

CSCI 1470/2470  
Spring 2022

Ritambhara Singh

January 31, 2024  
Wednesday

Perceptron cond. and Loss Functions

# Deep Learning

Lab1, Homework 1 and  
Quiz 1 out today!

Today's goal – Continue discussion on perceptron and learn about the loss functions

(1) Perceptron learning algorithm

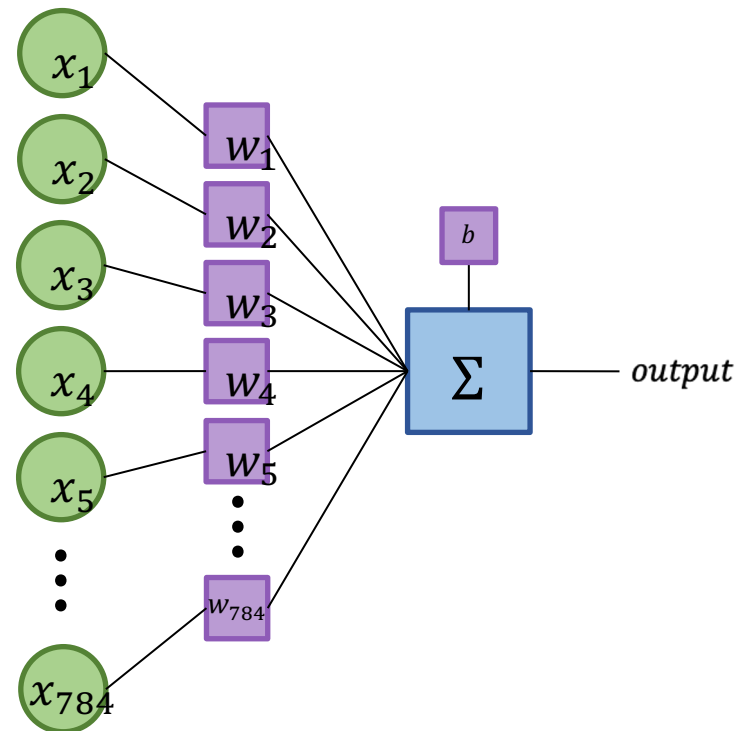
(2) Extending perceptron for Multi-class classification

(3) Loss functions for

- Regression
- Classification

# Recap: A Binary Perceptron for MNIST

- *Inputs*  $[x_1, x_2, \dots, x_n]$  are all positive
  - $n = 784$  (28×28 pixel values)
- *output* is either 0 or 1
  - 0  $\rightarrow$  input is not the digit type we're looking for
  - 1  $\rightarrow$  input *is* the digit type we're looking for

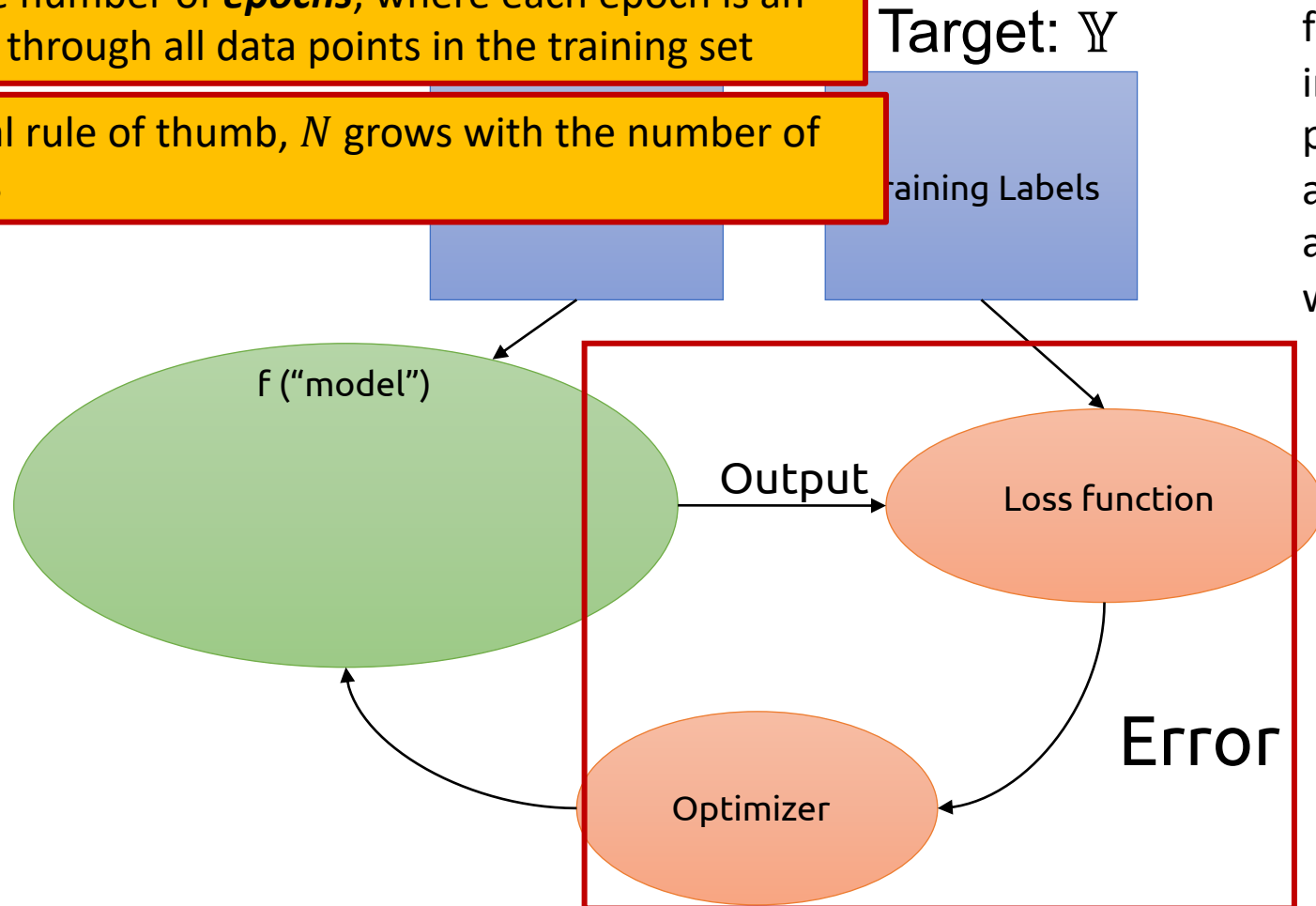




# Training a perceptron

$N$  is known as the number of *epochs*, where each epoch is an iteration of going through all data points in the training set

As a general rule of thumb,  $N$  grows with the number of parameters



0. set the parameters  $\Phi = \{w \cup b\}$  to 0
1. Iterate over training set several times, feeding in each training example into the model, producing an output, and adjusting the parameters according to whether that output was right or wrong

2. Stop once we either (a) get every training example right or (b) after  $N$  iterations, a number set by the programmer.

# The Perceptron Learning Algorithm

1. set  $w$ 's to 0.
2. for  $N$  iterations, or until the weights do not change:
  - a) for each training example  $\mathbf{x}^k$  with label  $y^k$ 
    - i. if  $y^k - f(\mathbf{x}^k) = 0$  continue
    - ii. else for all weights  $w_i$ ,  $\Delta w_i = (y^k - f(\mathbf{x}^k)) x_i^k$

- 
- $b$  = bias
  - $w$  = weights
  - $N$  = maximum number of training iterations
  - $\mathbf{x}^k$  =  $k^{\text{th}}$  training example
  - $y^k$  = label for the  $k^{\text{th}}$  example
  - $w_i$  = weight for the  $i^{\text{th}}$  input where  $i \leq n$
  - $n$  = number of pixels per image
  - $x_i^k$  =  $i^{\text{th}}$  input of the example where  $i \leq n$

# The Perceptron Learning Algorithm

1. set  $w$ 's to 0.
2. for  $N$  iterations, or until the weights do not change:
  - a) for each training example  $\mathbf{x}^k$  with label  $y^k$ 
    - i. **if  $y^k - f(\mathbf{x}^k) = 0$  continue**
    - ii. else for all weights  $w_i$ ,  $\Delta w_i = (y^k - f(\mathbf{x}^k)) x_i^k$

- 
- If the output of our model matches the label, we continue
  - If the correct label is 1, and our output is 1,  $1 - 1 = 0$
  - If the correct label is 0, and our output is 0,  $0 - 0 = 0$

# The Perceptron Learning Algorithm

1. set  $w$ 's to 0.
2. for  $N$  iterations, or until the weights do not change:
  - a) for each training example  $\mathbf{x}^k$  with label  $y^k$ 
    - i. if  $y^k - f(\mathbf{x}^k) = 0$  continue
    - ii. else for all weights  $w_i, \Delta w_i = (y^k - f(\mathbf{x}^k)) x_i^k$**

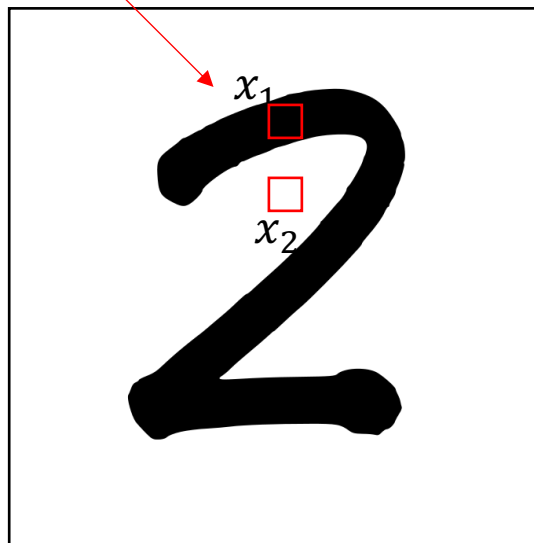
- 
- If our label  $y^k$  is a 1, and our model's output is a 0, we update the  $i^{th}$  weight by:
    - $(1 - 0) \cdot x_i^k = x_i^k$
    - Output was 0 and should have been 1, so make the output more positive
  - If our label  $y^k$  is a 0, and our model's output is a 1, we update the  $i^{th}$  weight by:
    - $(0 - 1) \cdot x_i^k = -x_i^k$
    - Output was 1 and should have been 0, so make the output more negative

Example: Predict whether a digit is a “2”



# Predict whether a digit is a “2”

Just look at the effect of these two pixels



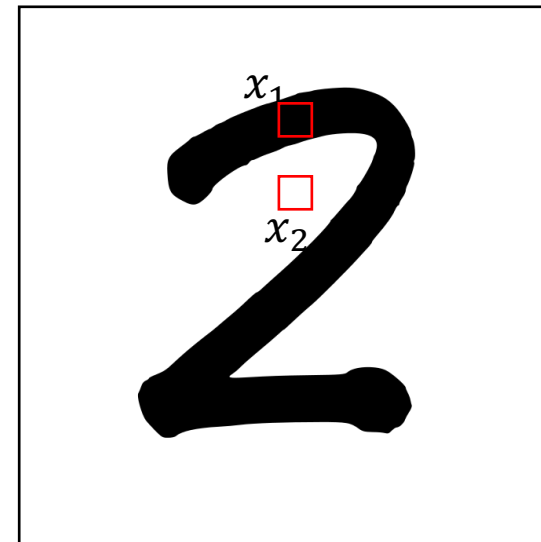
$$x_1 = 0.8$$

$$x_2 = 0$$

# Predict whether a digit is a “2”

- Start off training with all parameters as 0, so  $w_1 = 0$ ,  $w_2 = 0$ , and  $b = 0$
- $f(x) = (w_1 \cdot x_1 + w_2 \cdot x_2 + b \cdot 1)$
- $f(x) = (0 \cdot 0.8 + 0 \cdot 0 + 0 \cdot 1) = 0$ 
  - Return 0 because value is not greater than 0
- Predict that it is not a 2!
- Correct answer: it is a 2...
- Parameter update:
  - $\Delta w_1 = (1 - 0) \cdot 0.8 = 0.8$
  - $\Delta w_2 = (1 - 0) \cdot 0 = 0$
  - $\Delta b = (1 - 0) \cdot 1 = 1$
- Now
  - $w_1 = 0.8$
  - $w_2 = 0$
  - $b = 1$

True label = 1



$$x_1 = 0.8$$

$$x_2 = 0$$

# Predict whether a digit is a “2”

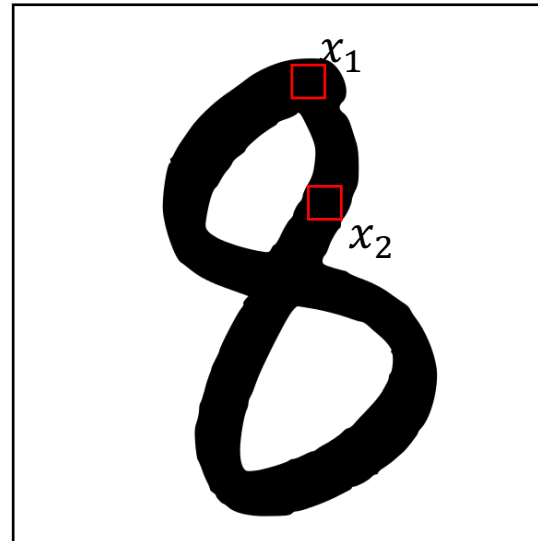
- Next example:

Remember the starting weights are now:

$$w_1 = 0.8$$

$$w_2 = 0$$

$$b = 1$$



$$x_1 = 0.9$$

$$x_2 = 0.9$$

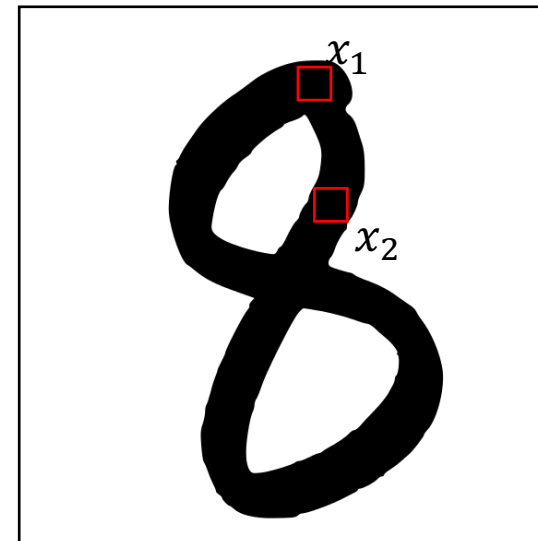
# Predict whether a digit is a “2”

- At end of last iteration:
  - $w_1 = 0.8, w_2 = 0,$  and  $b = 1$
- $f(x) = (w_1 \cdot x_1 + w_2 \cdot x_2 + b \cdot 1)$
- $f(x) = (0.8 \cdot 0.9 + 0 \cdot 0.9 + 1 \cdot 1) > 0$ 
  - Return 1 because value is greater than 0
- Predict that it is a 2!
- Correct answer: it is not a 2...
- Parameter update:
  - $\Delta w_1 = (0 - 1) \cdot 0.9 = -0.9$
  - $\Delta w_2 = (0 - 1) \cdot 0.9 = -0.9$
  - $\Delta b = (0 - 1) \cdot 1 = -1$
- Now
  - $w_1 = 0.8 - 0.9 = -0.1$
  - $w_2 = 0 - 0.9 = -0.9$
  - $b = 1 - 1 = 0$

Any questions?



True label = 0



$$x_1 = 0.9$$

$$x_2 = 0.9$$

# Multi-class problem


# Bringing back the complexity

Input:  $\mathbb{X}$

Classifying MNIST digits requires predicting 1 of 10 possible values

Target:  $\mathbb{Y}$

Pixel Grid

$x^{(1)} =$    
28x28 pixels

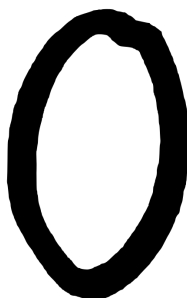
→ Function:  $f$  →

$f(\mathbb{X}) \rightarrow \mathbb{Y}$

Which digit is it?

$y^{(1)} = \text{"2"}$

How do we do that?

$x^{(2)} =$  

Rather than predicting whether a handwritten digit is of a certain class, we predict the class it is most likely in

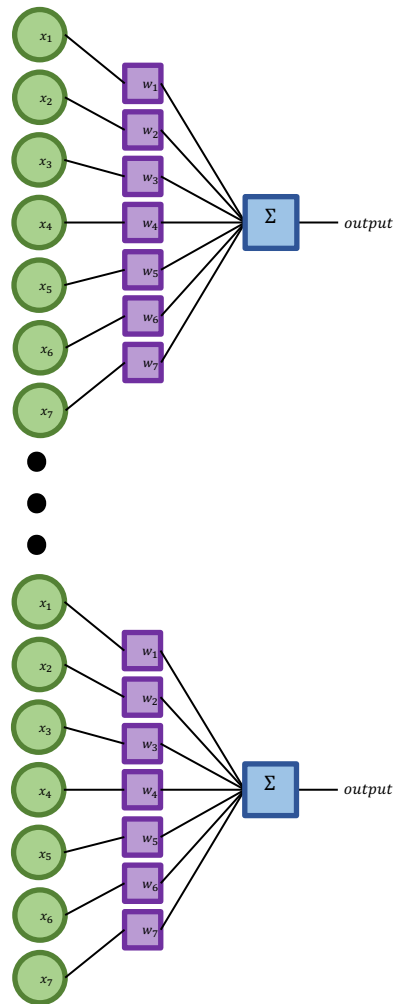
$y^{(2)} = \text{"0"}$



# Using multiple perceptrons

- We can extend perceptrons to multi-class problems by creating  $m$  perceptrons, where  $m =$  the number of classes
- For MNIST, we would have 10 perceptrons
- Each individual perceptron returns a value, so our model will return the class whose perceptron value is the highest.
  - Here, “perceptron value” refers to the value of the weighted sum before being thresholded.

# Using multiple perceptrons

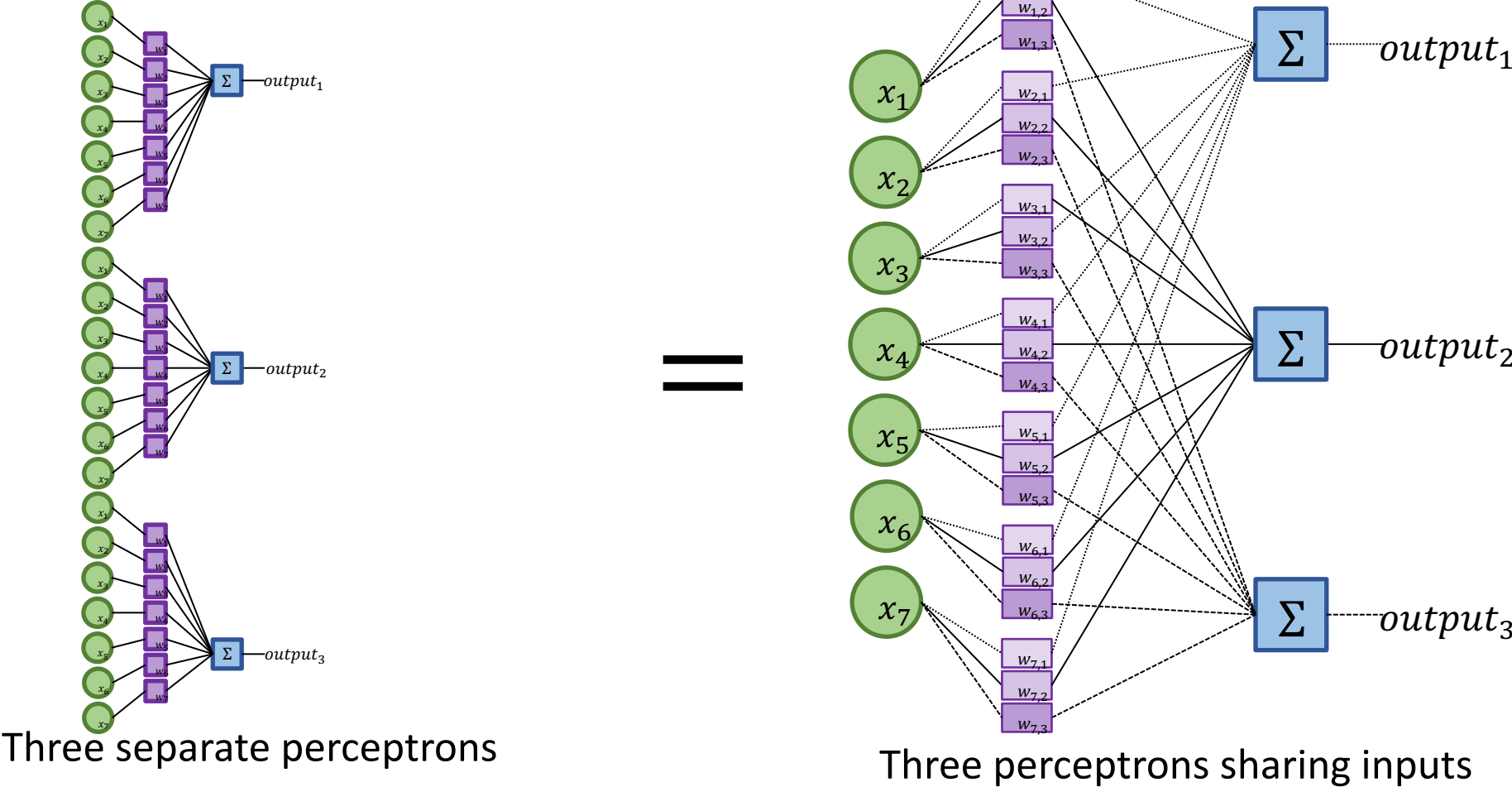


Perceptron for predicting whether handwritten digit is a 0

•  
•  
•

Perceptron for predicting whether handwritten digit is a 9

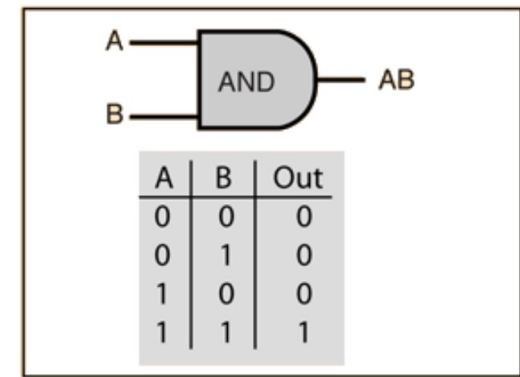
# Multi-class Perceptron



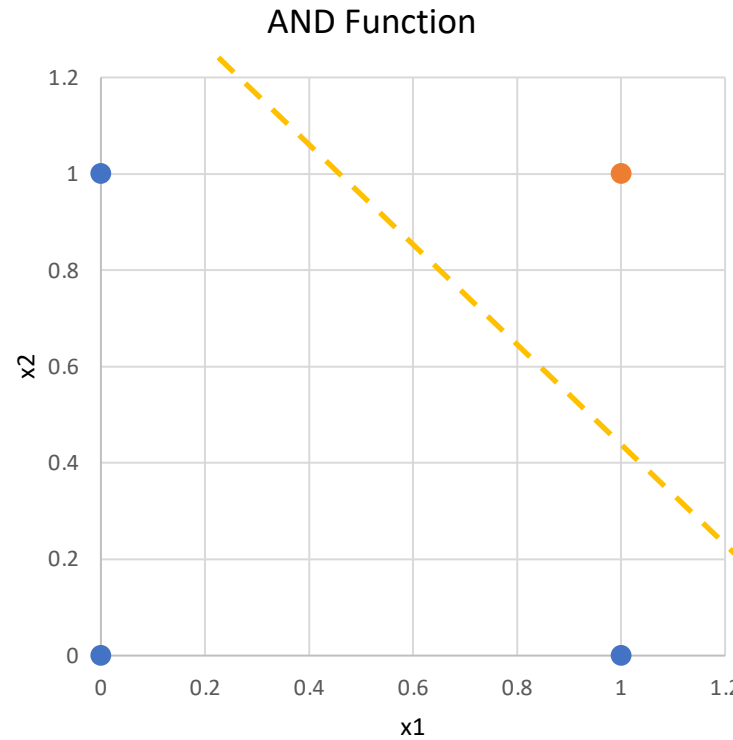
Is there anything the Perceptron can't learn?

# AND Function

Perceptrons work well with linearly separable data



AND Gate

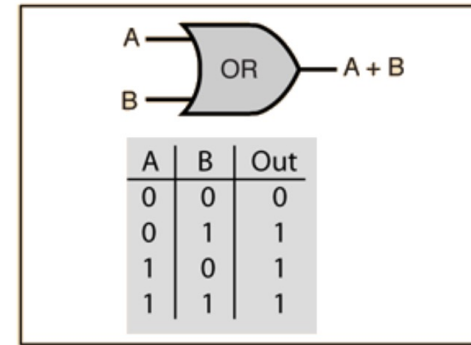


- Output = 0
- Output = 1

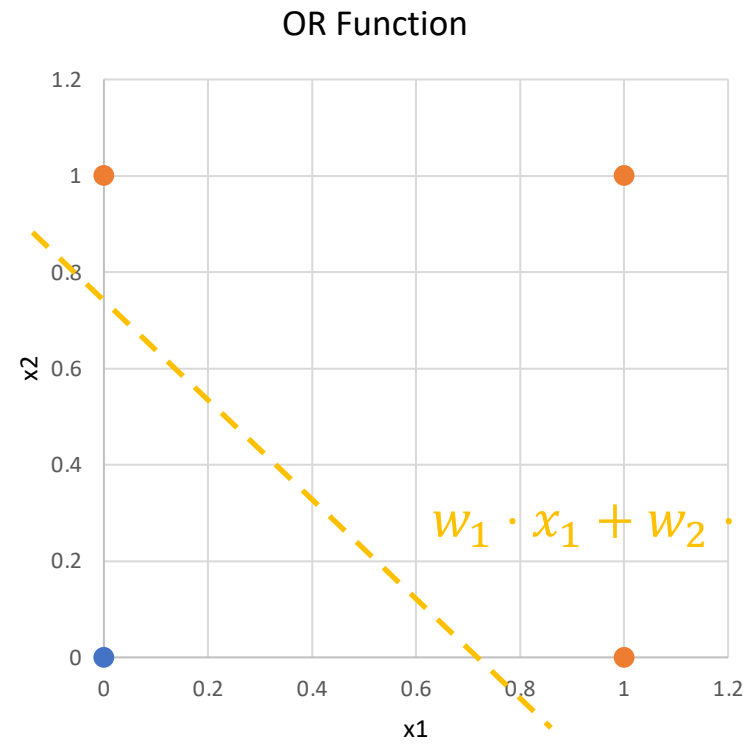
$$w_1 \cdot x_1 + w_2 \cdot x_2 + b > 0?$$

Linear decision boundary

# OR Function



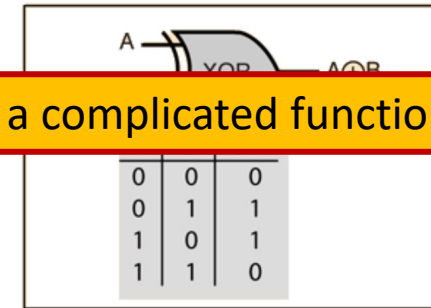
OR Gate



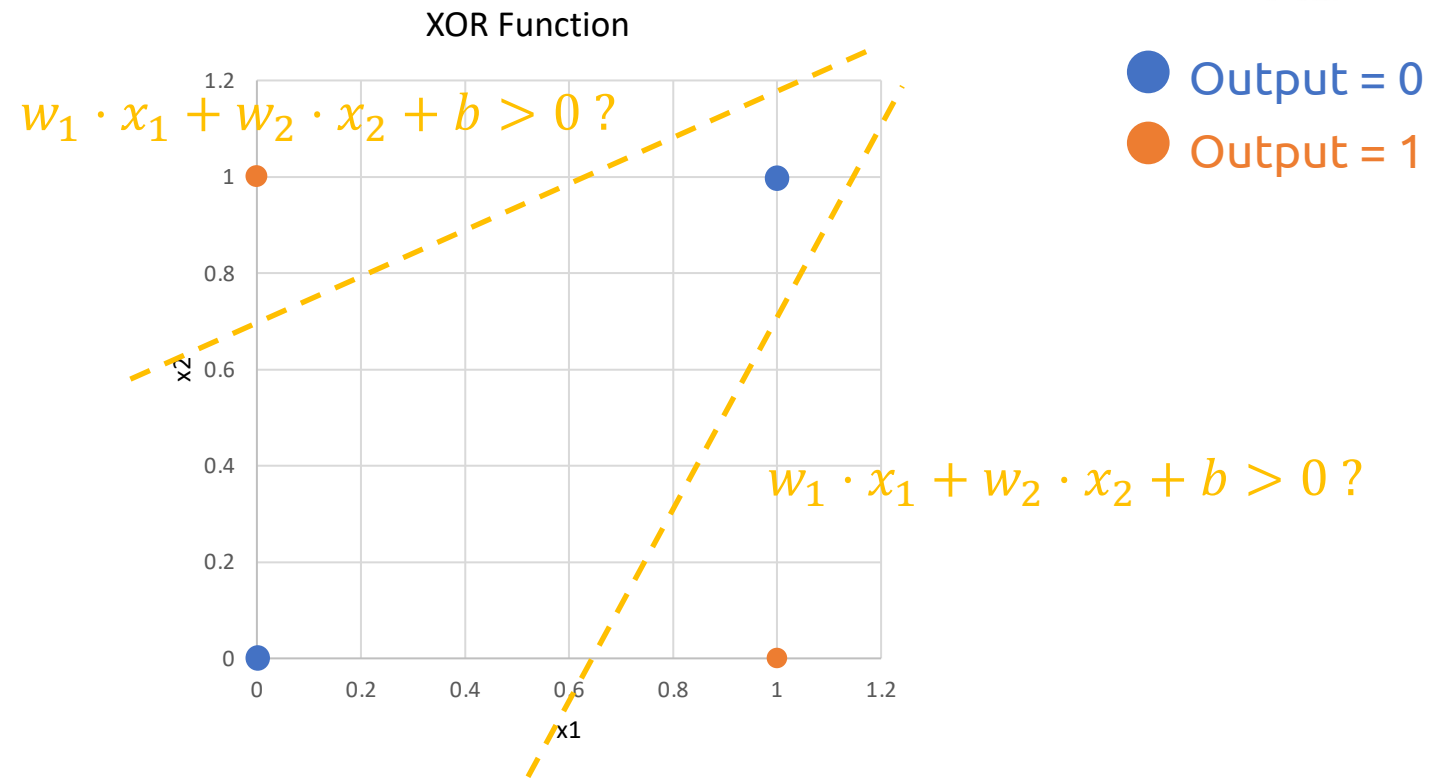
- Output = 0
- Output = 1

# XOR Function

Complicated data would need a complicated function!



XOR Gate



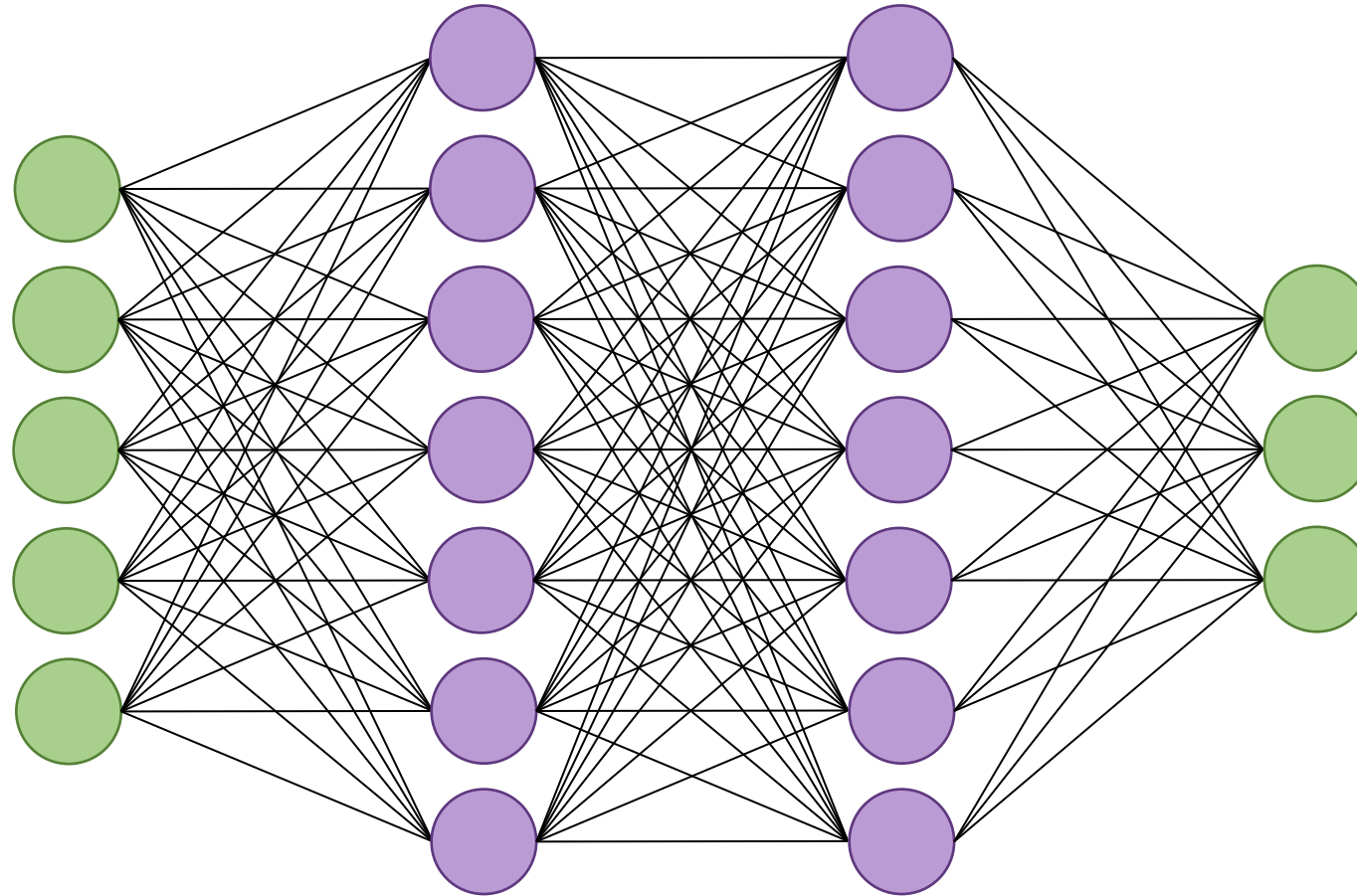
Need **two** linear decision boundaries to represent this function...



# Multi-Layered Neural Net

May see the term Multi-Layer Perceptron (MLP), HOWEVER "perceptrons" are not perceptrons in the strictest possible sense

We really don't use the threshold function of a perceptron but still use the linear function



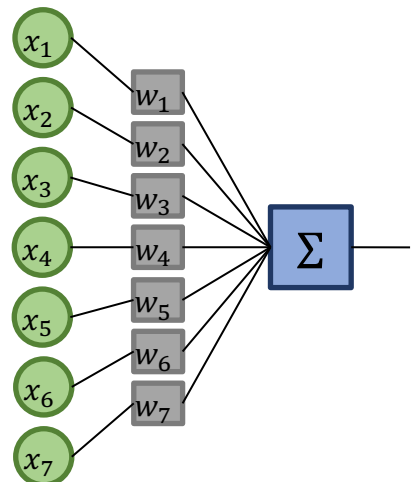
A Multi-Layered Neural Net

Any questions?



# How do we train multi-layer networks?

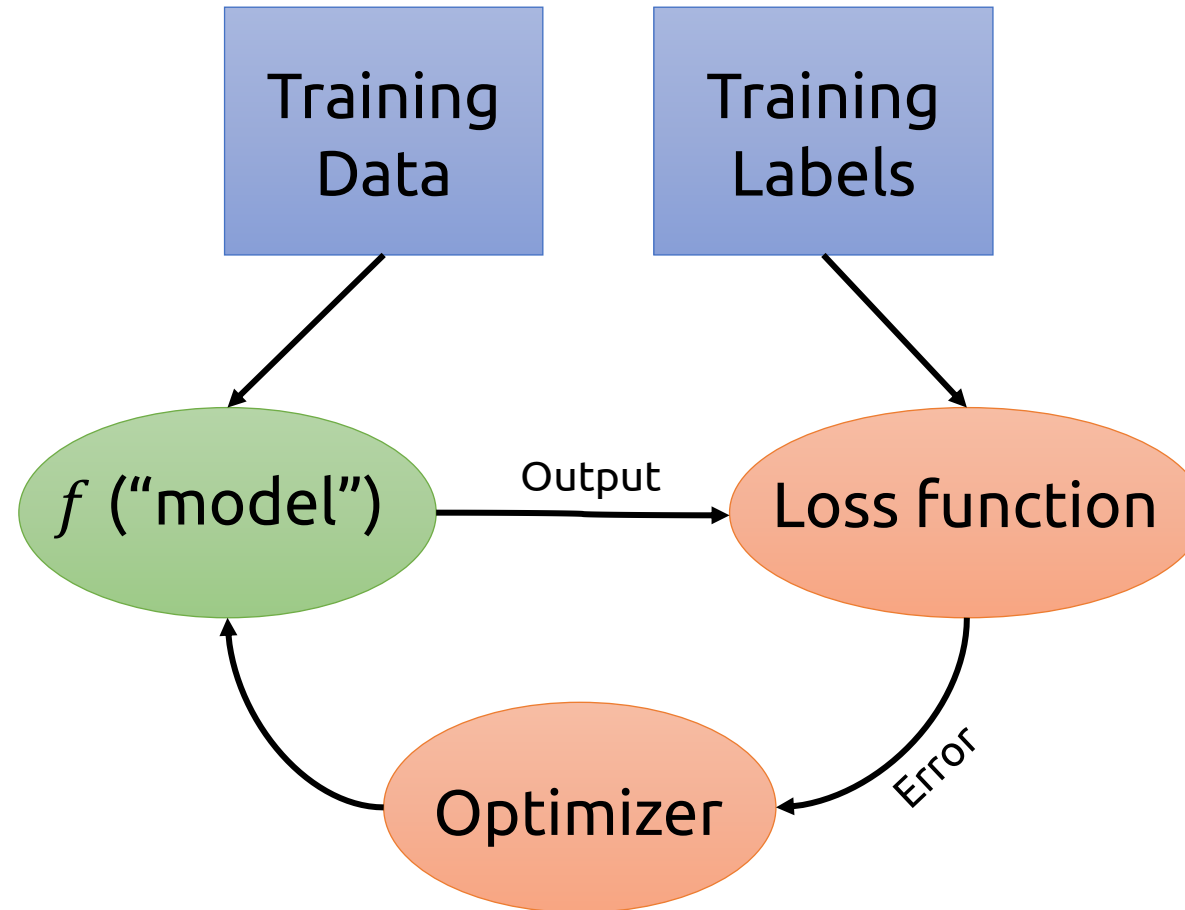
- Unfortunately, the perceptron algorithm doesn't generalize beyond the one-layer case... 😞
- We need a new algorithm...



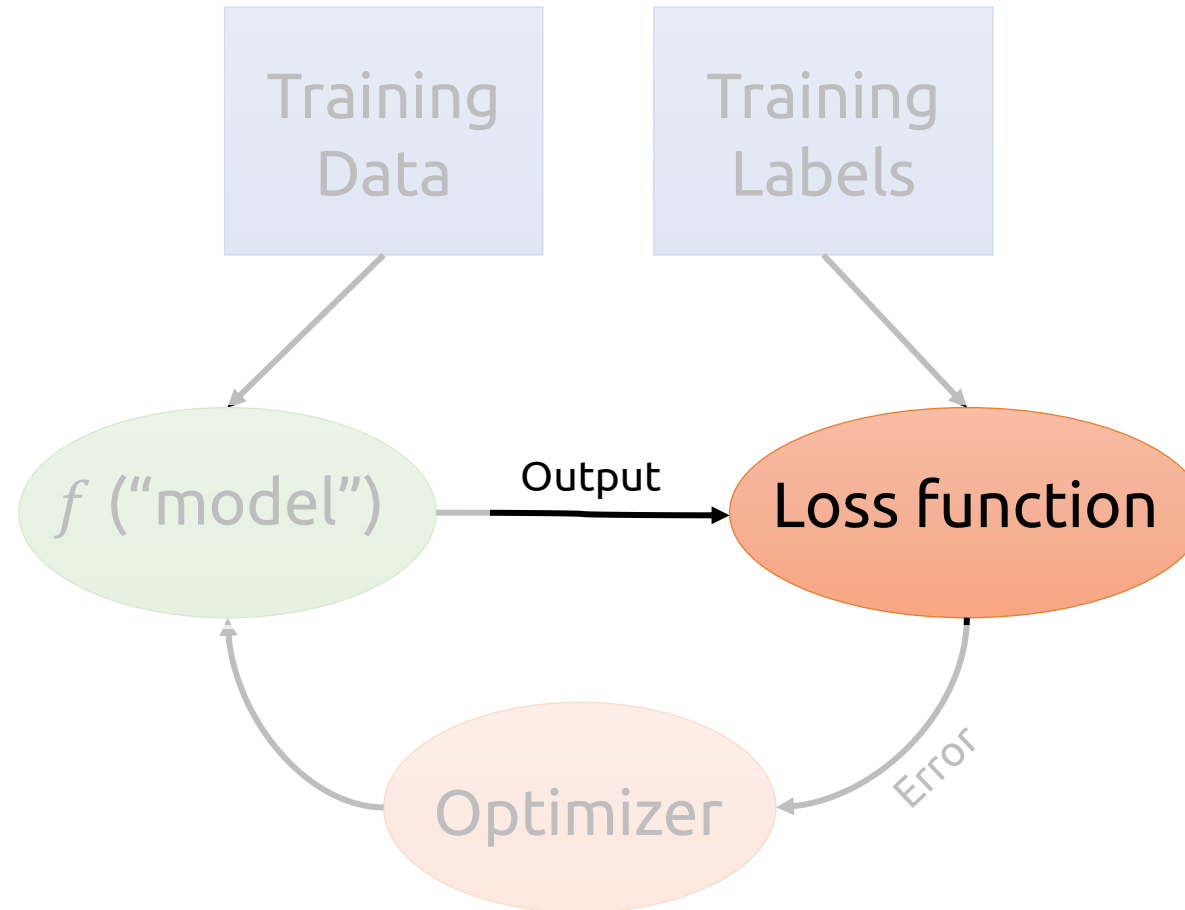
A one-layer binary perceptron

1. set  $w$ 's to 0.
2. for  $N$  iterations, or until the weights do not change:
  - a) for each training example  $\mathbf{x}^k$  with label  $a^k$ 
    - i. if  $a^k - f(\mathbf{x}^k) = 0$  continue
    - ii. else for all weights  $w_i$ ,  $\Delta w_i = (a^k - f(\mathbf{x}^k)) x_i$

# A critical ingredient for our new approach: Loss functions



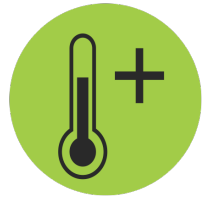
# A critical ingredient for our new approach: Loss functions



# What is a Loss Function?

- A function  $L$  which measures how “wrong” a network is
- $L$  is computed by comparing two values (predicted and true) that shows which is better
- Evaluate  $L$  and update parameters **to decrease  $L$** , making the network “less wrong”

# Recap – regression task



Input:  $X$   
"Temperature"

$$x^{(1)} = 100.1$$

$$X \in \mathbb{R}$$

$$x^{(2)} = 60.0$$

$$x^{(3)} = 30.3$$

Regression

Function:  $f$

$$f(X) \rightarrow Y$$

Target:  $Y$

"Profit made on selling  
lemonade"

$$y^{(1)} = 200.0$$

$$y^{(2)} = 160.5$$

$$y^{(3)} = 145.1$$

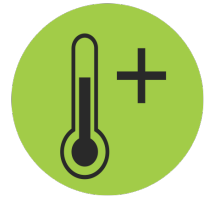
$$Y \in \mathbb{R}$$

(Numerical output)





# Recap: Learning function $f$



Input:  $\mathbb{X}$   
"Temperature"

$$x^{(1)} = 100.1$$

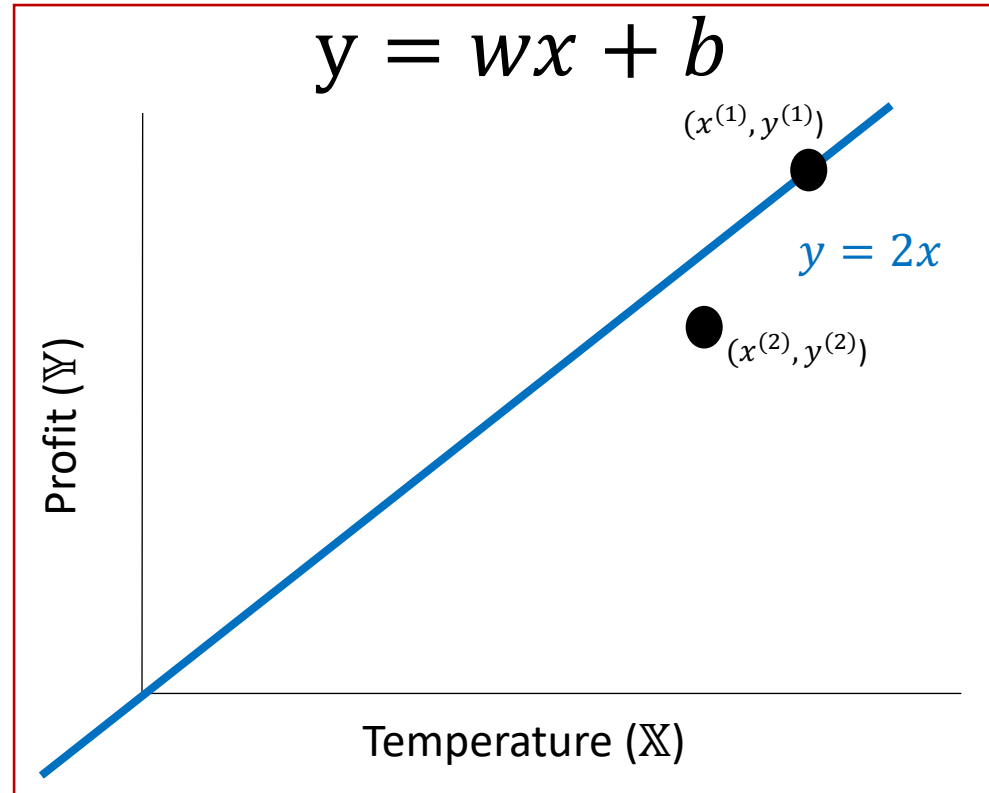
$$x^{(2)} = 60.0$$

$$x^{(3)} = 30.3$$

$\mathbb{X} \in \mathbb{R}$

What could be our loss function?

Linear function



Target:  $\mathbb{Y}$

"Profit made on selling lemonade"



$$y^{(1)} = 200.0$$

$$y^{(2)} = 160.5$$

$$y^{(3)} = 145.1$$

$\mathbb{Y} \in \mathbb{R}$   
(Numerical output)

# Mean Squared Error (MSE)

Average squared residual (residual: difference between predicted and true value)

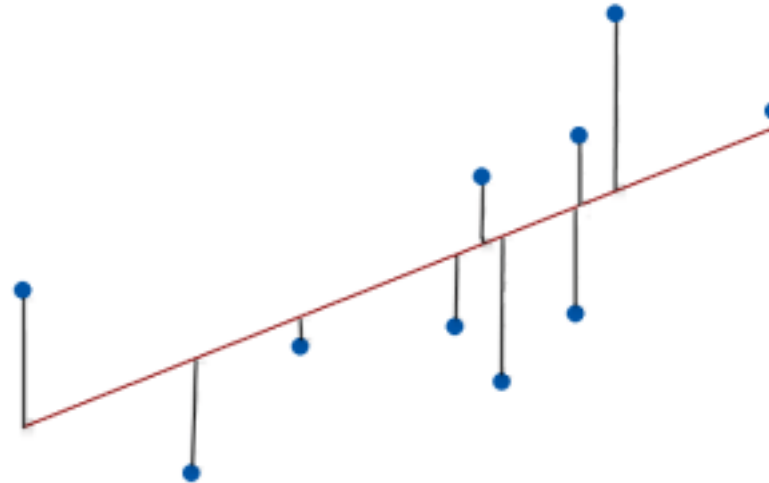
Decreasing the MSE = the model has less error = data points fall closer to the regression line

$$MSE = \frac{\sum_{k=1}^n (y^k - \hat{y}^k)^2}{n}$$

$y^k$ : true output value

$\hat{y}^k$ : predicted output value

$n$ : number of samples



MSE is the average squared distance between the observed and predicted values

What could be the purpose of squaring the distance?

# Mean Squared Error (MSE)

Average squared residual (residual: difference between predicted and true value)

Decreasing the MSE = the model has less error = data points fall closer to the regression line

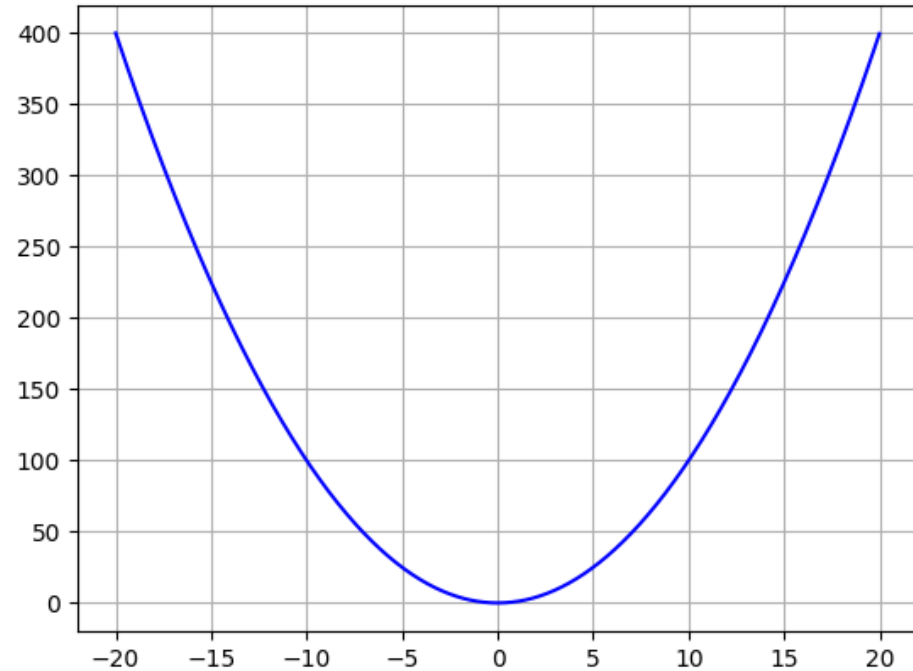
$$MSE = \frac{\sum_{k=1}^n (y^k - \hat{y}^k)^2}{n}$$

$y^k$ : true output value

$\hat{y}^k$ : predicted output value

$n$ : number of samples

What could be the purpose of squaring the distance?



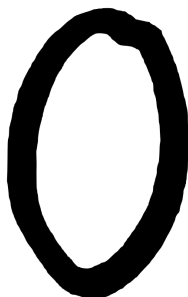
# Recap: Binary classification

Input:  $\mathbb{X}$

Pixel Grid

$x^{(1)} =$  

28x28 pixels

$x^{(2)} =$  

What is a good loss for our binary classification?

We want our network to produce the right answer with high probability

→ Function:  $f$  →

1. Make the network output a probability for class 1 (a value between 0 and 1)

2. Use this probability to compute a loss

Target:  $\mathbb{Y}$

Is it digit 2?

$y^{(1)} =$  "1"

$y^{(2)} =$  "0"

# Cross Entropy Loss (for Binary classification)

$y = \text{true label of class (0 or 1)}$

$p = \text{predicted probability of class 1}$

$\log(p)$

When the true label is 1  
we want higher predicted  
probability for a digit to  
be 1

When the true label is 0  
we want lower predicted  
probability for a digit to  
be 1

Some examples:

$$\log(0.9) = -0.105$$

$$\log(0.5) = -0.693$$

$$\log(0.001) = -6.908$$

# Cross Entropy Loss (for Binary classification)

$y = \text{true label of class (0 or 1)}$

$p = \text{predicted probability of class 1}$

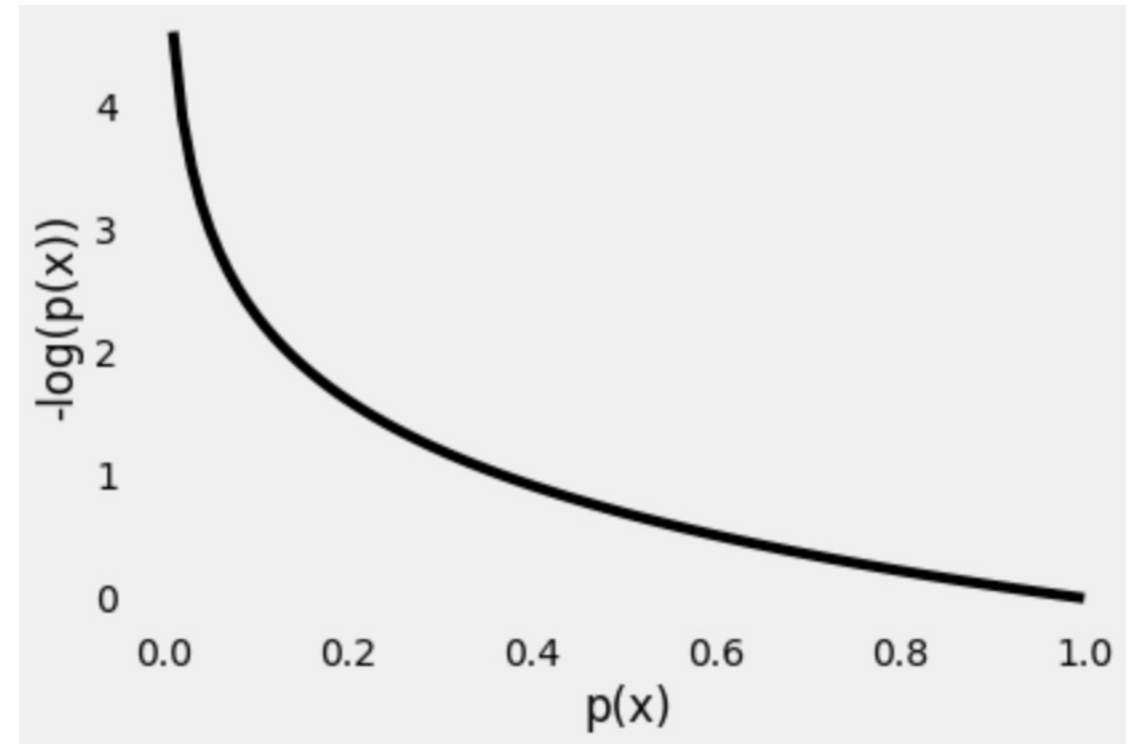
$$-\log(p)$$

Some examples:

$$\log(0.9) = -0.04$$

$$\log(0.5) = -0.3$$

$$\log(0.001) = -3$$



# Cross Entropy Loss (for Binary classification)

$y = \text{true label of class (0 or 1)}$

$p = \text{predicted probability of class 1}$

$$-(y \log(p))$$

Some examples:

$$\log(0.9) = -0.04$$

$$\log(0.5) = -0.3$$

$$\log(0.001) = -3$$

# Cross Entropy Loss (for Binary classification)

$y = \text{true label of class (0 or 1)}$

$p = \text{predicted probability of class 1}$

$$-(y \log(p) + (1 - y) \log(1 - p))$$

$$y = 1, p = 0.9$$

$$y = 0, p = 0.9$$

$$y = 1, p = 0.001$$

Some examples:

$$\log(0.9) = -0.04$$

$$\log(0.5) = -0.3$$

$$\log(0.001) = -3$$

Any questions?



We get this probability by using a Sigmoid function

$$y = 0, p = 0.001$$



# Recap: Multi-class classification

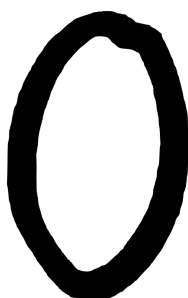
Input:  $\mathbb{X}$

Target:  $\mathbb{Y}$

Pixel Grid

$x^{(1)} =$  

28x28 pixels

$x^{(2)} =$  

We want our network to produce the right answer with high probability

→ Function:  $f$  →

Which digit is it?

$y^{(1)} = \text{"2"}$

1. Make the network output probabilities for all classes (values between 0 and 1)

2. Use these probabilities to compute a loss

$y^{(2)} = \text{"0"}$

# Cross Entropy Loss (for Multi-class classification)

$$-\sum_{j=1}^m y_j \log(p_j)$$

<b>p</b>	Classes (m)	<b>y</b>
0.3	"0"	0
0.2	"1"	0
0.5	"2"	1

2

Some examples:

$$\log(0.9) = -0.04$$

$$\log(0.5) = -0.3$$

$$\log(0.001) = -3$$

We want model to assign high probability to the true class and low to others

# Recap

Perceptron

Perceptron training w/ working example

Multi-class classification

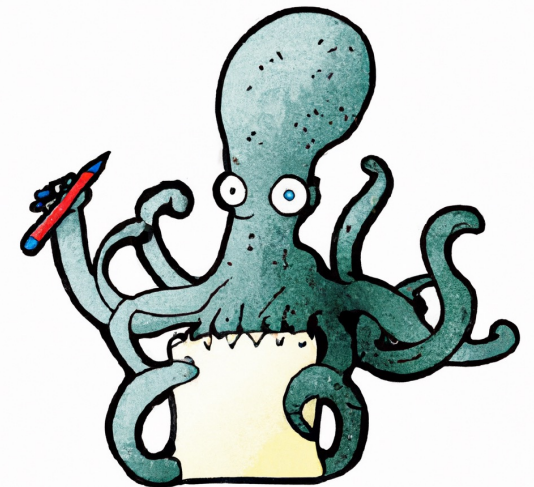
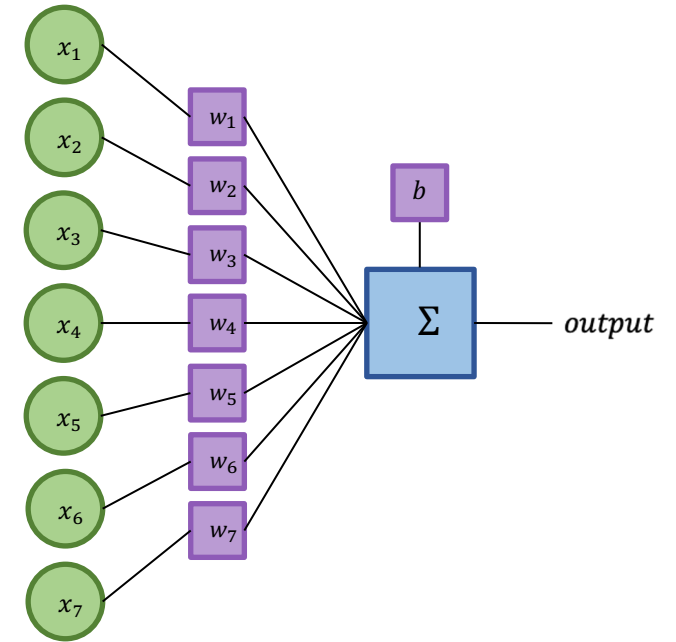
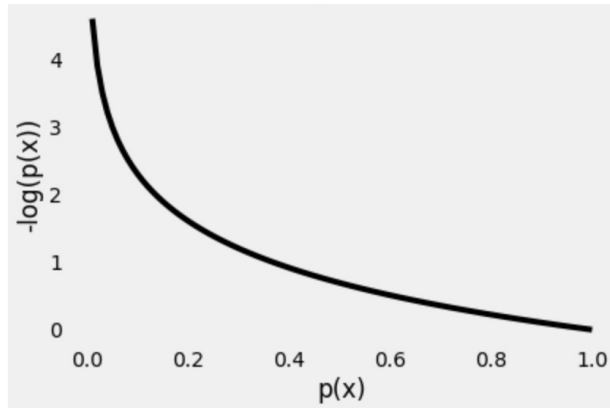
When perceptron fails

Loss functions

MSE loss for regression

Cross entropy loss for binary classification

Cross entropy loss for multi-class classification



# Some Trivia: The Fall of Perceptrons

- In 1969, Marvin Minsky and Seymour Papert released a book, *Perceptrons*, demonstrating that perceptrons are not able to learn the XOR function
- Many earlier researchers heavily focused on logical reasoning, a feature of high-level human cognition, so a machine's ability for logical reasoning was thought to indicate "artificial intelligence"
- Part of a funding battle: Minsky and Papert wanted federal AI funding to go to their kind of 'symbolic' AI, not the early neural net folks...