# Perceptron & Neural Networks
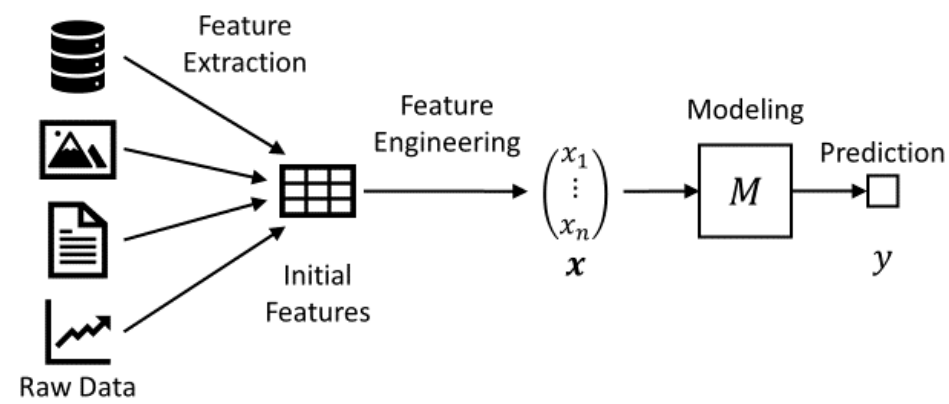
**9**

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

**A**

NUS | Computing
National University of Singapore

## Feature Extraction/Engineering → Modeling



Feature Extraction

Feature Engineering

Modeling

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

$M$

Prediction

$y$

Initial Features
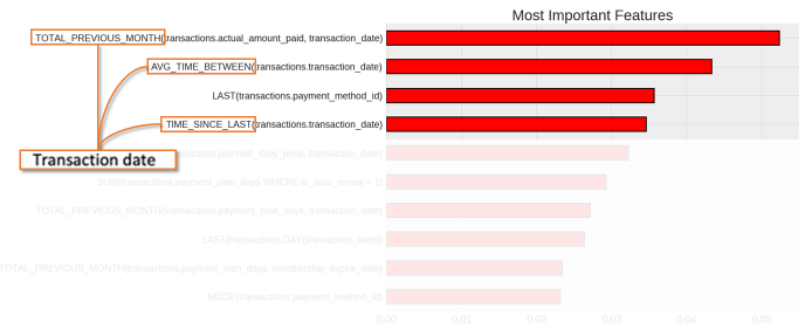
Raw Data

NUS CS3244: Machine Learning          9

## What did we learn?

1. 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 |

2. Describe **issues** when extracting features for various data types

NUS CS3244: Machine Learning          52

## 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)

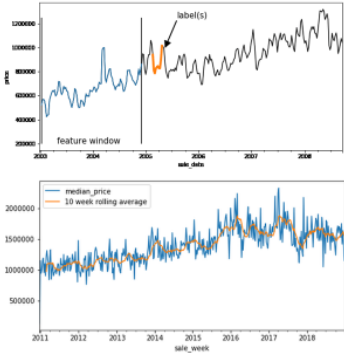label(s)

feature window

median_price
10 week rolling average

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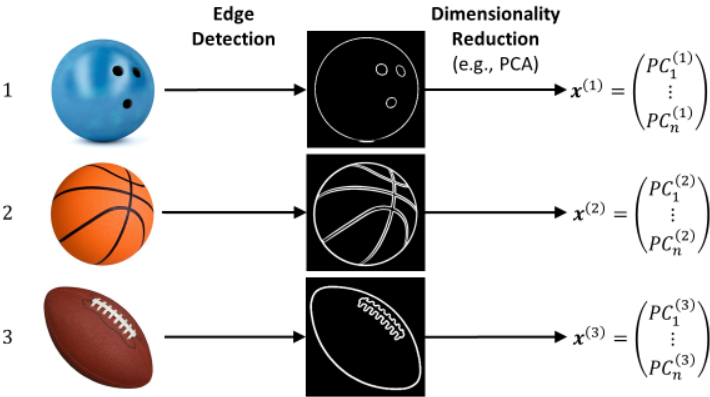
NUS CS3244: Machine Learning    19

## 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

## 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

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

$$x^{(1)} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad x^{(2)} = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

| # | 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 **Perceptron Learning Algorithm**

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

4. Describe how Neural Networks are *trained* with **gradient descent** and **backpropagation** [W09b]

# Week 09A: Lecture Outline

1. Perceptron

2. Perceptron Learning Algorithm (PLA)

3. Activation Functions

4. Gradient Descent

5. Neural Networks

6. Backpropagation [W09B]

AI HISTORY

# Perceptron

# Linear Classifiers

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

Line is a **Decision Boundary**

$x_2$

$x_1$

0

# Perceptron

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

# Perceptron



$x_1$

$w_0$

$w_1$

$\vdots$

$w_r$

$x_r$

$\Sigma \mid g$     $\hat{y}$

$\vdots$

$w_n$

$x_n$



Dendrite

Axon Terminal

Node of
Ranvier

Cell body

Axon

Schwann cell

Nucleus

Myelin sheath

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

$$x_2 = mx_1 + c$$

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

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

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

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



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

# Perceptron



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

Activation function $g = \text{sgn}$

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

$w_r$ 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)

# 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 $\hat{y}^{(i)} = \mathrm{sgn}\left(\boldsymbol{w}^{\top}\boldsymbol{x}^{(i)}\right)$

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

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

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

# Perceptron Weight Update

Old weight

Learning Error

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

New weight

Learning Rate

# Vector Distances and Similarity

**Cosine Similarity**

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

# Vector Distances and Similarity

**Cosine Curve**



$$\cos(\theta_{<90°}) > 0$$

$$\cos(\theta_{\geq 90°}) \leq 0$$

Image credit:
https://www.open.edu/openlearn/ocw/pluginfile.php/947914/mod_ouc
ontent/oucontent/48949/9eaffc43/9f8315d5/mfs_w4_fig4.jpg



**Cosine Similarity**

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

# Perceptron Weight Update



Consider this misclassification:

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

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

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.,

$$w' = w + x$$

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

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

- 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



$+\boldsymbol{x}$

$\boldsymbol{w'}$

$\boldsymbol{w}$

$\boldsymbol{x}$

$\boldsymbol{0}$

$y = +1$

$\hat{y} = -1$

$y - \hat{y} = +2$

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

$(y - \hat{y})\boldsymbol{x}$

$\boldsymbol{w}$

$\boldsymbol{x}$

$\boldsymbol{0}$

Old
weight

Learning
Error

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

New
weight

Learning
Rate

$-\boldsymbol{x}$

$\boldsymbol{w}$

$\boldsymbol{w'}$

$\boldsymbol{x}$

$\boldsymbol{0}$

$y = -1$

$\hat{y} = +1$

$y - \hat{y} = -2$

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

# 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 $\hat{y}^{(i)} = \text{sgn}\left(\boldsymbol{w}^\top \boldsymbol{x}^{(i)}\right)$

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

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

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

# Perceptron Learning Algorithm in Action

## Parameters

| | | | |
|---|---|---|---|
| random seed: | 11 | base width: | 100 |
| separator trail length: | 5 | width variance: | 100 |
| run delay (ms): | 200 | base height | 100 |
| cloud size: | 20 | height variance: | 50 |

apply params    step    run

## Current Classification

| TP | FP | TN | FN | Precision | Recall | F1 |
|---|---|---|---|---|---|---|
| 20 | 0 | 20 | 0 | 1.00 | 1.00 | 0.50 |

$w(x) = -0.49607365285675603 \, x + 217$

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

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

**Perceptron**



**Linear Support Vector Machine (SVM)**



- 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

- Perceptron is a linear classifier
- **Non-linear** classifiers
  - Other **activation functions**
    - Differentiable ones!
  - **Multiple** perceptrons / neurons
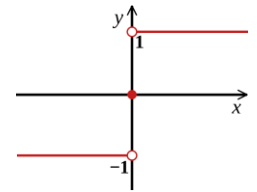    - Multi-Layer Perceptron (MLP)
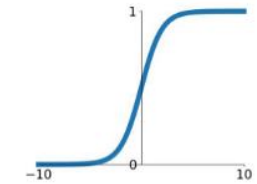
# Activation Functions

# Activation Functions

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

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

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

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

$x_1$

$w_1$

$w_0$

$w_r$

$x_r$

$\Sigma \mid g$

$\hat{y}$

$w_n$

$x_n$

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

**Step**

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

**Sigmoid**

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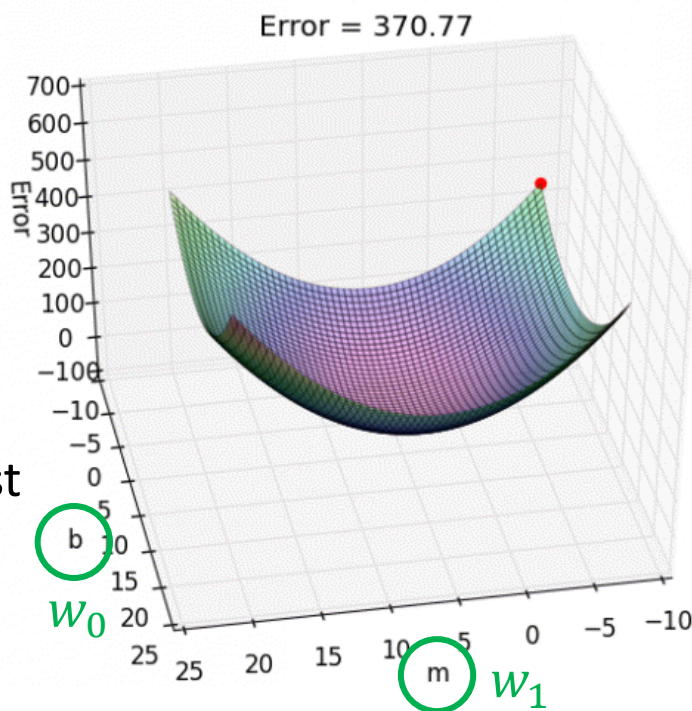
**tanh**

$\tanh(x)$

**ReLU**

$\max(0, x)$

Image Credit:
https://miro.medium.com/max/1400/0*sIJ-gbjlz0zrz8lb.png

# Gradient Descent

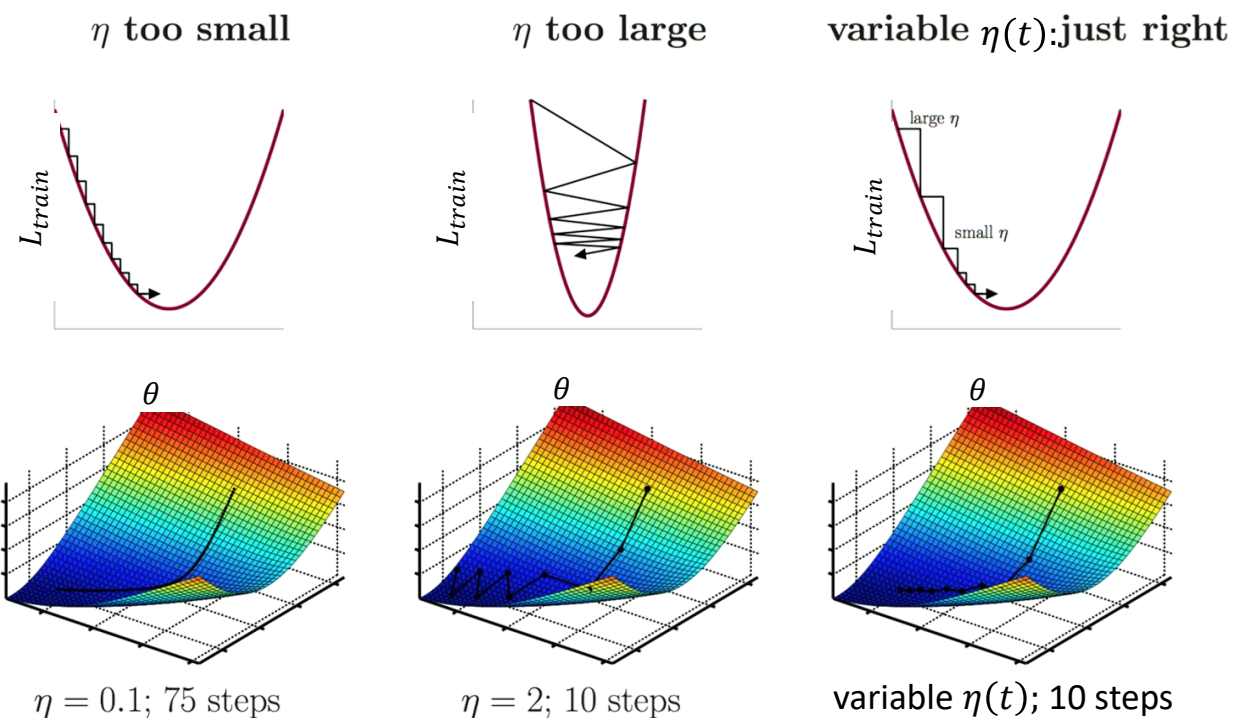## Optimization Goal:
# Iteratively find $w_r$ with minimum error $\varepsilon$



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

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

Credits: Alykhan Tejani's Medium Post

# Learning Rate $\eta$



$\eta$ **too small**

$\eta$ **too large**

**variable** $\eta(t)$:**just right**

$\eta = 0.1$; 75 steps

$\eta = 2$; 10 steps

variable $\eta(t)$; 10 steps

These graphs are also in the "weight space".

# Perceptron Weight Update

Old
weight

Learning
Error

$$\textcolor{green}{\boldsymbol{w}} \leftarrow \textcolor{green}{\boldsymbol{w}} + \textcolor{blue}{\eta} \textcolor{red}{(y - \hat{y})} \boldsymbol{x}$$

New
weight

Learning
Rate

# Gradient Descent Weight Update

Old weight

Direction of fastest error increase

Gradient of error

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

New weight

Learning Rate

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

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

Binary Cross-Entropy error
(for classification)

$$\varepsilon = -y \log \hat{y}$$

Square Error
(for regression)

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

# Chain Rule

## Consider composite function

**Lagrange notation**
Prime $'$ indicates first derivative relative to the function argument.
This can make writing derivatives more concise.
e.g., $y'(w) = dy/dw$

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

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

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

**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 a deeper composite function

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

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

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

# Chain Rule

Multivariate

For single neuron



$$\varepsilon(\boldsymbol{w}) = \varepsilon\Big(g(f(\boldsymbol{w}))\Big)$$

$$\varepsilon = \varepsilon(g), \qquad \hat{y} = g, \qquad g = g(f), \qquad f = f(\boldsymbol{w})$$

$$\nabla_{\boldsymbol{w}}\varepsilon = \frac{d\varepsilon}{d\boldsymbol{w}} = \frac{d\varepsilon}{dg}\frac{dg}{df}\frac{df}{d\boldsymbol{w}}$$

# Gradient of Weighted Sum

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

$$\frac{\partial f}{\partial w_r} = \frac{\partial}{\partial w_r}\left( w_r x_r + \sum_{\rho \neq r} w_\rho x_\rho \right)$$

$$= x_r + 0$$

$$\frac{\partial f}{\partial w_r} = x_r$$

$$f = \boldsymbol{w} \cdot \boldsymbol{x} = \boldsymbol{w}^\top \boldsymbol{x}$$

$$\frac{df}{d\boldsymbol{w}} = \sum_{r=0}^{n} \frac{\partial f}{\partial w_r} \boldsymbol{e}_r = \begin{pmatrix} \partial f/\partial w_0 \\ \vdots \\ \partial f/\partial w_r \\ \vdots \\ \partial f/\partial w_n \end{pmatrix} = \begin{pmatrix} x_0 \\ \vdots \\ x_r \\ \vdots \\ x_n \end{pmatrix}$$

$$\nabla_{\boldsymbol{w}} f = \frac{df}{d\boldsymbol{w}} = \boldsymbol{x}$$

# Chain Rule

Multivariate

For single neuron



$$\varepsilon(\boldsymbol{w}) = \varepsilon\Big(g\big(f(\boldsymbol{w})\big)\Big)$$

$$\varepsilon = \varepsilon(g), \qquad \hat{y} = g, \qquad g = g(f), \qquad f = \boldsymbol{w}^\top \boldsymbol{x}$$

$$\nabla_{\boldsymbol{w}}\varepsilon = \frac{d\varepsilon}{d\boldsymbol{w}} = \frac{d\varepsilon}{dg}\frac{dg}{df}\boldsymbol{x}$$

# Calculating gradient for single neuron

| Error $\varepsilon$ | | $\dfrac{d\varepsilon}{d\hat{y}}$ | Activation $\hat{y} = g(f)$ | | $\dfrac{dg}{df}$ | Weighted Sum $f(\boldsymbol{w})$ | $\dfrac{df}{d\boldsymbol{w}}$ |
|---|---|---|---|---|---|---|---|
| Square Error | $\dfrac{1}{2}(y - \hat{y})^2$ | $-(y - \hat{y})$ | Sigmoid | $\dfrac{1}{1 + e^{-f}}$ | $(1 - g)g$ | $\boldsymbol{w}^\top \boldsymbol{x}$ | $\boldsymbol{x}$ |
| Binary Cross Entropy | $-y \log \hat{y}$ | $-\dfrac{y}{\hat{y}}$ | tanh | $\tanh f$ | $1 - g^2$ | **Advanced:** <br> <u>further reading</u> | |
| | | | ReLU | $\max(0, f)$ | $[f > 0]$ | | |

# Derivative of sigmoid function $\sigma$

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

- Proof

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

  - Rewrite as compound function

    - $\sigma(\chi) = \frac{1}{1+\chi}$, where $\chi(x) = e^{-x}$
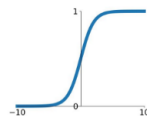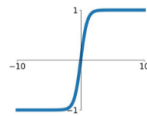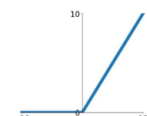
  - Using chain rule

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

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

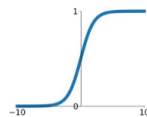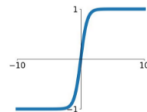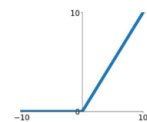  - Substituting back:

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

Example 1

# Calculating gradient for single neuron

| Error $\varepsilon$ | | $\dfrac{d\varepsilon}{d\hat{y}}$ | Activation $\hat{y} = g(f)$ | | $\dfrac{dg}{df}$ | Weighted Sum $f(\boldsymbol{w})$ | $\dfrac{df}{d\boldsymbol{w}}$ |
|---|---|---|---|---|---|---|---|
| Square Error | $\dfrac{1}{2}(y-\hat{y})^2$ | $-(y-\hat{y})$ | Sigmoid | $\dfrac{1}{1+e^{-f}}$ | $(1-g)g$ | $\boldsymbol{w}^\top\boldsymbol{x}$ | $\boldsymbol{x}$ |
| Binary Cross Entropy | $-y\log\hat{y}$ | $-\dfrac{y}{\hat{y}}$ | tanh | $\tanh f$ | $1-g^2$ | | |
| | | | ReLU | $\max(0,f)$ | $[f>0]$ | | |

$$\varepsilon = -y\log\hat{y} = -y\log g = -y\log(\max(0,f)) = -y\log(\max(0,\boldsymbol{w}^\top\boldsymbol{x}))$$

$$\nabla_{\boldsymbol{w}}\varepsilon = \frac{d\varepsilon}{d\hat{y}}\frac{dg}{df}\frac{df}{d\boldsymbol{w}} = -\frac{y}{\hat{y}}[f>0]\boldsymbol{x} = -\frac{y}{\max(0,\boldsymbol{w}^\top\boldsymbol{x})}[\boldsymbol{w}^\top\boldsymbol{x}>0]\boldsymbol{x}$$

43

Example 2

# Calculating gradient for single neuron

| Error $\varepsilon$ | | $\dfrac{d\varepsilon}{d\hat{y}}$ | Activation $\boxed{\hat{y} = g(f)}$ | | $\dfrac{dg}{df}$ | Weighted Sum $f(\boldsymbol{w})$ | $\dfrac{df}{d\boldsymbol{w}}$ |
|---|---|---|---|---|---|---|---|
| Square Error | $\dfrac{1}{2}(y - \hat{y})^2$ | $-(y - \hat{y})$ | Sigmoid | $\dfrac{1}{1 + e^{-f}}$ | $(1 - g)g$ | $\boldsymbol{w}^\mathsf{T}\boldsymbol{x}$ | $\boldsymbol{x}$ |
| Binary Cross Entropy | $-y \log \hat{y}$ | $-\dfrac{y}{\hat{y}}$ | tanh | $\tanh f$ | $1 - g^2$ | | |
| | | | ReLU | $\max(0, f)$ | $[f > 0]$ | | |

$$\varepsilon = \frac{1}{2}(y - \hat{y})^2 = \frac{1}{2}(y - g)^2 = \frac{1}{2}\left(y - \frac{1}{1 + e^{-f}}\right)^2 = \frac{1}{2}\left(y - \frac{1}{1 + e^{-\boldsymbol{w}^\mathsf{T}\boldsymbol{x}}}\right)^2$$

$$\nabla_{\boldsymbol{w}}\varepsilon = \frac{d\varepsilon}{d\hat{y}}\frac{dg}{df}\frac{df}{d\boldsymbol{w}} = -(y - \hat{y})(1 - \hat{y})\hat{y}\boldsymbol{x}$$
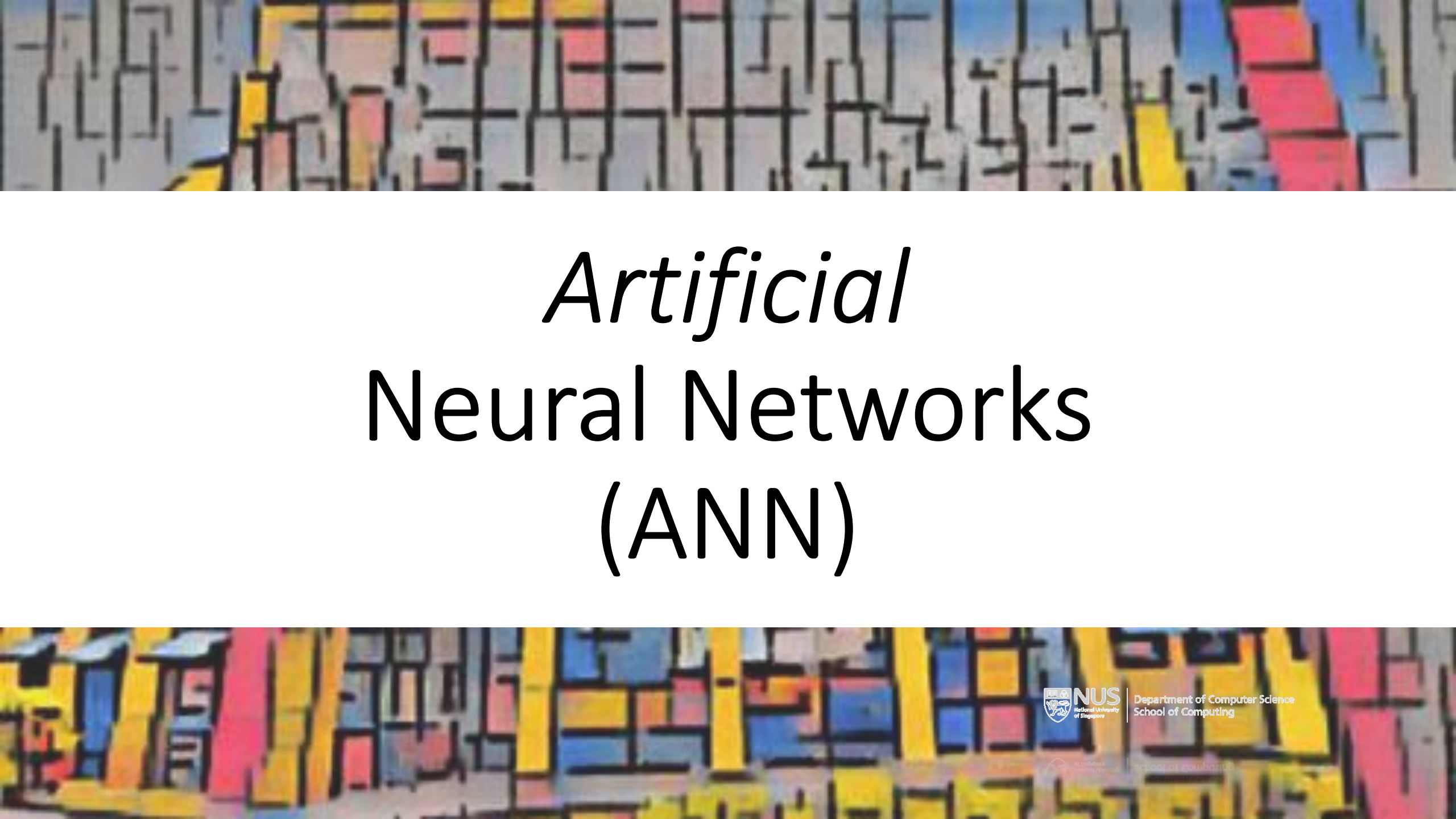
# Extending the Perceptron

- Perceptron is a linear classifier
- **Non-linear** classifiers
  - Other **activation functions**
    - Differentiable ones!
  - **Multiple** perceptrons / neurons      Feed-forward
    - Multi-Layer Perceptron (MLP)           Neural Network

# *Artificial* Neural Networks (ANN)

# 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

# **Artificial NNs** are *inspired* by, but *not mimicking* of **Human NNs**

**1** Layers and Uni-directional inference

Artificial NNs usually uses a sequence of layers with specific order to determine an output. In human brain, there is no fixed order and often async. … Artificial NN are mostly feed forward; output cannot then affect input. This is unlike Human NNs, which have cyclic loops in the neural structure.

**2** Non-diverse neurons and structures [W10]

Artificial NNs are made up of only one type of simple neuron. Human NNs are made up of many kinds of neurons.
*Remedy: Convolutional neuron for images, Recurrent neuron for sequence.*

**3** Data-specific [W10]

Human NNs can do many task (e.g., recognize sound, image, *and* text). Artificial NNs, it will only be able to suit one task at a time (e.g., either image *or* text).
*Remedy: CNNs for images and RNNs for text can be combined.*

**4** Deterministic

For the same input the Artificial NN will give the same output; but this may not apply to Human NNs.
*Research: Bayesian neural networks include randomization at inference to predict more robustly.*

**5** Energy efficiency of computations

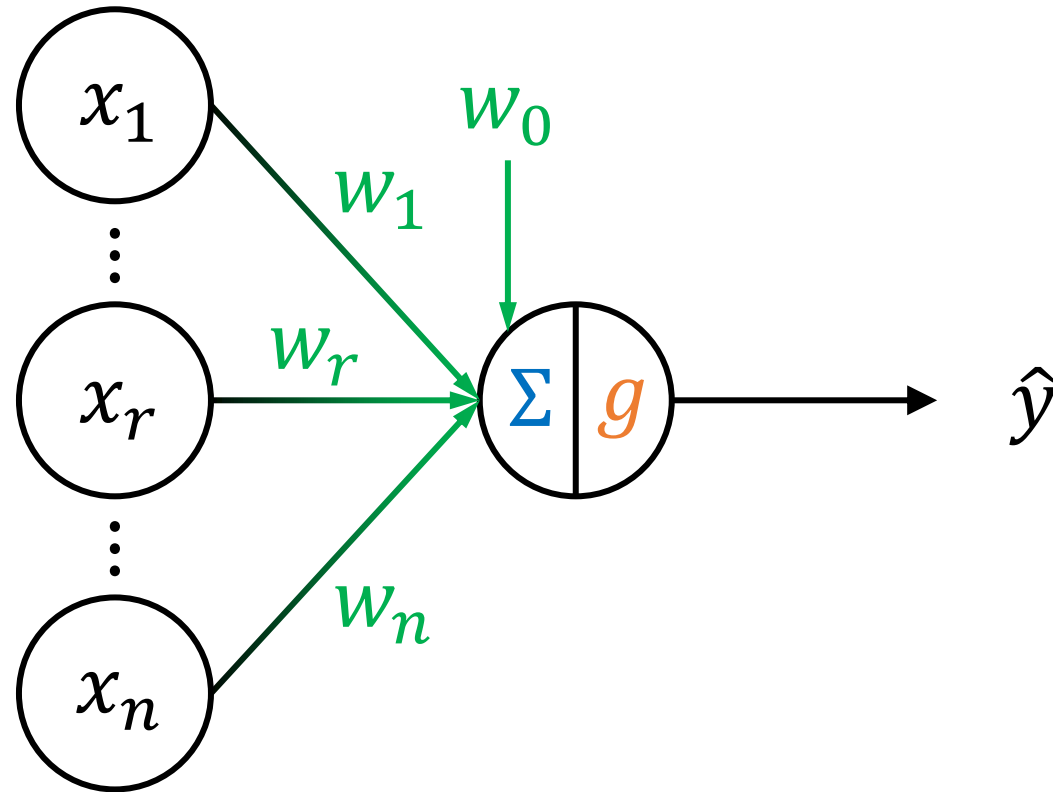Human NNs are more energy efficient than Artificial NN.
*Research: Spiking neural networks are being developed to be lower energy.*
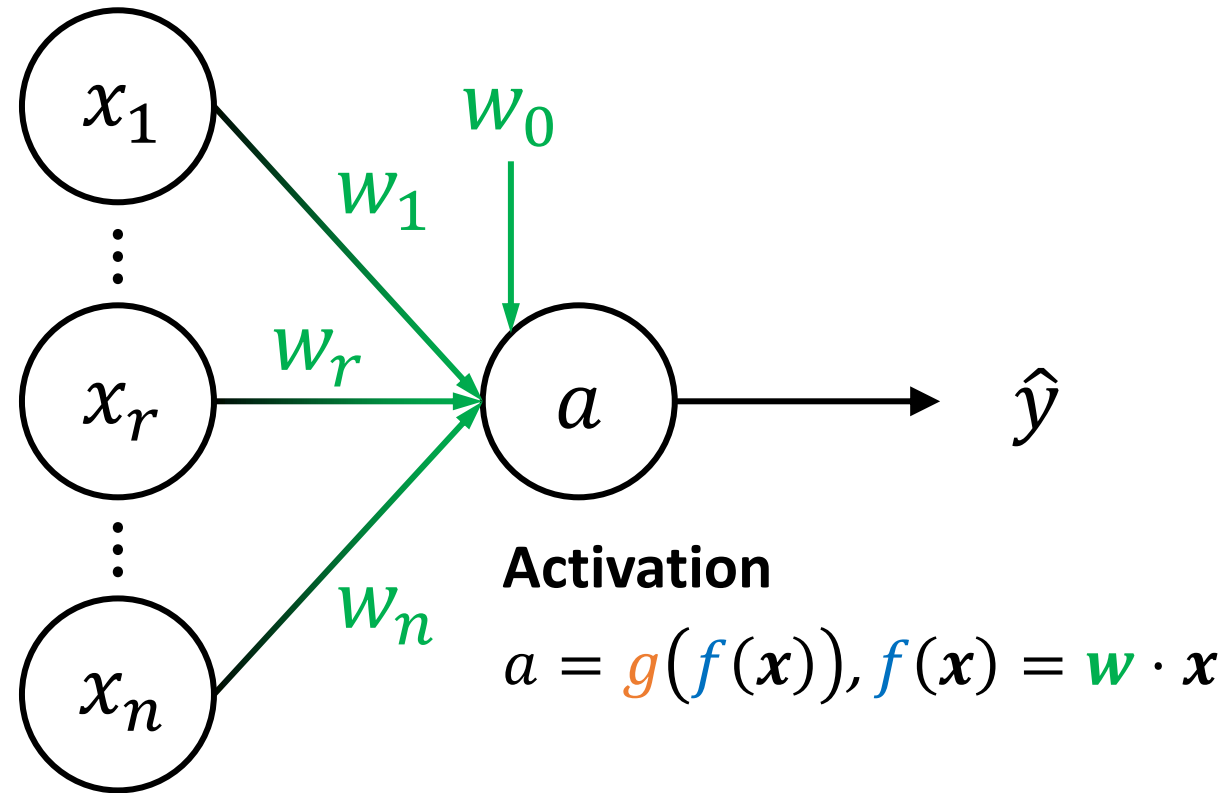
**6** Forgetting and Unlearning

Human can forget but Artificial NN will not. Once the Artificial NN is trained, it will remember what it learns and it becomes permanent knowledge.
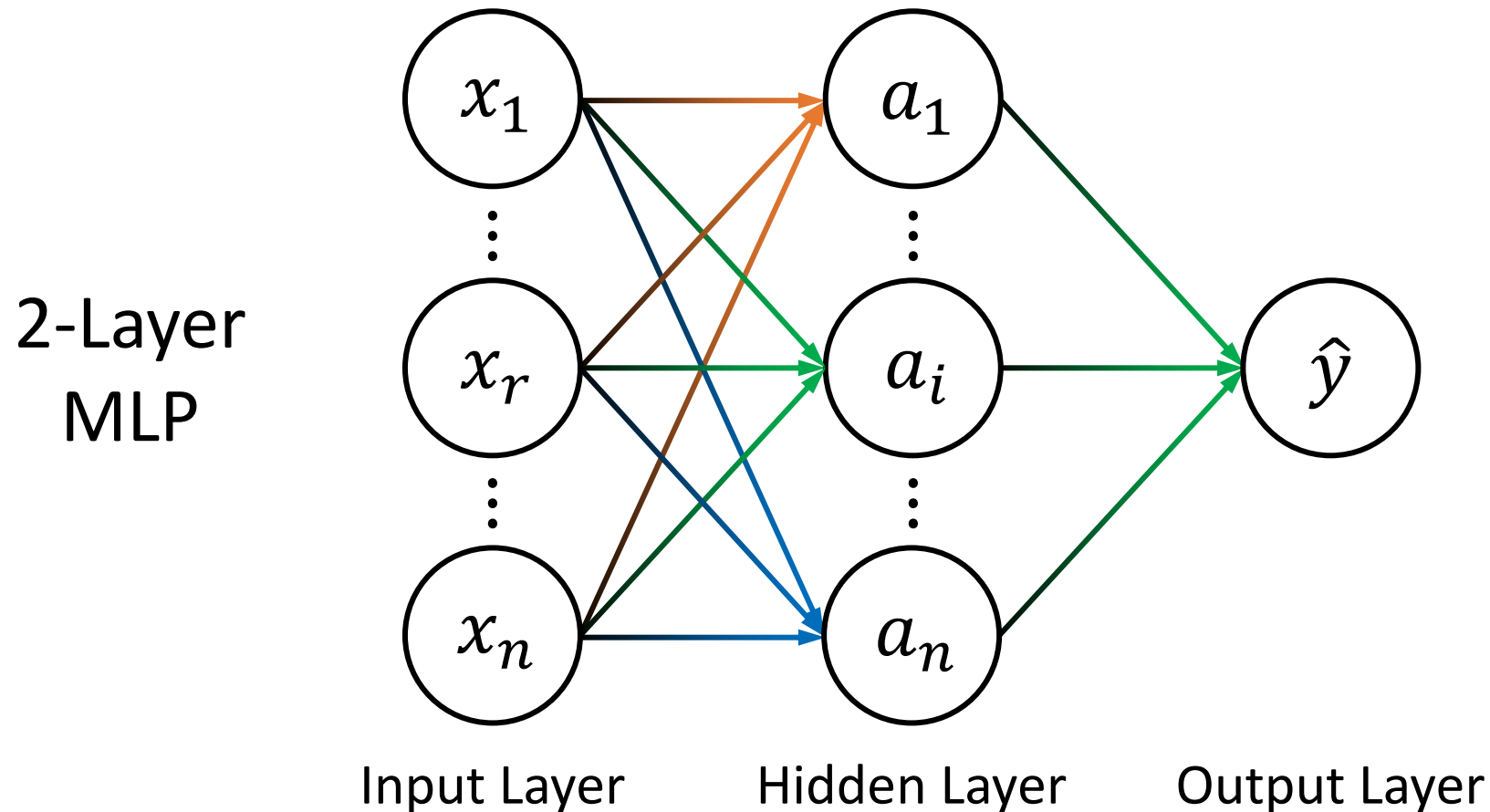*Research: Model unlearning can enable models to forget instances for legal privacy requirements.*
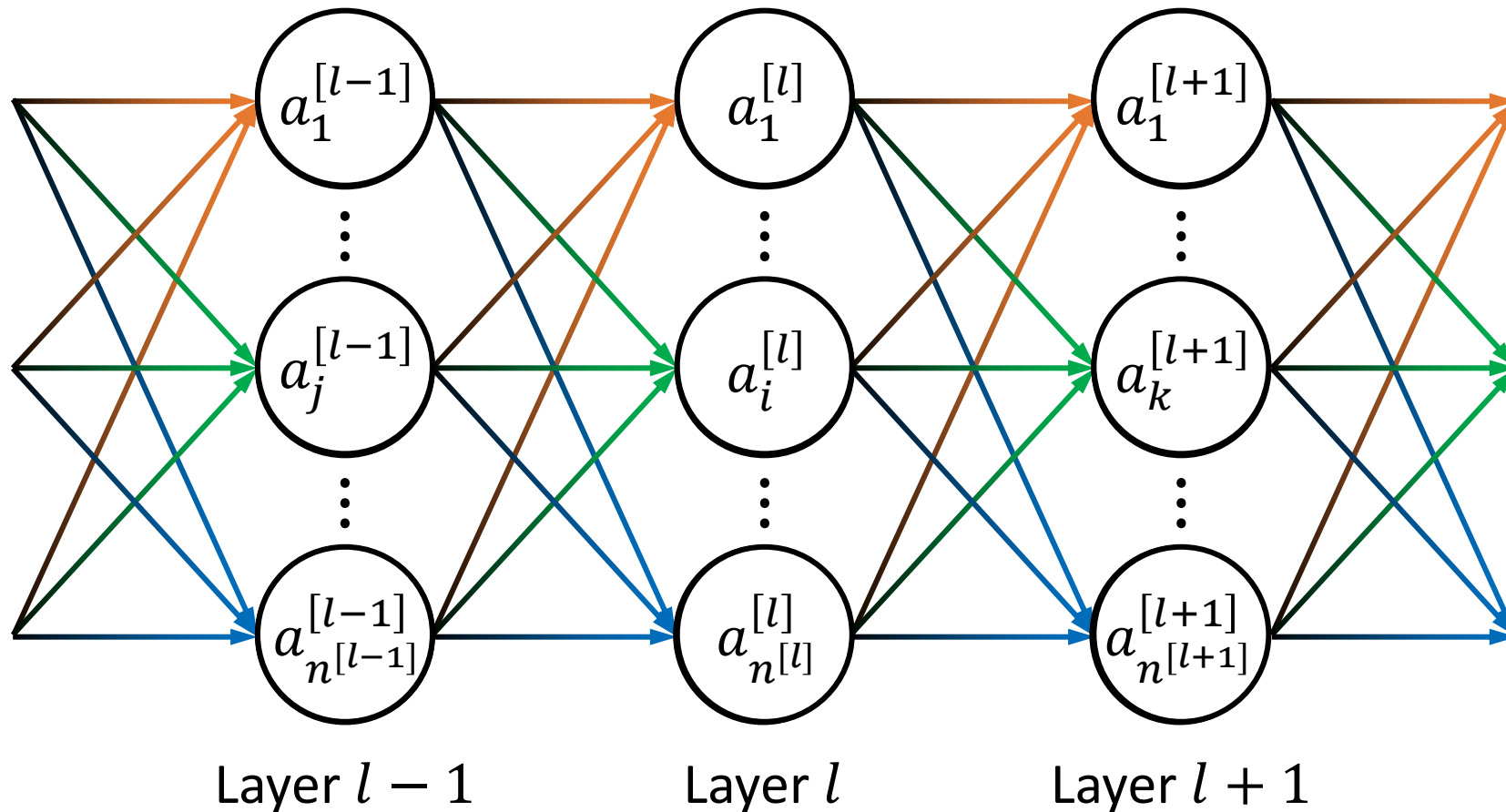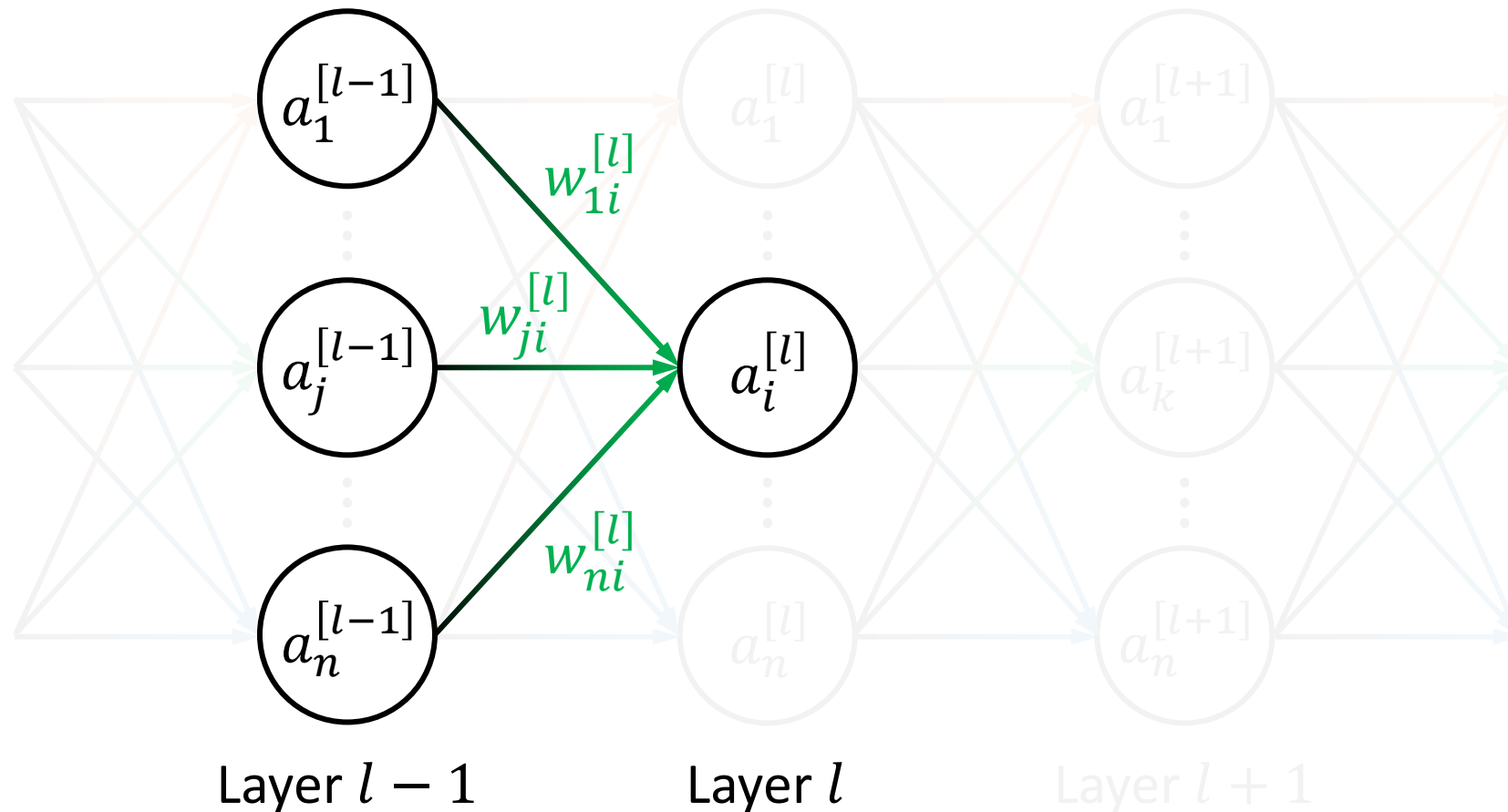
# Single-Layer Perceptron

# 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)

# Multi-Layer Perceptron (Neural Network)

# Multi-Layer Perceptron (Neural Network)



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

# Fitting **non-linear function** with MLP
# What model <u>weights</u> can model $\hat{y} = |x - 1|$?
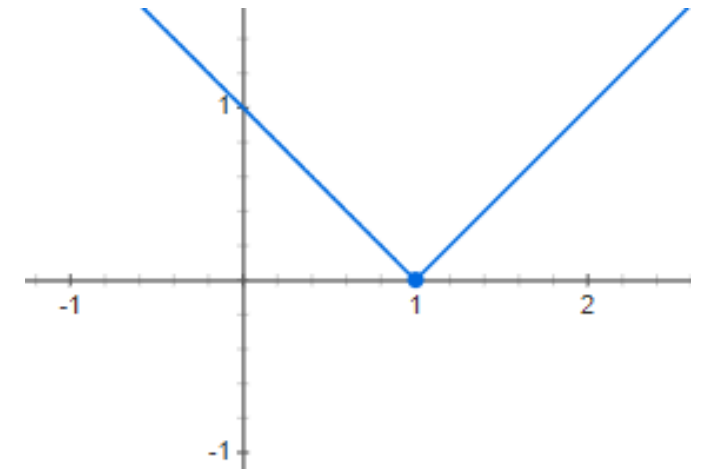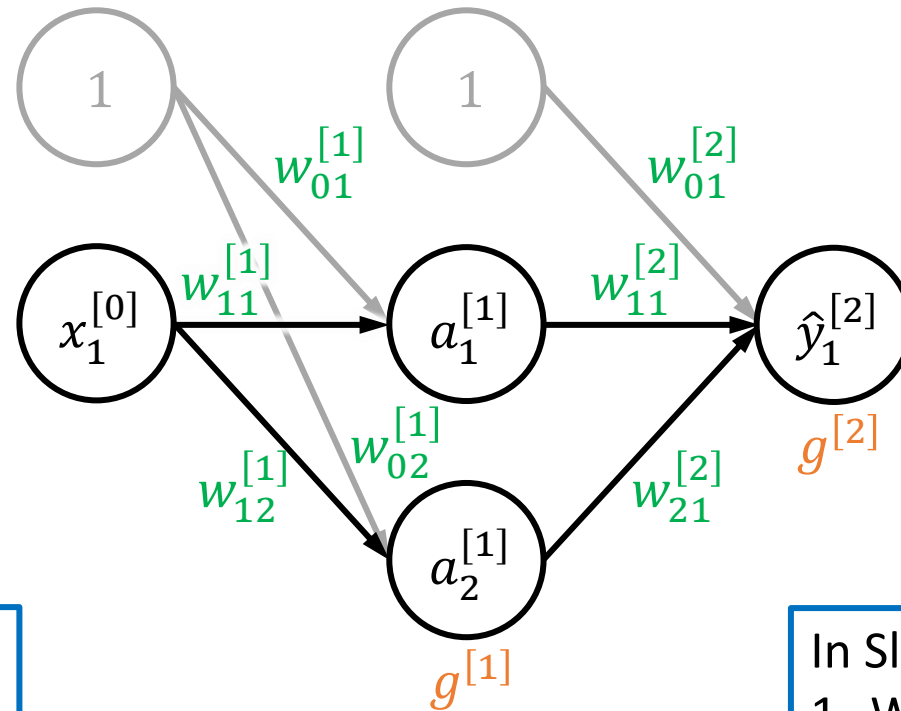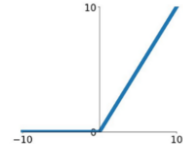
**Sigmoid**
$\sigma(x) = \frac{1}{1+e^{-x}}$

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

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

$w_{01}^{[1]}$

$w_{01}^{[2]}$

$w_{11}^{[1]}$

$w_{11}^{[2]}$

$x_1^{[0]}$

$a_1^{[1]}$

$\hat{y}_1^{[2]}$

$w_{02}^{[1]}$

$g^{[2]}$

$w_{12}^{[1]}$

$w_{21}^{[2]}$

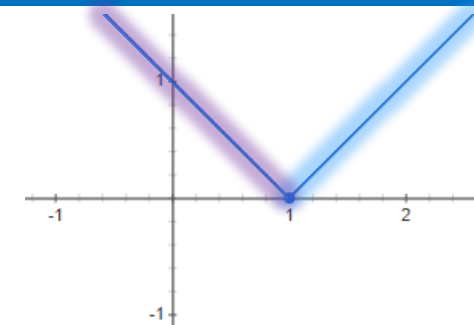$a_2^{[1]}$

$g^{[1]}$

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

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

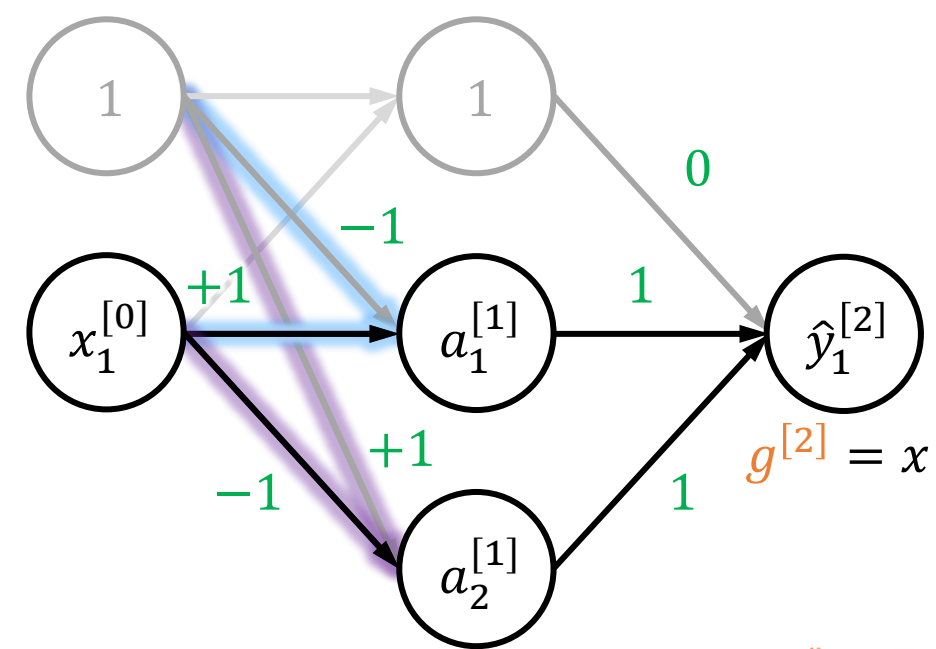# Fitting **non-linear function** with MLP
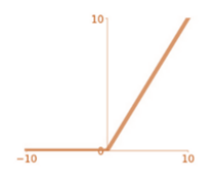## What model <u>weights</u> can model $\hat{y} = |x - 1|$?

$$\boldsymbol{W}^{[1]} = \begin{pmatrix} 1 & w_{01}^{[1]} & w_{02}^{[1]} \\ 0 & w_{11}^{[1]} & w_{12}^{[1]} \end{pmatrix}$$
$$= \begin{pmatrix} 1 & -1 & 1 \\ 0 & 1 & -1 \end{pmatrix}$$

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

$$\boldsymbol{a}^{[1]} = g^{[1]} \left( \left( \boldsymbol{W}^{[1]} \right)^{\top} \boldsymbol{x}^{[0]} \right)$$
$$= \begin{pmatrix} 1 \\ \text{ReLU}\left( -1 + x_1^{[0]} \right) \\ \text{ReLU}\left( 1 - x_1^{[0]} \right) \end{pmatrix}$$

$$\boldsymbol{W}^{[2]} = \begin{pmatrix} w_{01}^{[2]} \\ w_{11}^{[2]} \\ w_{21}^{[2]} \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}$$

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

$$\hat{y}^{[2]} = g^{[2]} \left( \left( \boldsymbol{W}^{[2]} \right)^{\top} \boldsymbol{a}^{[1]} \right)$$
$$= 0$$
$$+ \text{ReLU}\left( -1 + x_1^{[0]} \right)$$
$$+ \text{ReLU}\left( 1 - x_1^{[0]} \right)$$

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

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

$g^{[2]} = x$

$g^{[1]} = \text{ReLU}$
$= \max(0, x)$

Questions!

**desmos**

1. $y = |x - 1|$

2. $y = \max(0, x - 1) + \max(0, 1 - x) + a$

3. $a = 0$

0 ———●————————————— 2



powered by
**desmos**

Source: https://www.desmos.com/calculator/f05dzu2rdv

Web Viewer Terms | Privacy & Cookies    Edit

# Neural Network (vector notation)



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

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

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

Layer $l$ in
Neural Network

Layer $l$
Activations
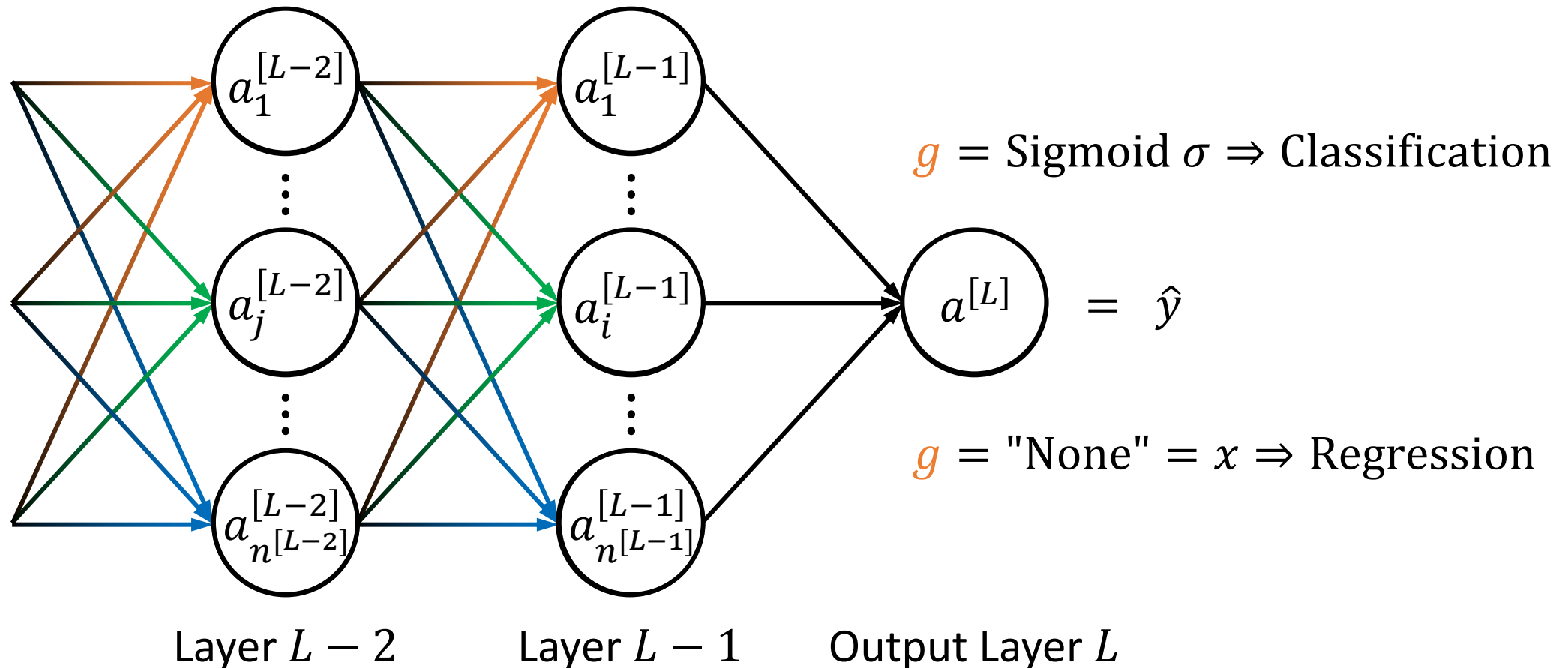
Layer $l-1$
Activations

# Multiclass / Multilabel Neural Networks

# Single-output Neural Network
(binary classification or scalar regression)



$g = \text{Sigmoid } \sigma \Rightarrow \text{Classification}$

$a^{[L]} = \hat{y}$

$g = \text{"None"} = x \Rightarrow \text{Regression}$

Layer $L - 2$      Layer $L - 1$      Output Layer $L$

# Multiple outputs Neural Network
## (Multiple classifications or multivariate regression)



Layer $L - 2$  Layer $L - 1$  Output Layer $L$

**Multi-label**

# Multiple outputs Neural Network
## (Multiclass classification)

$$\sum_{c=1}^{|C|} \hat{y}_c = 1$$



Softmax

$\sigma$

Layer $L-2$      Layer $L-1$      Layer $L$      Output

**Multi-class** classification

# Sigmoid vs. Softmax

| **Sigmoid** | **Multiple Sigmoids** | **Softmax** |
|---|---|---|
| • *Is prediction true?* | • *Which predictions are true?* | • *Which class is most probably true?* |
| • For **binary** classification | • For **multi-label** classification | • For **multiclass** classification |
| • $\sigma(z) = \frac{e^z}{1+e^z}$ | • $\sigma(z_k) = \frac{e^{z_k}}{1+e^{z_k}}, \ \forall k$ | • $\boldsymbol{\sigma}(z_c) = \frac{e^{z_c}}{\sum_{c=1}^{|C|}(1+e^{z_c})} \boldsymbol{e}_c$ |

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

$$\sigma(\boldsymbol{z}) = \begin{pmatrix} e^{z_1}/(1 + e^{z_1}) \\ \vdots \\ e^{z_N}/(1 + e^{z_N}) \end{pmatrix}$$

$$\boldsymbol{\sigma}(\boldsymbol{z}) = \frac{\begin{pmatrix} e^{z_1} \\ \vdots \\ e^{z_{|C|}} \end{pmatrix}}{\sum_{c=1}^{|C|}(1 + e^{z_c})} = \frac{e^{\boldsymbol{z}}}{(\boldsymbol{1} + e^{\boldsymbol{z}}) \cdot \boldsymbol{1}}$$

# Multiple outputs Neural Network
(Multiclass classification)



Softmax

One-Hot

$a_1^{[L-1]}$

$a_i^{[L-1]}$

$\boldsymbol{\sigma}$

$a_{n^{[L-1]}}^{[L-1]}$

$z_1$

$z_k$

$z_{|c|}$

$\hat{y}_1$

$\hat{y}_c$

$\hat{y}_{|c|}$

$y_1$

$y_c$

$y_{|c|}$

Layer $L-1$

Layer $L$

Output

Ground Truth

$$\varepsilon = \textbf{Categorical Cross Entropy} = -\boldsymbol{y} \cdot \log \hat{\boldsymbol{y}}$$
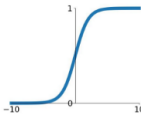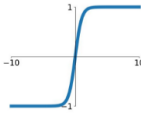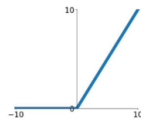
# Multiple outputs Neural Network
(Vector regression)

$\varepsilon = \textbf{Euclidean distance} = \|\boldsymbol{y} - \hat{\boldsymbol{y}}\|_2$

# Derivatives for calculating gradient

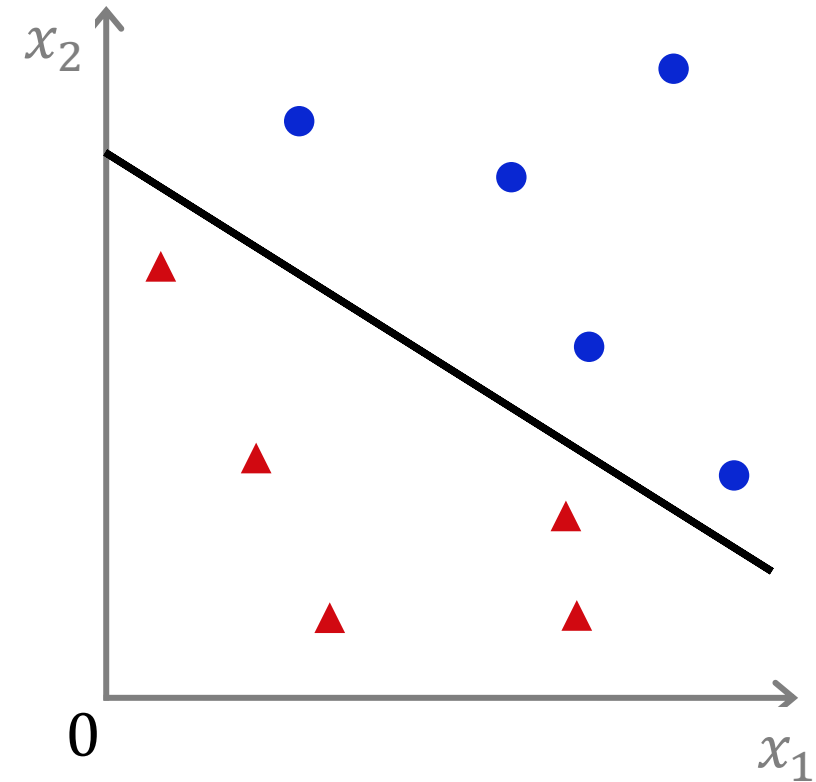| Error $\varepsilon(\hat{y})$ | | $\dfrac{d\varepsilon}{d\hat{y}}$ | Activation $g(f)$ | | $\dfrac{dg}{df}$ | Weighted Sum $f(\boldsymbol{w})$ | $\dfrac{df}{d\boldsymbol{w}}$ |
|---|---|---|---|---|---|---|---|
| Square Error | $\dfrac{1}{2}(y-\hat{y})^2$ | $-(y-\hat{y})$ | Sigmoid | | $\dfrac{1}{1+e^{-f}}$ $(1-g)g$ | $\boldsymbol{w}^{\top}\boldsymbol{x}$ | $\boldsymbol{x}$ |
| Binary Cross Entropy | $-y\log\hat{y}$ | $-\dfrac{y}{\hat{y}}$ | tanh | | $\tanh f$ $\quad 1-g^2$ | | |
| Error $\varepsilon(\boldsymbol{z})$ | | $\dfrac{d\varepsilon}{d\boldsymbol{z}}$ | ReLU | | $\max(0,f)$ $\quad[f>0]$ | | |
| Categorical Cross Entropy | $-\boldsymbol{y}\cdot\log\widehat{\boldsymbol{y}}$ | $\widehat{\boldsymbol{y}}-\boldsymbol{y}$ | Softmax | | $\dfrac{e^{\boldsymbol{z}}}{(\mathbf{1}+e^{\boldsymbol{z}})\cdot\mathbf{1}}$ $\quad\dfrac{dg}{d\boldsymbol{z}}$ | **Advanced:** <u>further reading</u> | |

# Perceptron → Neural Network

- Linear classifiers
- **Non-linear** classifiers
  - Other **activation functions**
    - Differentiable ones!
  - **Multiple** perceptrons / neurons
    - Multi-Layer Perceptron (MLP)



Can only model
**straight** lines

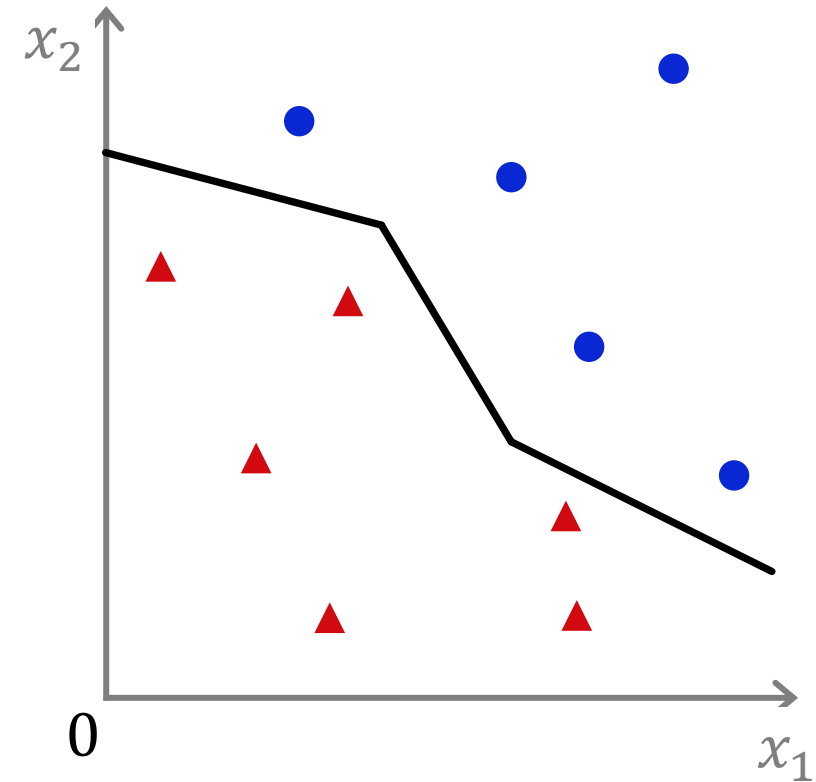# Perceptron → Neural Network

- Linear classifiers

- **Non-linear** classifiers
  - Other **activation functions**
    - Differentiable ones!
  - **Multiple** perceptrons / neurons
    - Multi-Layer Perceptron (MLP)



Can model **multiple** straight lines

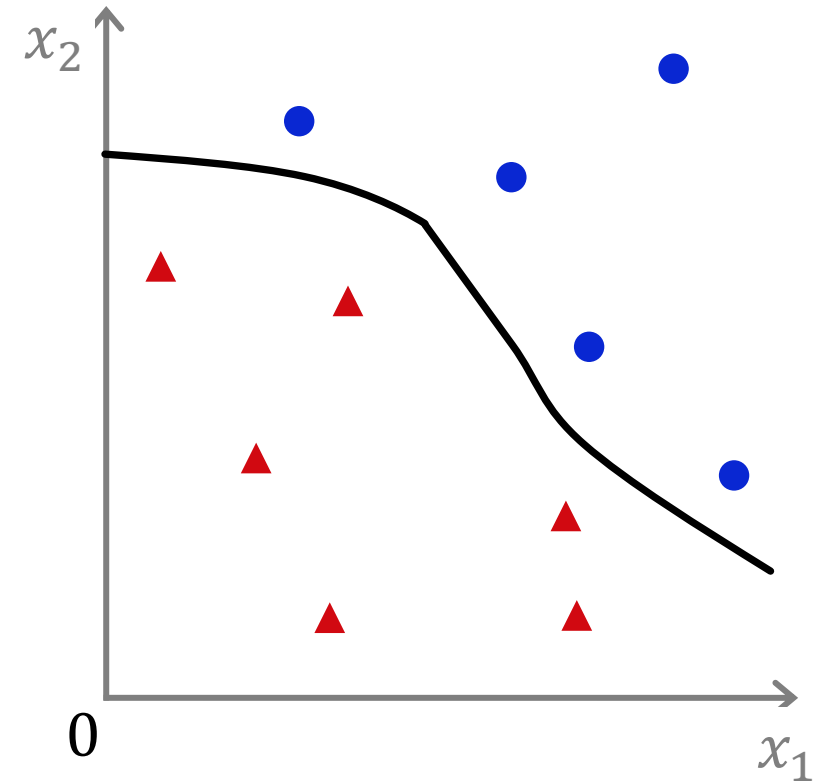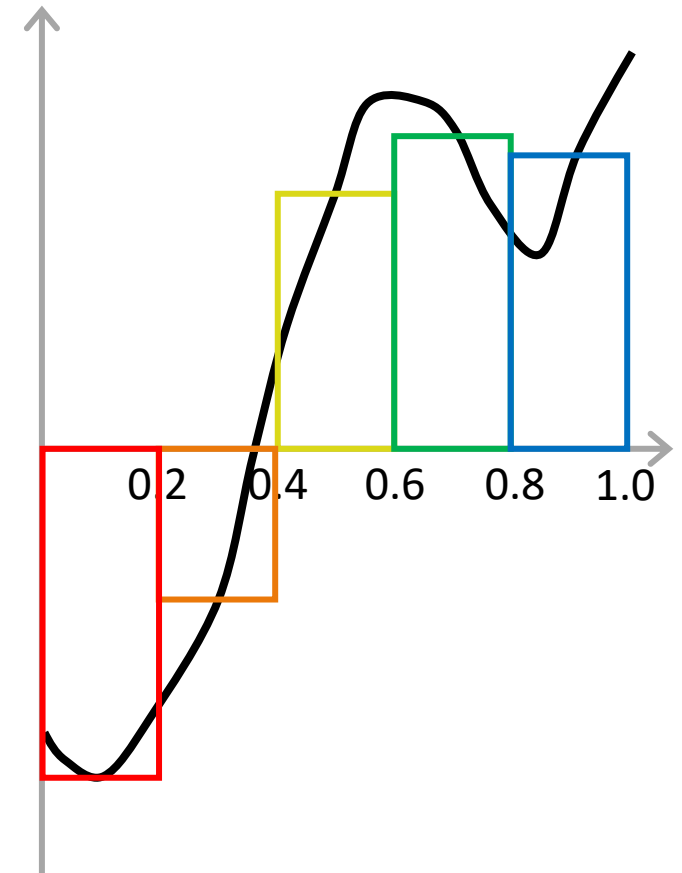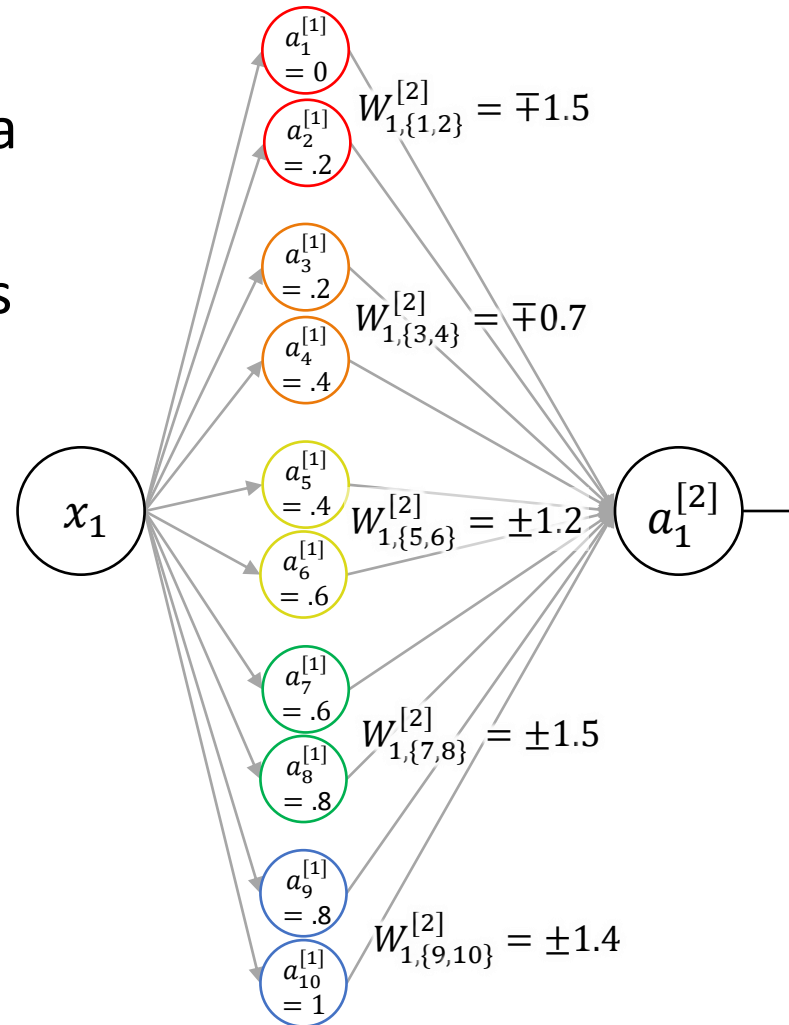# Perceptron → Neural Network

- Linear classifiers

- **Non-linear** classifiers
  - Other **activation functions**
    - Differentiable ones!
  - **Multiple** perceptrons / neurons
    - Multi-Layer Perceptron (MLP)

Can model **multiple** *curved* lines

# Universal Approximation Theorem

- Each neuron contributes a **piecewise function**

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



$a_1^{[1]} = 0$

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

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

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

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

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

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

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

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

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

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

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

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

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

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

$x_1$

$a_1^{[2]}$

# Gradient Descent
## for Neural Networks
# Backpropagation

# W09B Thursday: Backpropagation

Image credit: https://www.mensjournal.com/wp-content/uploads/mf/main-turn-your-training-around-with-backward-running.jpg

NUS CS3244: Machine Learning

# Wrapping Up

# What did we learn?