# What are the differences?
# Perceptron vs. Linear SVM

In Slack #general
1. Write to thread to suggest feature
2. Emote (👍 :+1:) to vote for feature
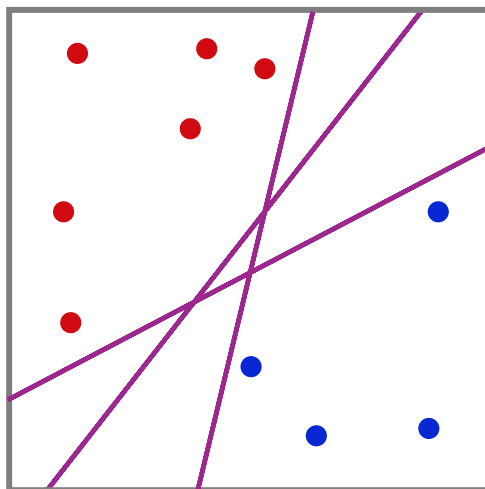
**Perceptron**

**Linear Support Vector Machine (SVM)**



$\xi^{(2)}$

$\xi^{(1)}$

- Can select any model to linear => not robust
  (learns different weights for different initializations)
- Cannot converge on non-linearly separable data

- Perceptron of "optimal stability"
- Maximizes margin
- Soft-margin: allows soft error => can learn from non-linearity separable data

# Extending the Perceptron

- Other **activation functions**
  - Differentiable ones!
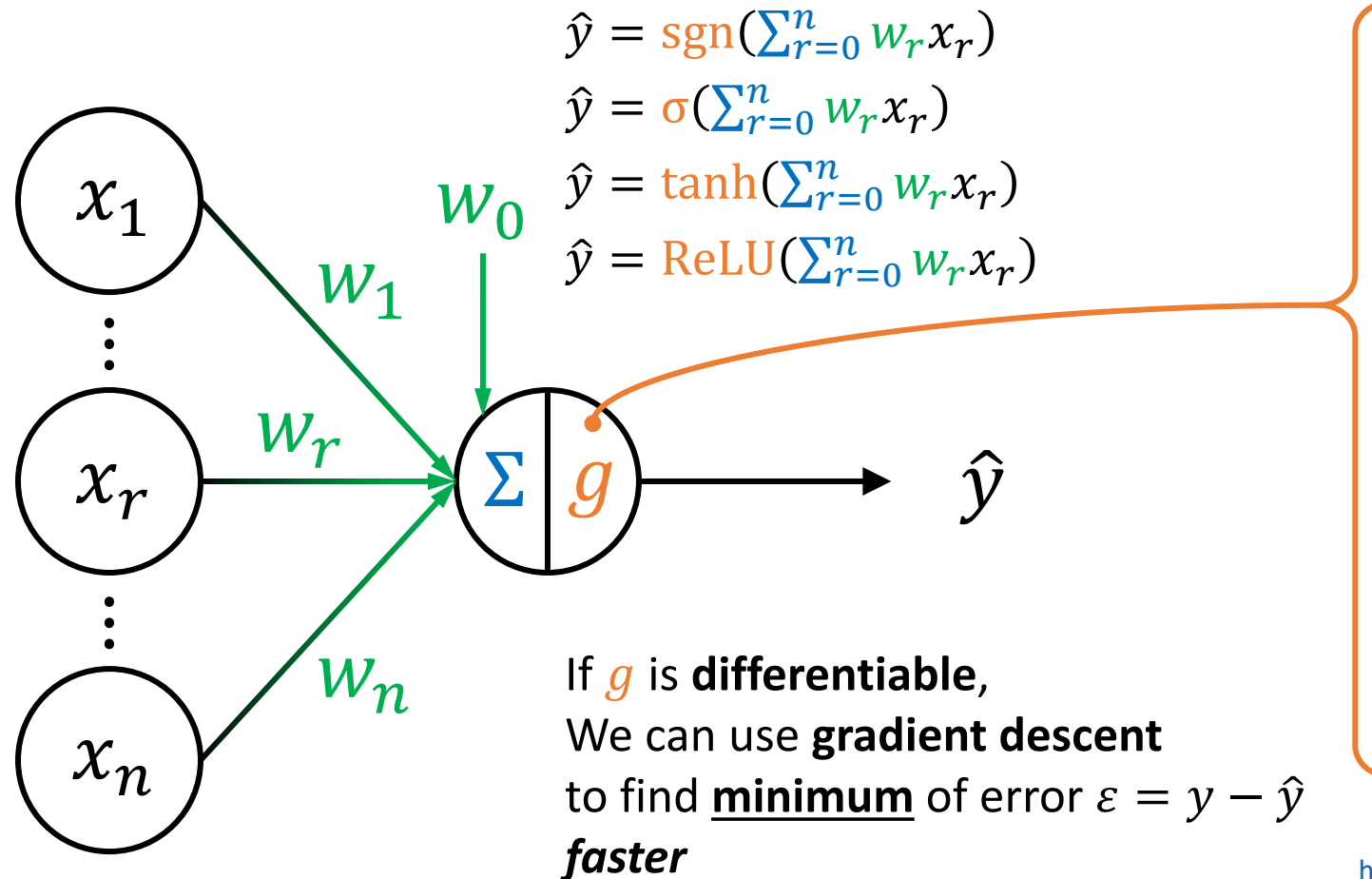- **Multiple** perceptrons
  - Multi-Layer Perceptron (MLP)

  Feed-forward
  Neural Network
- **Non-linear** classifiers
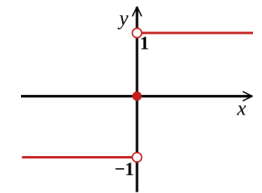
# Activation Functions

# Differentiable Activation Functions

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

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

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

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

$w_0$

$x_1$

$w_1$

$x_r$

$w_r$

$x_n$

$w_n$

$\Sigma$ | $g$

$\hat{y}$

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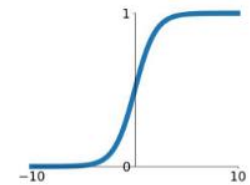*faster*

**Step**

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

**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-gbjlz0zrz8lb.png

# Gradient Descent

## Optimization Goal:
# Iteratively find $w_r$ with minimum error $\varepsilon$



**"Weight space"**
To search for best parameter (weight) values

Iterative **steps** in direction towards (local) minimum

Credits: Alykhan Tejani's Medium Post

# Learning Rate $\eta$



These graphs are also in the "weight space".

# Perceptron Weight Update

Old
weight

Learning
Error

$$\boldsymbol{w} \leftarrow \boldsymbol{w} + \eta(y - \hat{y})\boldsymbol{x}$$

New
weight

Learning
Rate

# Gradient Descent Weight Update

MSE error

$$\varepsilon = \frac{1}{2}(\hat{y} - y)^2$$

Old weight    Learning Error

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

New weight    Learning Rate

Gradient of error

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

# Chain Rule

Consider composite function

$$g(x) = g\big(f(x)\big)$$

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

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

**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)

# Chain Rule

Consider composite function

$$h(x) = h\left(g(f(x))\right)$$

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

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

Optimization Goal:

# Iteratively find $w_r$ with minimum error $\varepsilon$

$$\hat{y} = g(f(w_r)), \quad f(w_r) = \sum_{r=0}^{n} w_r x_r$$

$$\varepsilon = \frac{1}{2}(\hat{y} - y)^2 \quad \text{MSE error}$$

**"Weight space"**
To search for best parameter (weight) values

$$\frac{d\varepsilon}{dw_r} = \frac{d\varepsilon}{d\hat{y}}\frac{d\hat{y}}{dw_r} = (\hat{y} - y)\left(\frac{dg}{df}\frac{df}{dw_r}\right)$$

$$\frac{dg}{df} = g'(f) \qquad \frac{df}{dw_r} = \frac{d(\sum_{r=0}^{n} w_r x_r)}{dw_r}$$

Weight Update

$$\boxed{w_r \leftarrow w_r - \eta\frac{d\varepsilon}{dw_r} = w_r - \eta(\hat{y} - y)g'(f)x_r}$$

$$= \frac{d(w_r x_r)}{dw_r} + \sum_{\rho \neq r}^{n}\left(\frac{d(w_\rho x_\rho)}{dw_r}\right)$$

$$= x_r$$

# What is the derivative of sigmoid function $\sigma$?

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

- Rewrite as compound function

  - $\sigma(z) = \frac{1}{1+z}, \ z(x) = e^{-x}$

- Using chain rule

  - $\sigma'(x) = \frac{dg}{dx} = \frac{dg}{dz}\frac{dz}{dx} = \frac{-1}{(1+z)^2}(-e^{-x}) = \frac{e^{-x}}{(1+e^{-x})^2} = \sigma(x)\frac{e^{-x}}{1+e^{-x}}$

- Notice: $1 - \sigma(x) = 1 - \frac{1}{1+e^{-x}} = \frac{1+e^{-x}-1}{1+e^{-x}} = \frac{e^{-x}}{1+e^{-x}}$

- Substituting back:

  - $\sigma'(x) = \sigma(x)\big(1 - \sigma(x)\big)$

Optimization Goal:

# Iteratively find $w_r$ with minimum error $\varepsilon$

$$\hat{y} = g\big(f(w_r)\big), \; f(w_r) = \sum_{r=0}^{n} w_r x_r$$

$$\varepsilon = \frac{1}{2}(\hat{y} - y)^2 \quad \text{MSE error}$$

$$\frac{d\varepsilon}{dw_r} = \frac{d\varepsilon}{d\hat{y}} \frac{d\hat{y}}{dw_r} = (\hat{y} - y)\left(\frac{dg}{df} \frac{df}{dw_r}\right)$$

Weight
Update

$$\frac{dg}{df} = \sigma'(f) = \hat{y}(1 - \hat{y}) \qquad\qquad \frac{df}{dw_r} = x_r$$

$$w_r \leftarrow w_r - \eta \frac{d\varepsilon}{dw_r} = w_r - \eta(\hat{y} - y)\hat{y}(1 - \hat{y})x_r$$

Optimization Goal:
# Iteratively find $\boldsymbol{w}$ with minimum error $\varepsilon$

$$\hat{y} = g\big(f(\boldsymbol{w})\big), \qquad f(\boldsymbol{w}) = \boldsymbol{w} \cdot \boldsymbol{x}$$

$$\varepsilon = \frac{1}{2}(\hat{y} - y)^2 \quad \text{MSE error}$$

Gradient $\quad \nabla_{\boldsymbol{w}}\varepsilon = \dfrac{d\varepsilon}{d\boldsymbol{w}} = \dfrac{d\varepsilon}{d\hat{y}}\dfrac{d\hat{y}}{d\boldsymbol{w}} = (\hat{y} - y)\left(\dfrac{dg}{df}\dfrac{df}{d\boldsymbol{w}}\right)$

Weight
Update
$$\frac{dg}{df} = \sigma'(f) = \hat{y}(1 - \hat{y}) \qquad \nabla_{\boldsymbol{w}}f = \frac{df}{d\boldsymbol{w}} = \frac{d(\boldsymbol{w} \cdot \boldsymbol{x})}{d\boldsymbol{w}}$$

$$\boxed{\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta\frac{d\varepsilon}{d\boldsymbol{w}} = \boldsymbol{w} - \eta(\hat{y} - y)\hat{y}(1 - \hat{y})\boldsymbol{x}}$$

$$= \begin{pmatrix} \vdots \\ d(w_r x_r)/dw_r \\ \vdots \end{pmatrix} = \begin{pmatrix} \vdots \\ x_r \\ \vdots \end{pmatrix}$$

$$= \boldsymbol{x}$$

Questions!

# Extending the Perceptron

- Other **activation functions**
  - Differentiable ones!
- **Multiple** perceptrons
  - Multi-Layer Perceptron (MLP)    } Feed-forward Neural Network
- **Non-linear** classifiers

# Neural Networks

# W09 Pre-Lecture Task (due before next Mon)

**Watch**

1. [But what is a neural network? | Chapter 1, Deep learning](#) (~20 min) by [3Blue1Brown](#)

2. [The Nervous System, Part 1: Crash Course A&P #8](#) (~10 min) by [CrashCourse](#)

**Discuss**

1. Reflect on how **artificial** neural networks are different from **human** neural networks.

2. Identify **one** point (no need to write several).

3. Post a 1–2 sentence answer to the topic in your tutorial group: #tg-xx

A difference is that artificial neural networks compute layer by layer, whereas the human neural network can process/compute information asynchronously

**Difference**: Artificial neural networks tend to function well(or even better) within its own domain where data and outputs are available for learning while human neural networks have long displayed the ability to automatically find and understand new features overtime regardless of any context.

human neural networks are non deterministic

An interesting difference is that a human neural network is much more efficient in terms of power consumption, with a human brain operating on about 20 watts, while an artificial neural network (running on a modern GPU) may take up about 200~300 watts.

The biggest difference for me is that ANN relies on and can be described completely with mathematics. A biological neural network however works by sending electrical impulses / frequencies to each other as elaborated in Part 2 of the CrashCourse series, and would appear to be much harder to model.
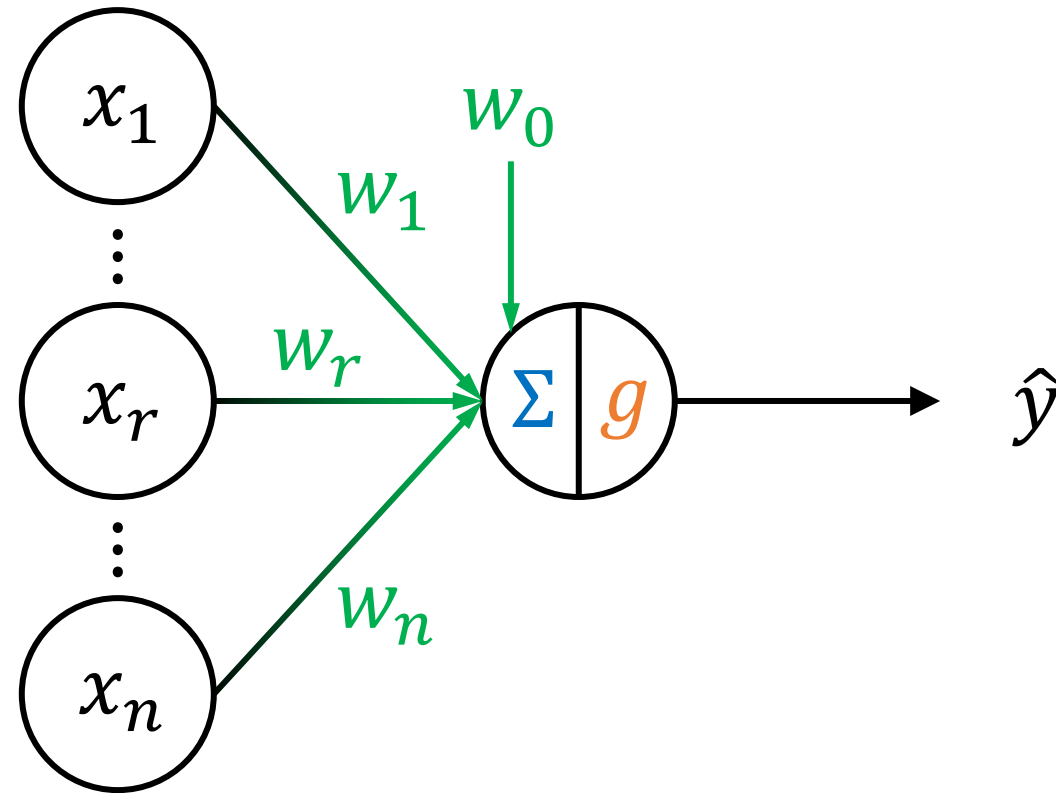
The artificial neural network is mainly uni directional where the transmissions mainly traverse from one layer to another later, from input to output layers, whereas for human nervous system, the transmissions are bi directional and they can traverse from central nervous system to peripheral nervous system that can transmit signals back (edited)

Human neural networks contain different types of neurons that are highly specialized and are of various shapes and sizes. Artificial neural networks contain largely the same types of neurons that are used to store weights with the main difference being their activation functions.
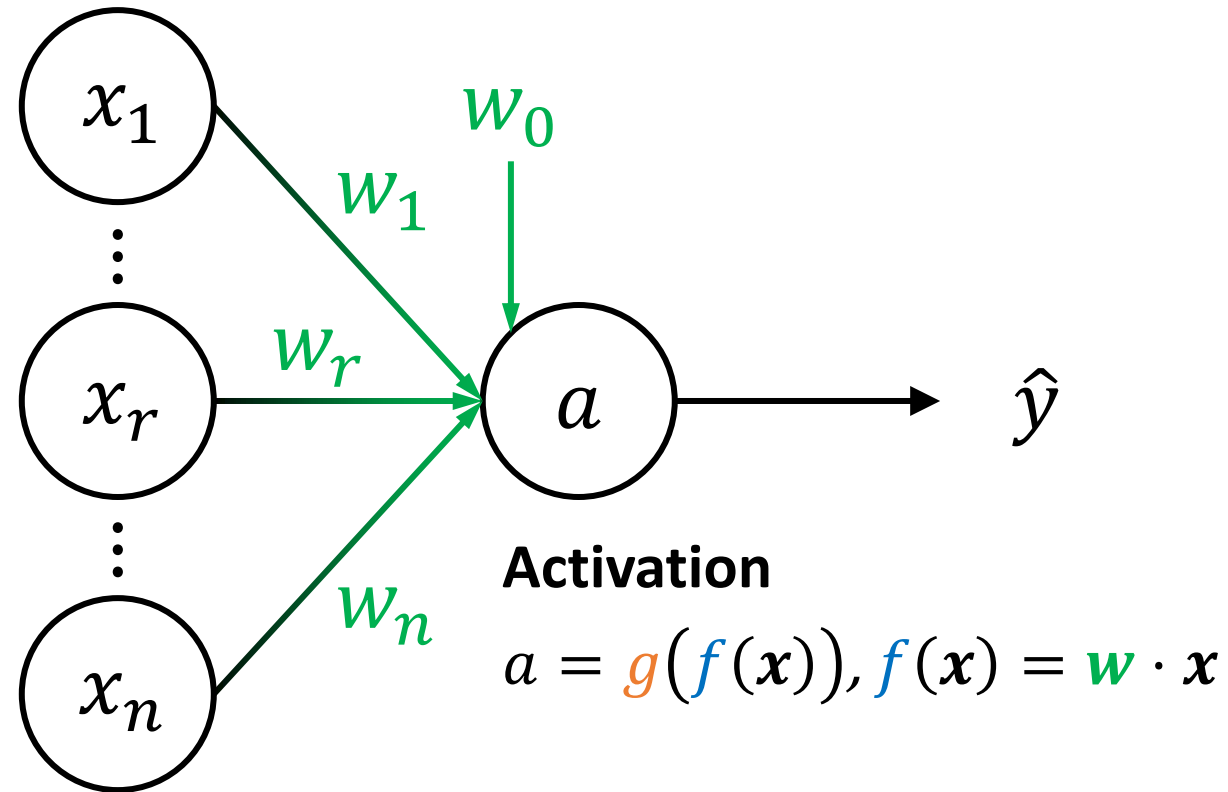
# Neural Networks

- What are Neural Networks (Multi-Layer Perceptron)?
- How to train?

# Single-Layer Perceptron

# Single-Layer Perceptron



**Activation**

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

# Multi-Layer Perceptron (Neural Network)



2-Layer
MLP

Input Layer          Hidden Layer          Output Layer

# Multi-Layer Perceptron (Neural Network)



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

# Multi-Layer Perceptron (Neural Network)



Layer $l-1$      Layer $l$      Layer $l+1$
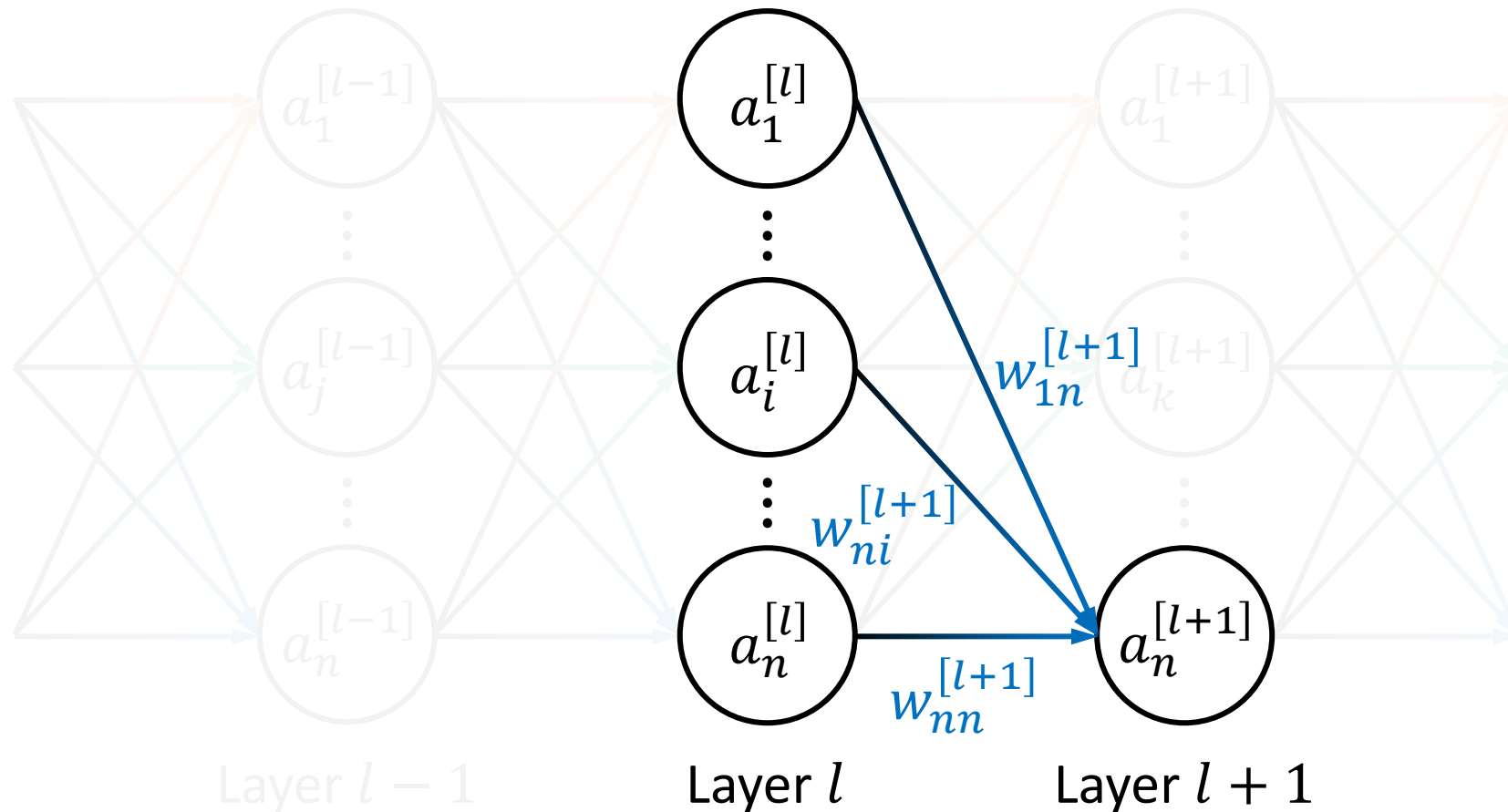
# Multi-Layer Perceptron (Neural Network)

# Multi-Layer Perceptron (Neural Network)

# Fitting **non-linear function** with MLP
## What model <u>weights</u> can model $\hat{y} = |x - 1|$?

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

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

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

$w_{10}^{[1]}$

$w_{10}^{[2]}$

$x_1^{[0]}$  $w_{11}^{[1]}$  $a_1^{[1]}$  $w_{11}^{[2]}$  $\hat{y}_1^{[2]}$

$g^{[2]}$

$w_{20}^{[1]}$

$w_{21}^{[1]}$

$w_{12}^{[2]}$

$a_2^{[1]}$

$g^{[1]}$

**Bonus Question:**
What <u>activation function(s)</u> $g$
should you use for each layer?

In Slack #general
1. <u>Write</u> to thread to suggest feature
2. <u>Emote</u> (👍 :+1:) to vote for feature

# Fitting **non-linear function** with MLP
## What model <u>weights</u> can model $\hat{y} = |x - 1|$?

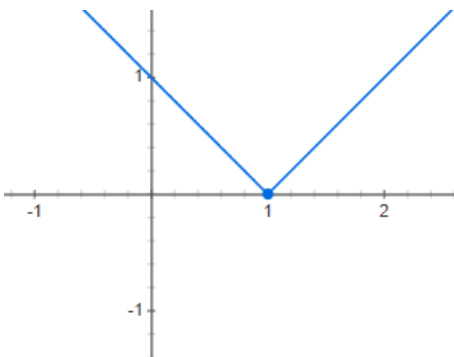$$W^{[1]} = \begin{pmatrix} w_{10}^{[1]} & w_{20}^{[1]} \\ w_{11}^{[1]} & w_{21}^{[1]} \end{pmatrix} = \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix}$$

$$x^{[0]} = \begin{pmatrix} 1 \\ x_1^{[0]} \end{pmatrix}$$

$$a^{[1]} = g^{[1]}\left(\left(W^{[1]}\right)^{\top} x^{[0]}\right)$$

$$= \begin{pmatrix} \text{ReLU}\left(x_1^{[0]} - 1\right) \\ \text{ReLU}\left(1 - x_1^{[0]}\right) \end{pmatrix}$$

Network diagram nodes: $1$, $1$, $x_1^{[0]}$, $a_1^{[1]}$, $a_2^{[1]}$, $\hat{y}_1^{[2]}$

Edge labels: $-1$, $0$, $+1$, $1$, $+1$, $-1$, $1$

$g^{[2]} = I$

$g^{[1]} = \text{ReLU}$

$\max(0, x)$

$$W^{[2]} = \begin{pmatrix} w_{10}^{[2]} \\ w_{11}^{[2]} \\ w_{12}^{[2]} \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}$$

$$a^{[1]} = \begin{pmatrix} 1 \\ \text{ReLU}\left(x_1^{[0]} - 1\right) \\ \text{ReLU}\left(1 - x_1^{[0]}\right) \end{pmatrix}$$

$$\hat{y}^{[2]} = g^{[2]}\left(\left(W^{[2]}\right)^{\top} a^{[1]}\right)$$

$$= 0$$

$$+ \text{ReLU}\left(x_1^{[0]} - 1\right)$$

$$+ \text{ReLU}\left(1 - x_1^{[0]}\right)$$

Questions!

# 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:// | www.desmos.com/calculator/f05dzu2rdv

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

# Universal Approximation Theorem

- Each neuron contributes a **piecewise function**

- Many piecewise functions can **approximate a curve**



$a_1^{[1]} = 0$

$a_2^{[1]} = .2$

$W_{1,\{1,2\}}^{[2]} = \mp 1.5$

$a_3^{[1]} = .2$

$a_4^{[1]} = .4$

$W_{1,\{3,4\}}^{[2]} = \mp 0.7$

$x_1$

$a_5^{[1]} = .4$

$a_6^{[1]} = .6$

$W_{1,\{5,6\}}^{[2]} = \pm 1.2$

$a_1^{[2]}$

$a_7^{[1]} = .6$

$a_8^{[1]} = .8$

$W_{1,\{7,8\}}^{[2]} = \pm 1.5$

$a_9^{[1]} = .8$

$a_{10}^{[1]} = 1$

$W_{1,\{9,10\}}^{[2]} = \pm 1.4$

Adapted based on Michael A. Nielsen "NN and DL" 2015, Determination Press, CC By-NC 3.0

# Weighted Sum

**Summation Series** = Scalar

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

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

**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 Vector Multiplication** = Scalar

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

**Transposed Matrix Multiplication** = Vector

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

# Backpropagation "backprop"

# Backpropagation

Backpropagation **efficiently** computes the **gradient** by

- Avoiding duplicate calculations

- Not computing unnecessary intermediate values,

- Computing the **gradient** of *each* **layer**

Specifically, the gradient of the weighted input of each layer is calculated from back $[l + 1]$ to front $[l]$):

$$\hat{y}'\left(W^{[l]}\right) = a^{[l-1]}\delta^{[l]} \quad \Big| \quad \delta^{[l]} = g'^{[l]}\left(f^{[l]}\right)W^{[l+1]}\delta^{[l+1]}$$

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

Fitting **non-linear function** with MLP
What model weights can model $\hat{y} = |x - 1|$?

$$\boldsymbol{a}^{[1]} = g^{[1]}\left(\left(\boldsymbol{W}^{[1]}\right)^{\top}\boldsymbol{x}^{[0]}\right)$$

$$\hat{y}^{[2]} = g^{[2]}\left(\left(\boldsymbol{W}^{[2]}\right)^{\top}\boldsymbol{a}^{[1]}\right)$$

# Layer Activation

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

Single-Layer
Perceptron

Layer $l$
Activation
Function

Layer $l$
Weights

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

Layer $l$ in
Neural Network

Layer $l$
Activations

Layer $l-1$
Activations

# Neural Network

$$\boldsymbol{W}^{[l-1]} \quad \boldsymbol{a}^{[l-1]} \quad \boldsymbol{W}^{[l]} \quad \boldsymbol{a}^{[l]} \quad \boldsymbol{W}^{[l+1]} \quad \boldsymbol{a}^{[l+1]} \quad \boldsymbol{W}^{[l+2]}$$

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

# Forward Propagation

$$\boldsymbol{a}^{[l]} \equiv g^{[l]}(f^{[l]}), \; f^{[l]} \equiv \left(W^{[l]}\right)^{\top} \boldsymbol{a}^{[l]}$$

$$g^{[1]}\left(f^{[1]}(x^{[0]})\right) = \boldsymbol{a}^{[1]} \qquad g^{[l]}\left(f^{[l]}(\boldsymbol{a}^{[l-1]})\right) = \boldsymbol{a}^{[l]} \qquad g^{[L]}\left(f^{[L]}(\boldsymbol{a}^{[L-1]})\right) = \boldsymbol{a}^{[L]}$$

$$g^{[1]}\left(\left(W^{[1]}\right)^{\top} x^{[0]}\right) = \boldsymbol{a}^{[1]} \qquad g^{[l]}\left(\left(W^{[l]}\right)^{\top} \boldsymbol{a}^{[l-1]}\right) = \boldsymbol{a}^{[l]} \qquad g^{[L]}\left(\left(W^{[L]}\right)^{\top} \boldsymbol{a}^{[L-1]}\right) = \boldsymbol{a}^{[L]}$$



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

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

# Gradient Descent Weight Update

MSE error

$$\varepsilon = \frac{1}{2}(\hat{y} - y)^2$$

Old weight   Learning Error

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

New weight   Learning Rate

Gradient of error

$$\nabla \varepsilon = \frac{\partial \varepsilon}{\partial \boldsymbol{W}} = \frac{\partial \varepsilon}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \boldsymbol{W}}$$

Even more Chain Rule:
# Gradient of Neural Network

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

**Gradient** relative to $\boldsymbol{W}$

$$\hat{y}'(\boldsymbol{W}^{[L-1]}) = \frac{dg^{[L]}}{d\boldsymbol{W}^{[L-1]}} = \frac{df^{[L]}}{d\boldsymbol{W}^{[L-1]}}\boxed{\frac{dg^{[L]}}{df^{[L]}}} \quad \delta^{[L-1]}$$

$$\hat{y}'(\boldsymbol{W}^{[l+1]}) = \frac{dg^{[L]}}{d\boldsymbol{W}^{[l+1]}} = \qquad \frac{df^{[l+1]}}{d\boldsymbol{W}^{[l+1]}}\boxed{\frac{dg^{[l+1]}}{df^{[l+1]}}\cdots\frac{df^{[L]}}{dg^{[L-1]}}\frac{dg^{[L]}}{df^{[L]}}} \quad \delta^{[l+1]}$$

$$\hat{y}'(\boldsymbol{W}^{[l]}) = \frac{dg^{[L]}}{d\boldsymbol{W}^{[l]}} = \frac{df^{[l]}}{d\boldsymbol{W}^{[l]}}\frac{dg^{[l]}}{df^{[l]}}\underbrace{\frac{df^{[l+1]}}{dg^{[l]}}\frac{dg^{[l+1]}}{df^{[l+1]}}\cdots\frac{df^{[L]}}{dg^{[L-1]}}\frac{dg^{[L]}}{df^{[L]}}}$$

**Recursive**

$$\hat{y}'(\boldsymbol{W}^{[l]}) = \frac{dg^{[L]}}{d\boldsymbol{W}^{[l]}} = \frac{df^{[l]}}{d\boldsymbol{W}^{[l]}}\frac{dg^{[l]}}{df^{[l]}}\frac{df^{[l+1]}}{dg^{[l]}}\delta^{[l+1]}$$

$$\hat{y}'(\boldsymbol{W}^{[1]}) = \frac{dg^{[L]}}{d\boldsymbol{W}^{[1]}} = \frac{df^{[1]}}{d\boldsymbol{W}^{[1]}}\frac{dg^{[1]}}{df^{[1]}}\cdots\frac{dg^{[l]}}{df^{[l]}}\frac{df^{[l+1]}}{dg^{[l]}}\frac{dg^{[l+1]}}{df^{[l+1]}}\cdots\frac{df^{[L]}}{dg^{[L-1]}}\frac{dg^{[L]}}{df^{[L]}}$$

Even more Chain Rule:
# Gradient of Neural Network

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

**Gradient** relative to $\boldsymbol{W}^{[l]}$

$$\hat{y}'(\boldsymbol{W}^{[l]}) = \frac{dg^{[L]}}{dW^{[l]}} = \frac{df^{[l]}}{dW^{[l]}}\frac{dg^{[l]}}{df^{[l]}}\frac{df^{[l+1]}}{dg^{[l]}}\delta^{[l+1]}$$

Reference

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

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

$$\frac{df^{[l]}}{dW^{[l]}} = \boldsymbol{a}^{[l-1]}$$

$$\frac{dg^{[l]}}{df^{[l]}} = g'^{[l]}(f^{[l]})$$

$$\frac{df^{[l+1]}}{dg^{[l]}} = \frac{df^{[l+1]}}{d\boldsymbol{a}^{[l]}} = W^{[l+1]}$$

$$\hat{y}'(\boldsymbol{W}^{[l]}) = \boldsymbol{a}^{[l-1]} \quad g'^{[l]}(f^{[l]}) \quad \boldsymbol{W}^{[l+1]} \quad \delta^{[l+1]} \quad \text{---} \quad \delta^{[l]}$$

$$\boxed{\hat{y}'(\boldsymbol{W}^{[l]}) = \boldsymbol{a}^{[l-1]}\delta^{[l]} \quad \bigg| \quad \delta^{[l]} = g'^{[l]}(f^{[l]})\boldsymbol{W}^{[l+1]}\delta^{[l+1]}}$$

# Backward Propagation



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

# 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** (not simple weighted sum)
  - They need to be **differentiable** to be able to calculate their **gradients**
  - Otherwise, backprop cannot update weights accurately

# Wrapping Up
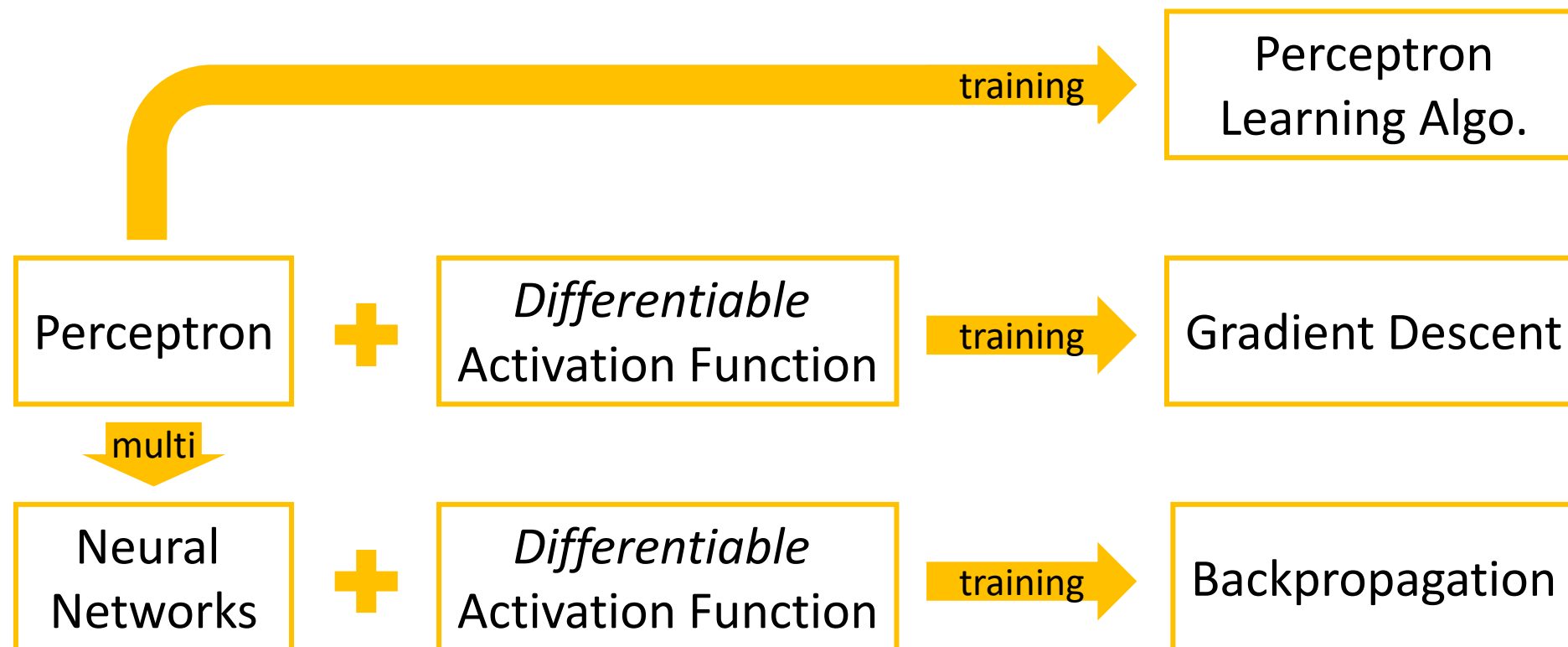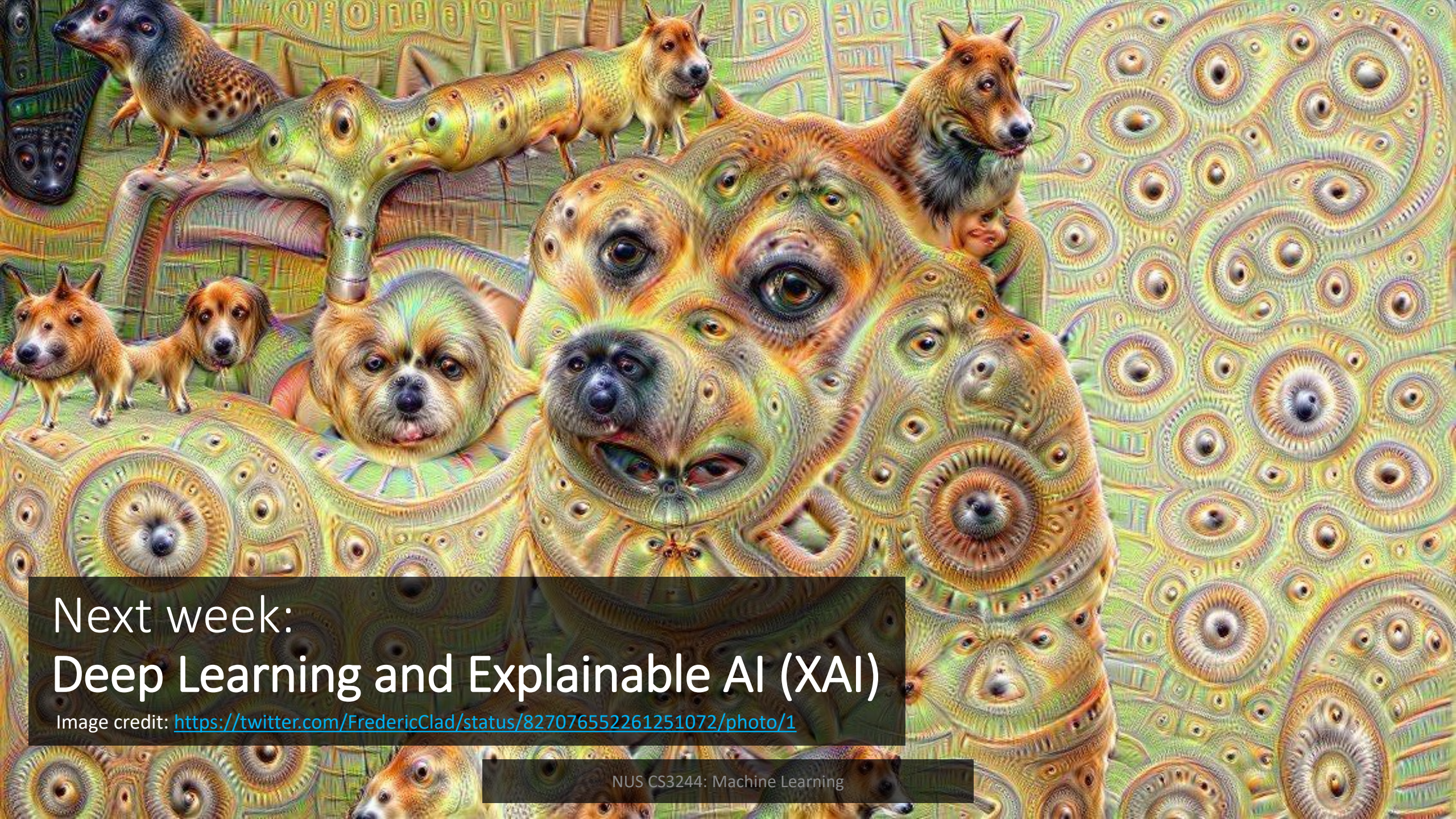
AI
HISTORY

# What did we learn?

| | | | | |
|---|---|---|---|---|
| | | | **training** → | Perceptron Learning Algo. |
| Perceptron | **+** | *Differentiable* Activation Function | **training** → | Gradient Descent |
| ↓ multi | | | | |
| Neural Networks | **+** | *Differentiable* Activation Function | **training** → | Backpropagation |

Next week:
Deep Learning and Explainable AI (XAI)

Image credit: https://twitter.com/FredericClad/status/827076552261251072/photo/1

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