



The Joy of *Finally* Learning

www.mldawn.com

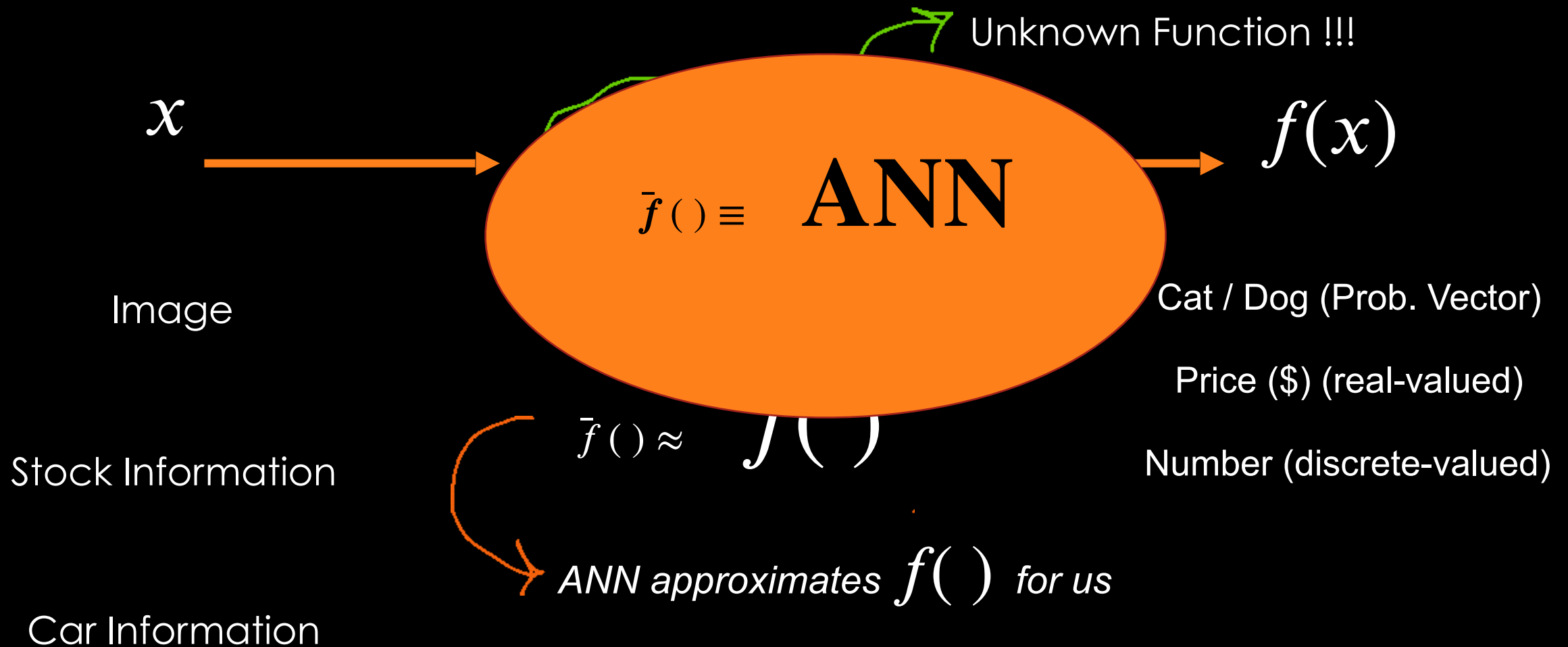
OUTLINE OF THE TALK

- What is an Artificial Neural Network (ANN)?
- What is a Perceptron?
- Bayes' Rule and Birth of Error Functions (with Derivations)
- The Birth of Binary-Cross Entropy Function (with Derivations)
- Forward Pass and Backward Pass in a Binary Classifier
- A Demo

WHAT IS A NEURAL NETWORK?

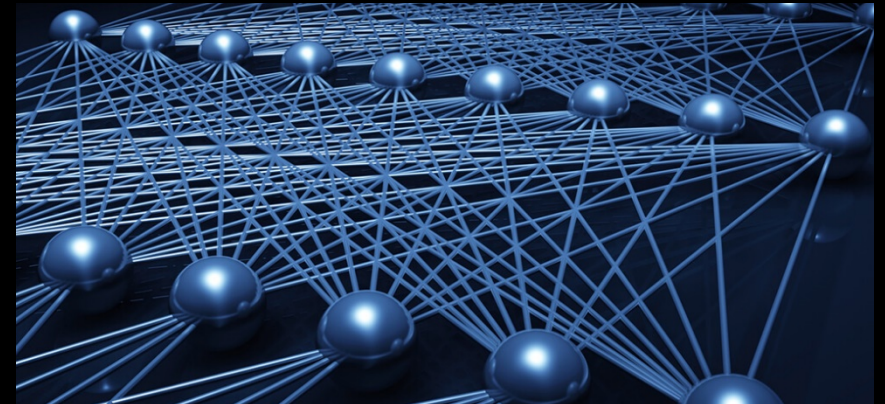
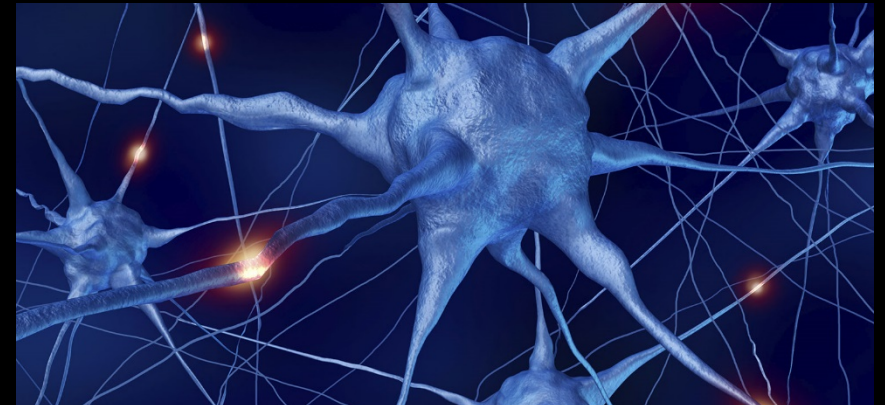
- An ANN is a function approximator!
 - What is a function?
 - Examples of functions: *real-valued, discrete-valued, and vector-valued*
- Applications of ANNs:
 - Recognizing hand-written characters
 - Recognizing spoken words
 - Recognizing faces

WHAT IS A NEURAL NETWORK?



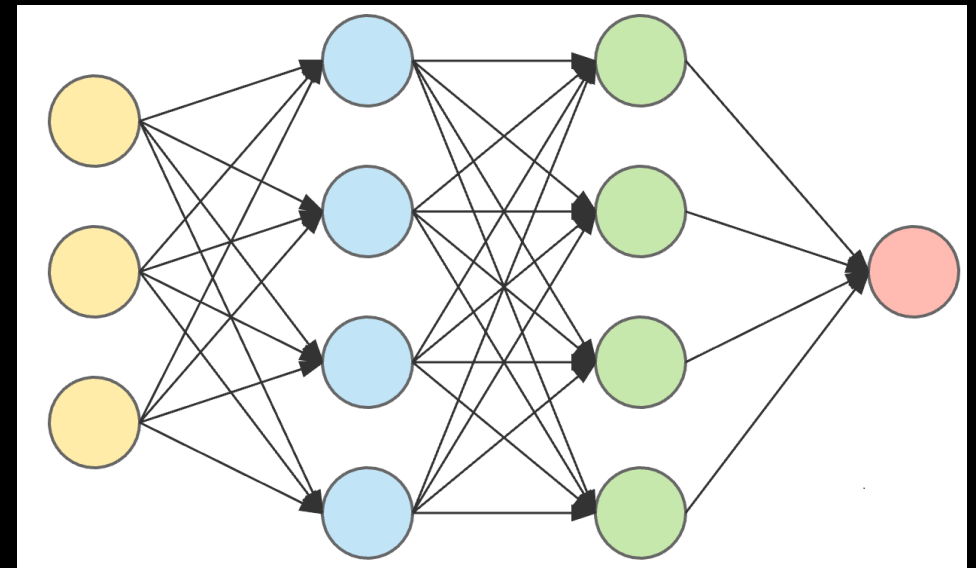
ANN BIOLOGICAL MOTIVATION

- ANNs are inspired by Biological Learning Systems (BLS)
 - BLS: A massive interconnection of neurons
- ANNs:
 - Densely interconnected set of simple unites



ANN BIOLOGICAL MOTIVATION

- For each unit in ANN:
 - **Possibly** several inputs
 - A Single real-valued output
- BLS processes are highly parallelized
 - Over nearly 10^{11} neurons
- ANNs do NOT fully reflect a BLS
 - Output of neurons are different
- Researchers try to:
 1. Model BLS processes
 2. Develop effective ANNs (no regards to BLS compatibility)



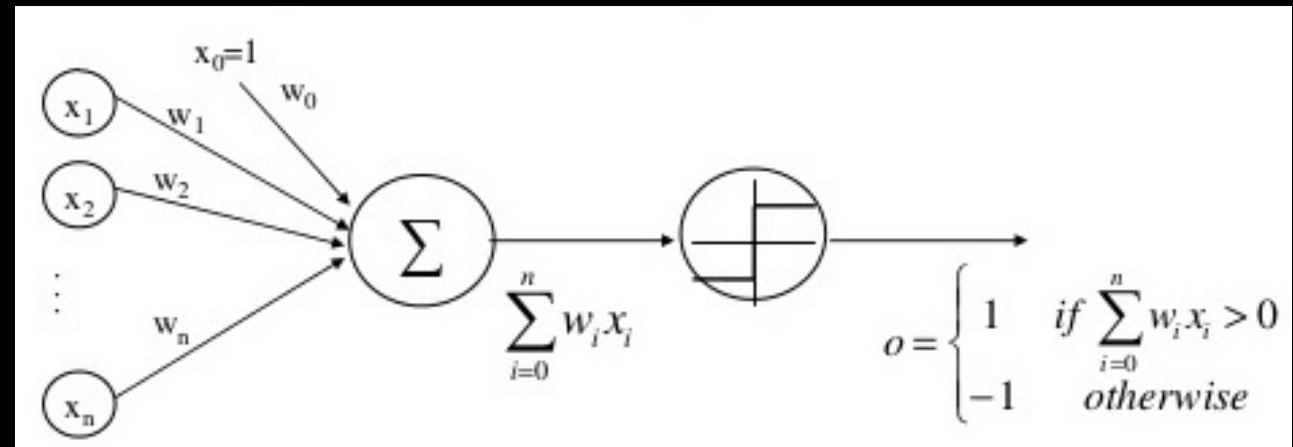


SIMPLEST ANN

A PERCEPTRON

A PERCEPTRON

- A type of ANN that is based on a unit called **perceptron**
- There is a rule about input/output of a perceptron:
 - Input: A vector of real-valued data
 - Linearly combine the components of the input



$$o(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n > 0 \\ -1 & \text{Otherwise} \end{cases}$$

A PERCEPTRON

- What are the weights? What is that w_0 ?

- Let's imagine an additional input

- $x_0 = 1$

- This simplifies the inequality:

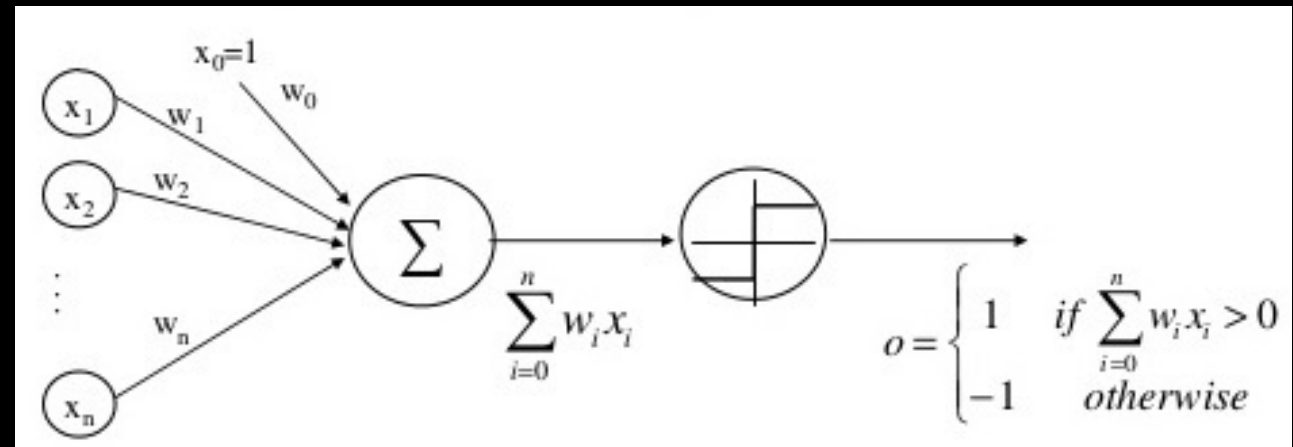
- $$\sum_{i=0}^n w_i x_i > 0$$

- Or in vector form: $\vec{w} \cdot \vec{x} > 0$

- Can we show the perceptron smarter?

- $o(\vec{x}) = \text{sgn}(\vec{w} \cdot \vec{x})$

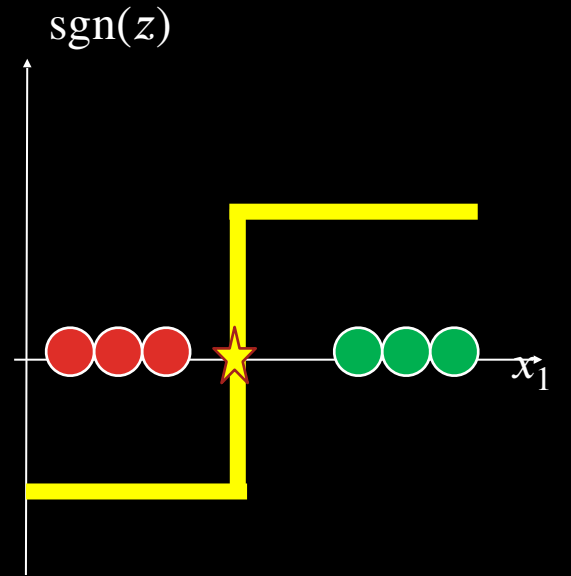
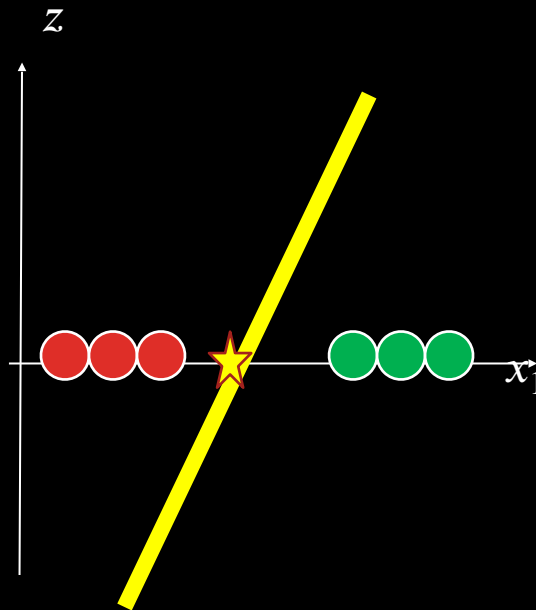
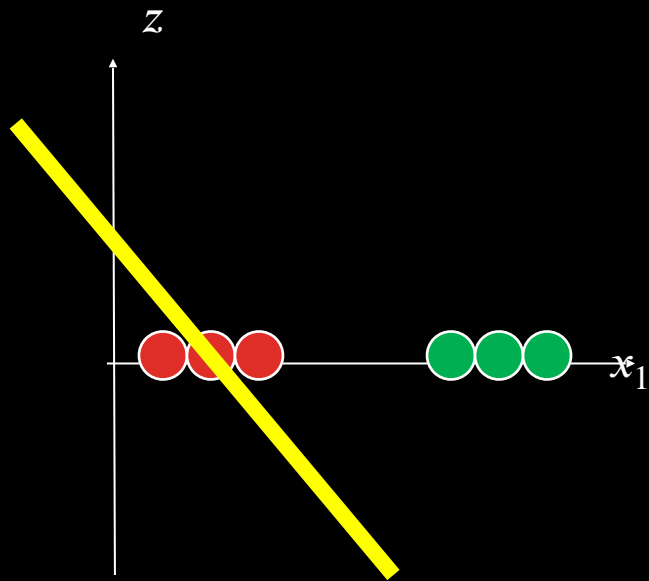
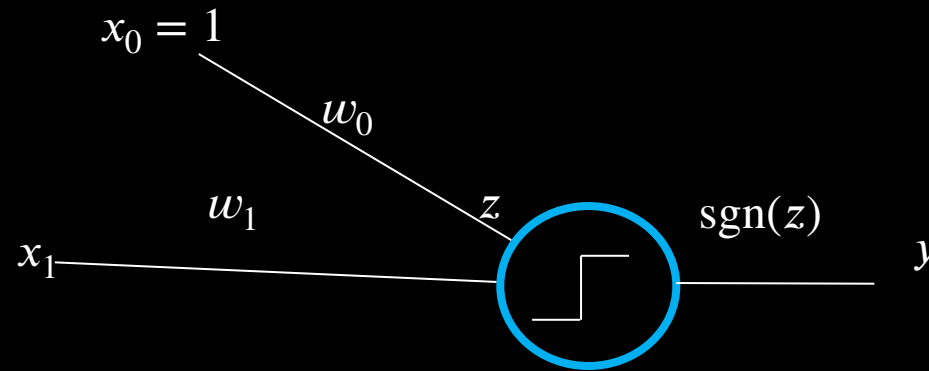
- Where:
$$\text{sgn}(z) = \begin{cases} +1 & \text{if } z > 0 \\ -1 & \text{otherwise} \end{cases}$$



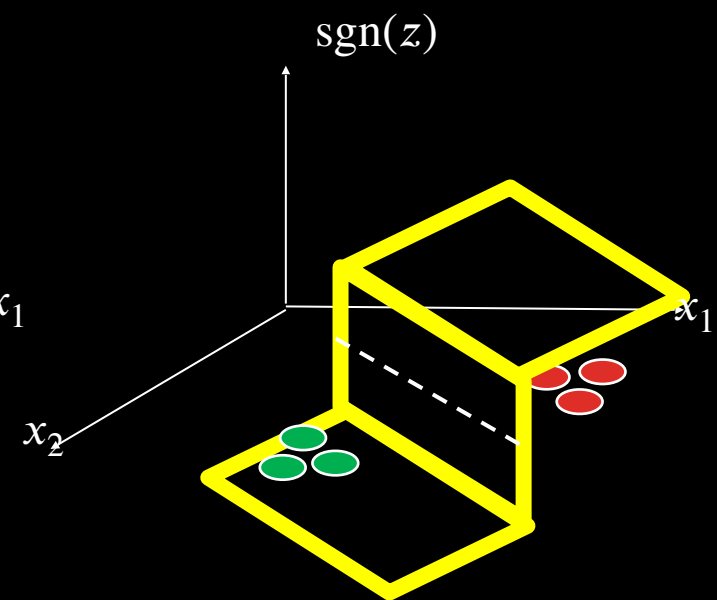
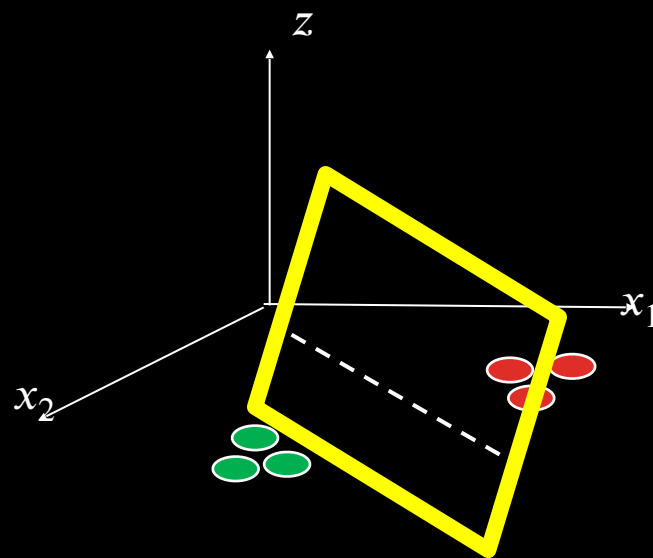
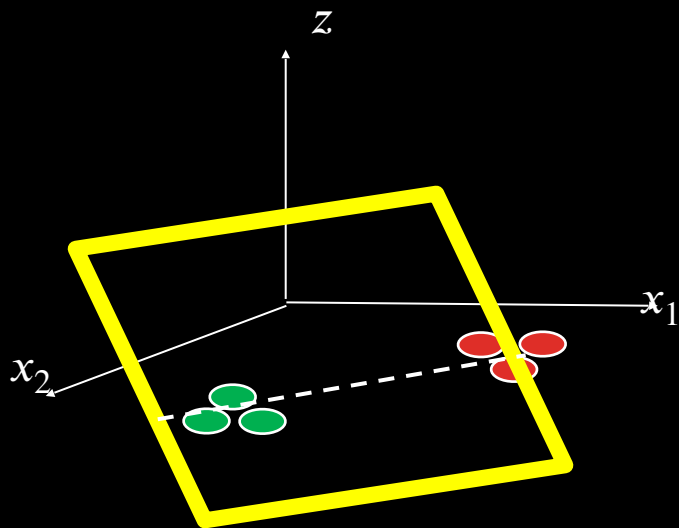
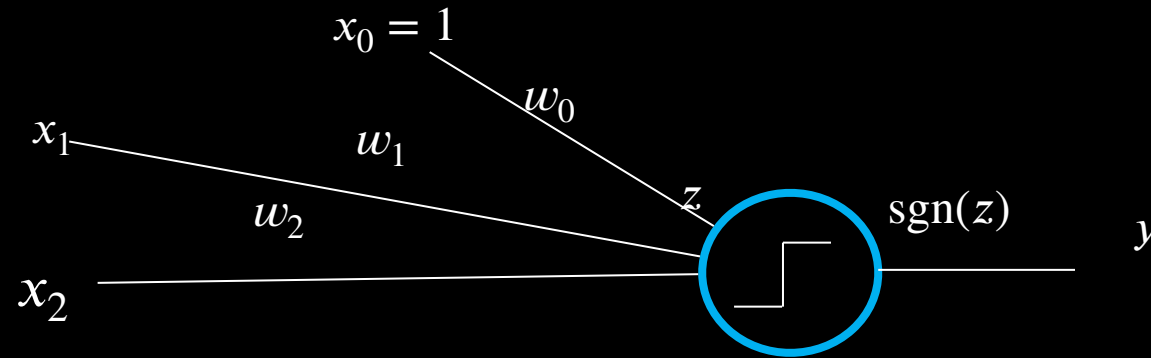
- Learning a perceptron = Choosing the weights

- The hypothesis space: $H = \{ \vec{w} \mid \vec{w} \in \mathbb{R}^{(n+1)} \}$

A PERCEPTRON: A DEEPER LOOK



A PERCEPTRON: A DEEPER LOOK





THE BIRTH OF ERROR
FUNCTIONS

DO I *TRUST* MY ANN?

BAYES' RULE

- Let's assume the following notation:
 - M : Our Model
 - D : Our Evidence
- Then the Bayes' Rule:

$$P(M | D) = \frac{P(D | M) P(M)}{P(D)}$$

TAKING THE LOGS() : A SIMPLIFICATION

$$P(M|D) = \frac{P(D|M)P(M)}{P(D)} \xrightarrow{\log\left(\frac{a}{b}\right) = \log(a) - \log(b)} \ln P(M|D) = \ln [P(D|M)P(M)] - \ln P(D)$$

$$\log(a \times b) = \log(a) + \log(b)$$

$$\ln P(M|D) = \ln P(D|M) + \ln P(M) - \ln P(D)$$

THE BAYESIAN MACHINERY (MAPE)

$$\ln P(M | D) = \ln P(D | M) + \ln P(M) - \ln P(D)$$



Maximize $P(M | D) \equiv \text{Minimize } -P(M | D)$

$$-\ln P(M | D) = -\ln P(D | M) - \ln P(M) + \ln P(D)$$

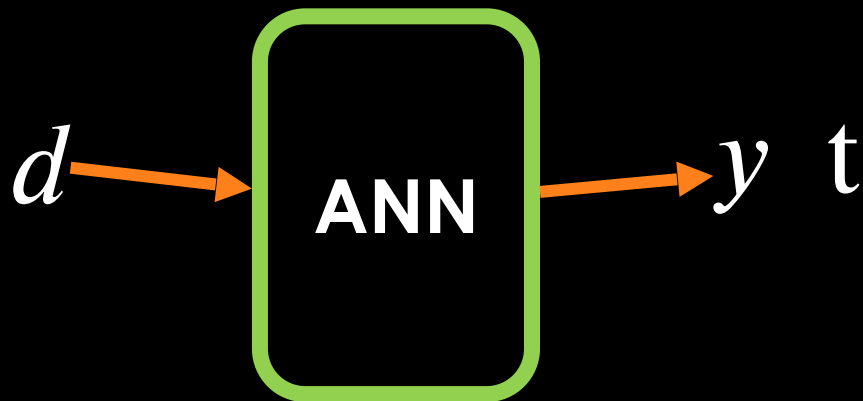


$$E = -\ln P(M | D)$$

This is how Bayes' Rule connects to Error Functions

MAPE AND ANNS : SUPERVISED LEARNING

- How can we apply the MAPE to an ANN with weights w ?
- Let's assume that we have a set of examples (d_i, t_i) where d_i is the i -th input and t_i the corresponding target.
- Then we can rewrite the likelihood as:



$$P(D | M) \text{ and } D \approx (d_i, t_i)$$

$$P(D | M) = P((d_i, t_i) | M)$$

$$P((d_i, t_i) | M) = P(t_i | d_i, M) \times \underbrace{P(d_i | M)}_{P(d_i)}$$

MAPE AND ANNS : SUPERVISED LEARNING

$$-\ln P(M|D) = -\ln P(D|M) - \ln P(M) + \ln P(D)$$

Model indeed means our weights W

$$-\ln P(W|D) = -\ln P(D|W) - \ln P(W) + \ln P(D)$$

$$P(t_i | d_i, W) \times P(d_i)$$

P(D) and P(d_i) are both independent of our weights W. Ignore them during Maximization

$$-\ln P(W|D) = -\ln \prod_{i=1}^N \left[P(t_i | d_i, W) \times P(d_i) \right] - \ln P(W) + \ln P(D)$$

$$-\ln P(W|D) = -\sum_{i=1}^N \ln \left[P(t_i | d_i, W) \times P(d_i) \right] - \ln P(W) + \ln P(D)$$

$$-\ln P(W|D) = -\sum_{i=1}^N \ln \left[P(t_i | d_i, W) \right] - \sum_{i=1}^N \ln P(d_i) - \ln P(W) + \ln P(D)$$

MAPE AND ANNS : SUPERVISED LEARNING

$$-\ln P(W|D) = - \sum_{i=1}^N \ln \left[P(t_i | d_i, W) \right] - \ln P(W)$$

Can be used for regularizing W to avoid overfitting

$$P(W) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(w-\mu)^2}{2\sigma^2}} \quad \ln(\) \longrightarrow \approx -(w_1^2 + w_1^2 + \dots + w_N^2)$$

$$E = -\ln P(W|D) = - \sum_{i=1}^N \ln \left[P(t_i | d_i, W) \right]$$



BINARY CROSS-ENTROPY
ERROR FUNCTION

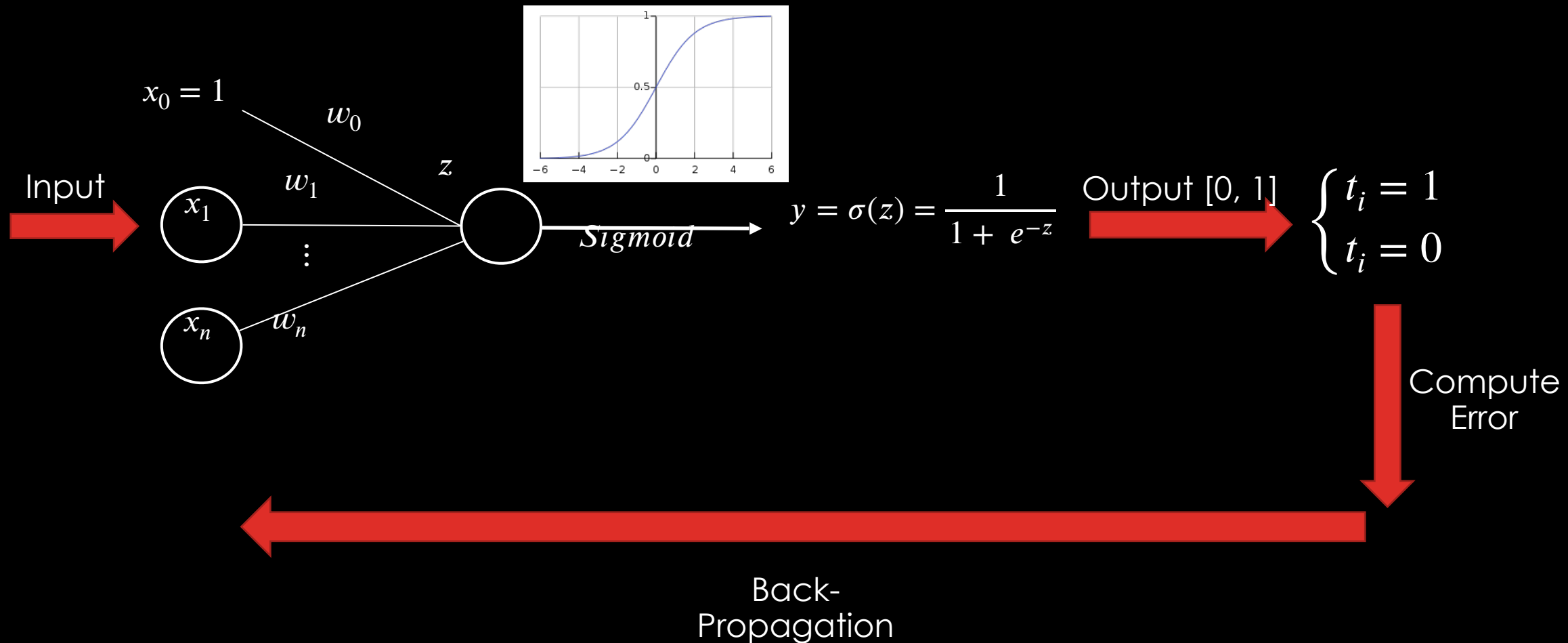
THE DERIVATION

BINARY CLASSIFICATION

- We have 2 classes: Positive and Negative
- Binary Classification:
 - $t_i = 1$ means that x_i belongs to the **positive** class
 - $t_i = 0$ means that x_i belongs to the **negative** class
- Artificial Neural Networks (ANNs): Probabilistic approach
 - ANN computes the probability of the **Positive Class**: $P(y_i = 1 \mid x_i)$
 - Could a linear model work? Say, $y = w^T x + b \rightarrow$ **Not Bounded!**

We need a Squasher!

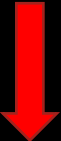
BINARY CLASSIFICATION



BINARY CLASSIFICATION

- The output of our ANN is probability. Let's denote:
 - y_i : Probability x_i belongs to the Positive class (i.e., $t_i = 1$)
 - $(1 - y_i)$: Probability x_i belongs to the Negative class (i.e., $t_i = 0$)

$$P(t_i | d_i, w) = \begin{cases} y_i & \text{if } t_i = 1 \\ 1 - y_i & \text{if } t_i = 0 \end{cases}$$


$$P(t_i | d_i, w) = y_i^{t_i} (1 - y_i)^{1-t_i}$$

**Bernoulli
Distribution**



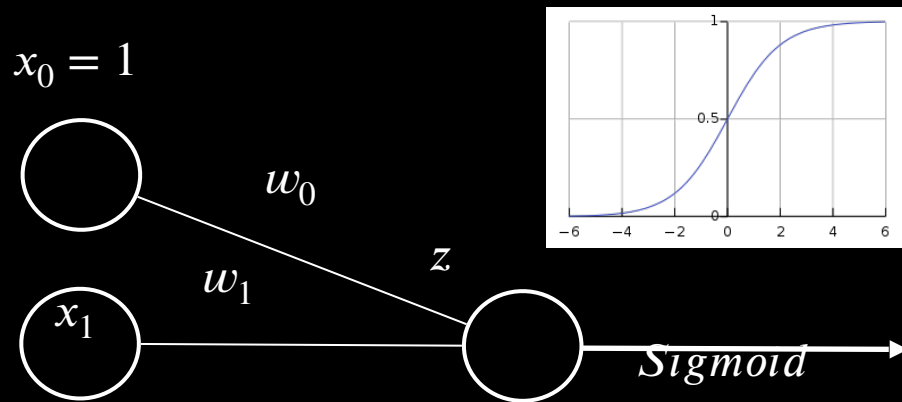
$$P(t_i | d_i, w) = y_i^{t_i} (1 - y_i)^{1-t_i}$$

BINARY CLASSIFICATION

$$\begin{aligned} E &= - \sum_{i=1}^N \ln [P(t_i | d_i, \mathbf{W})] = - \sum_{i=1}^N \ln [y_i^{t_i} (1 - y_i)^{1-t_i}] \\ &= - \sum_{i=1}^N [\ln y_i^{t_i} + \ln(1 - y_i)] \end{aligned}$$

$$E = - \sum_{i=1}^N [t_i \ln y_i + (1 - t_i) \ln(1 - y_i)]$$

FORWARD PASS



$$y = \sigma(z) = \frac{1}{1 + e^{-z}} \rightarrow E = -t \times \ln y - (1 - t) \times \ln(1 - y)$$

$$\text{Error} = E(\underbrace{\sigma(Z)}_y) = E(\underbrace{\sigma(w_1 x_1 + w_0)}_Z)$$

Binary Cross-Entropy Error
Function

BACKWARD PASS

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial y} \times \frac{\partial y}{\partial Z} \times \frac{\partial Z}{\partial w_1} = \left[\frac{-t}{y} + \frac{(1-t)}{(1-y)} \right] \times y(1-y) \times x_1 = (y-t) \times x_1$$

$$\frac{\partial E}{\partial w_0} = \frac{\partial E}{\partial y} \times \frac{\partial y}{\partial Z} \times \frac{\partial Z}{\partial w_0} = \left[\frac{-t}{y} + \frac{(1-t)}{(1-y)} \right] \times y(1-y) \times x_0 = (y-t)$$

- So the learning rule for both our learnable parameters (i.e., weights) is as follows:

$$w_1 = w_1 - \eta \times \frac{\partial E}{\partial w_1} = w_1 - \eta \times [(y-t) \times x_1]$$

$$w_0 = w_0 - \eta \times \frac{\partial E}{\partial w_0} = w_0 - \eta \times (y-t)$$



The Joy of *Finally* Learning

 www.mldawn.com

 @MLDawn2018

 Mehran Bazargani

 ML Dawn

Course Available at: <https://www.mldawn.com/course/the-birth-of-error-functions/>

Demo's Code: <https://www.mldawn.com/binary-classification-from-scratch-using-numpy/>