# Deep Learning

CSCI 1470

Eric Ewing

Monday,
1/27/25

Day 3: Linear Regression, Perceptrons, and MNIST

# *B*rown *AI* *S*afety *T*eam
# It's BAIST.

Our mission is to mitigate large-scale risks from advanced AI through **technical research** and effective **policy interventions**.

Apply to our policy and technical programs or become a member at **BAIST.ai**
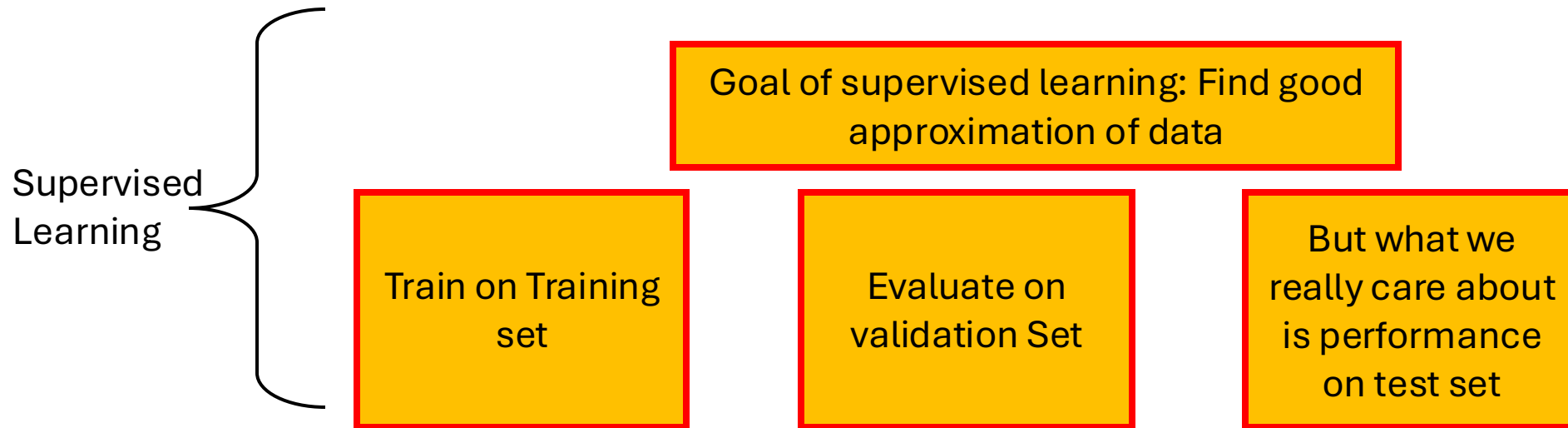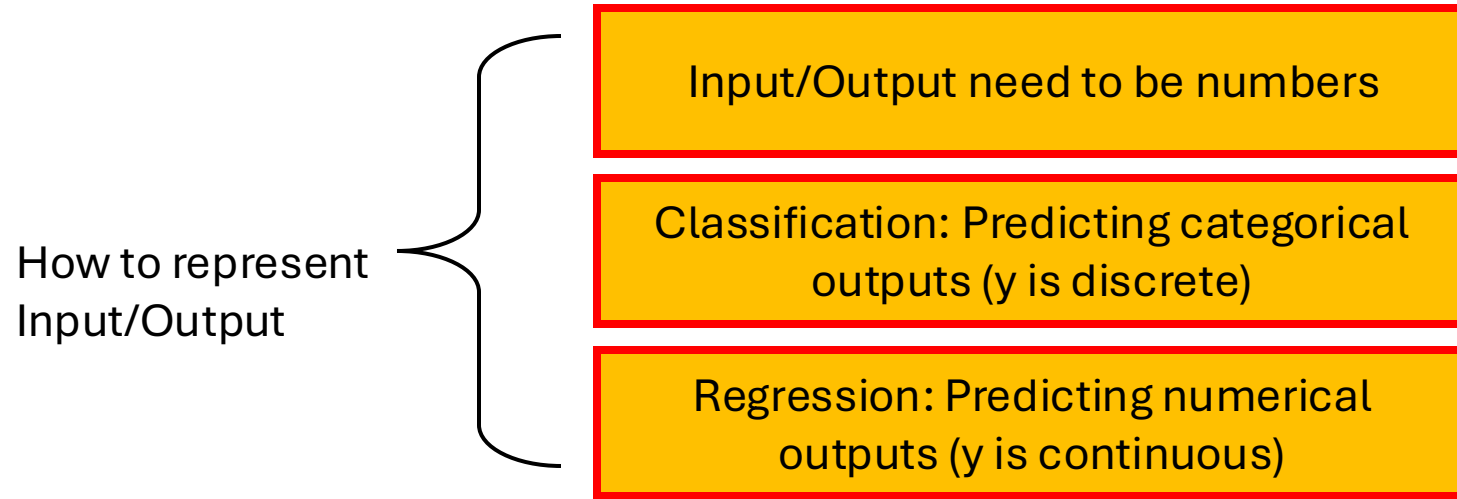
# deepseek

# DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning

DeepSeek-AI

research@deepseek.com

- Much less Supervised Fine Tuning (SFT) than previous models (e.g., GPT4)
- Uses Reinforcement Learning heavily (final part of this course)

# Key Ideas Review

How to represent Input/Output

- Input/Output need to be numbers
- Classification: Predicting categorical outputs (y is discrete)
- Regression: Predicting numerical outputs (y is continuous)

Supervised Learning

- Goal of supervised learning: Find good approximation of data
- Train on Training set
- Evaluate on validation Set
- But what we really care about is performance on test set
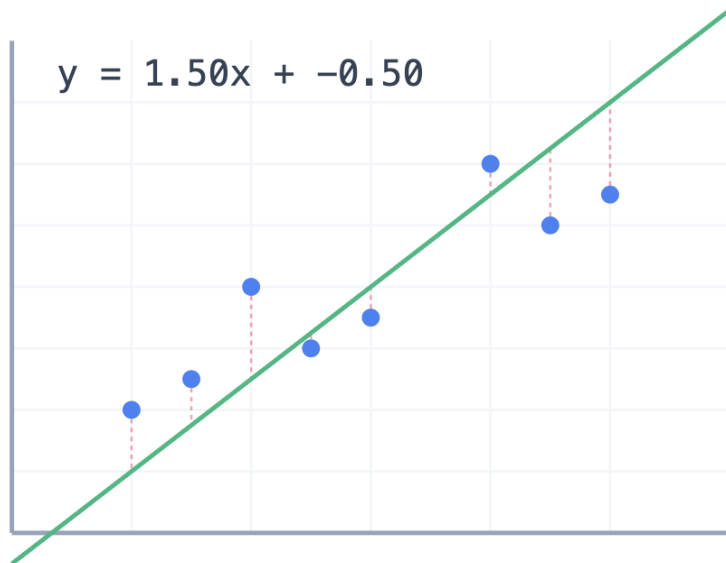
# Today's Goal: Learn about Perceptrons, the first building block of Neural Networks

- Optimization

- Perceptrons

- Introduction to MNIST
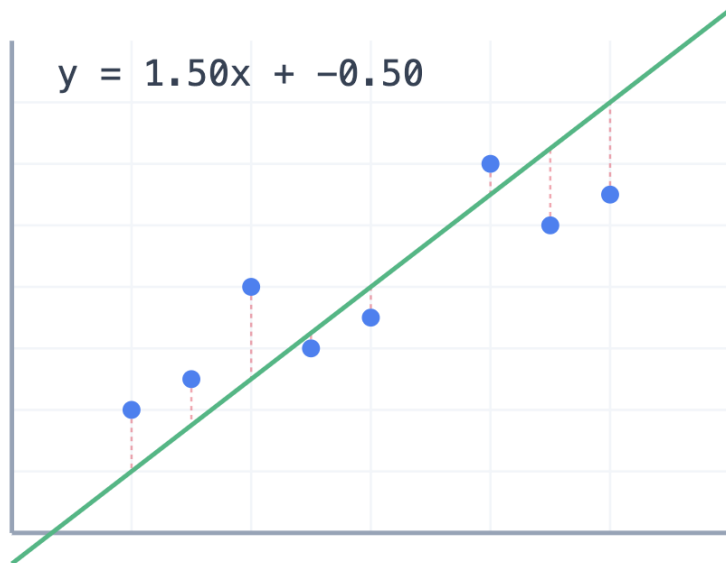
# Linear Regression

$$y = mx + b$$

y = 1.50x + −0.50

With 1 input feature, 2 **parameters**
-   m (slope)
-    b (bias)

# Linear Regression

$$y = mx + b$$

y = 1.50x + −0.50
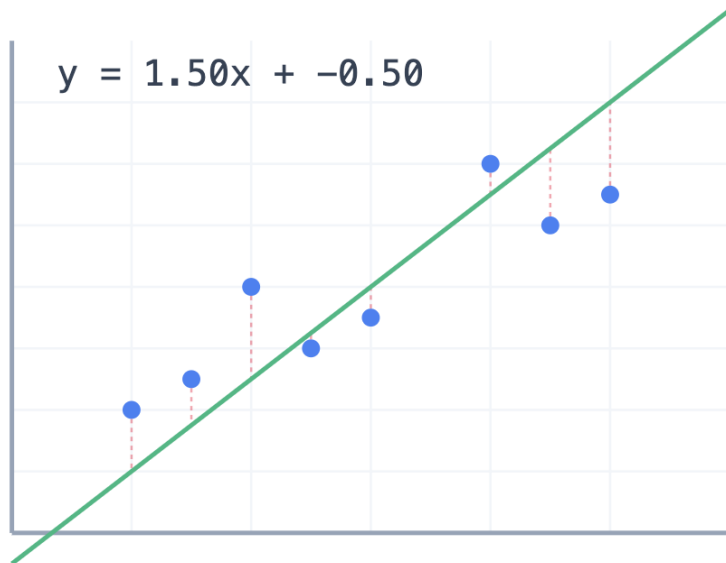
With 1 input feature, 2 **parameters**
- m (slope)
- b (bias)

|  | Input Features | | | Output Target |
| --- | --- | --- | --- | --- |
|  | $x_1$ | $x_2$ | $x_3$ | $y$ |
|  | Temperature | Sunny? | Day of Week | Profit |
| $x^{(1)}$ | 90 | Yes | Sat | $200 |
| $x^{(2)}$ | 80 | No | Mon | $91 |
| $x^{(3)}$ | 62 | No | Wed | $54 |

# Linear Regression

$$y = mx + b$$

y = 1.50x + −0.50

With 1 input feature, 2 **parameters**
- m (slope)
- b (bias)

| | Input Features | | | | Output Target |
| --- | --- | --- | --- | --- | --- |
| | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
| | Temperature | Sunny? | Day of Week | Constant | Profit |
| $x^{(1)}$ | 90 | Yes | Sat | 1 | $200 |
| $x^{(2)}$ | 80 | No | Mon | 1 | $91 |
| $x^{(3)}$ | 62 | No | Wed | 1 | $54 |

# Linear Regression

$$y = mx + b$$

y = 1.50x + −0.50

With 1 input feature, 2 **parameters**
- m (slope)
- b (bias)

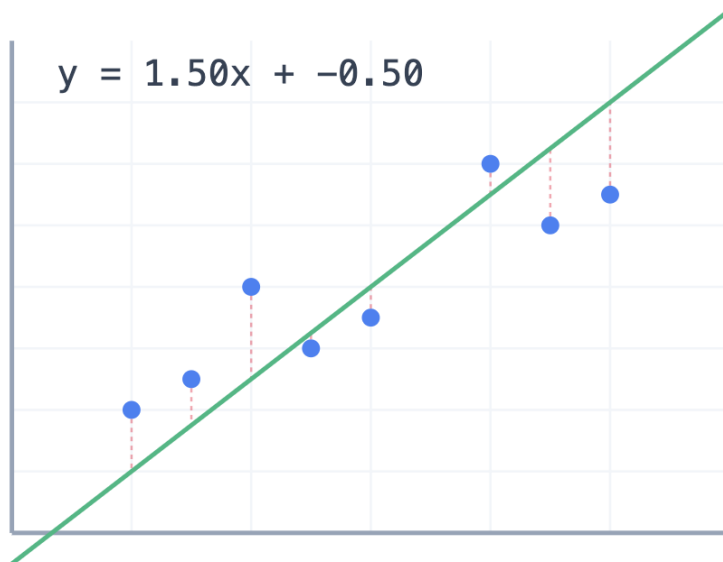|  | Input Features | | | | Output Target |
|---|---|---|---|---|---|
|  | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|  | Temperature | Sunny? | Day of Week | Constant | Profit |
| $x^{(1)}$ | 90 | Yes | Sat | 1 | $200 |
| $x^{(2)}$ | 80 | No | Mon | 1 | $91 |
| $x^{(3)}$ | 62 | No | Wed | 1 | $54 |

With multiple input features:
- Need a weight parameter $w_i$ for each feature $x_i$
- $y = x_1^{(i)} \cdot w_1 + x_2^{(i)} \cdot w_2^{(i)} + \cdots + x_d^{(i)} \cdot w_d$
- Can be rewritten: $y = \vec{x} \cdot \vec{w}$

# Linear Regression

$$y = mx + b$$

y = 1.50x + -0.50

With 1 input feature, 2 **parameters**
- m (slope)
- b (bias)

| | Input Features | | | | Output Target |
|---|---|---|---|---|---|
| | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
| | Temperature | Sunny? | Day of Week | Constant | Profit |
| $x^{(1)}$ | 90 | Yes | Sat | 1 | $200 |
| $x^{(2)}$ | 80 | No | Mon | 1 | $91 |
| $x^{(3)}$ | 62 | No | Wed | 1 | $54 |

With multiple input features:
- Need a weight parameter $w_i$ for each feature $x_i$
- $y = x_1^{(i)} \cdot w_1 + x_2^{(i)} \cdot w_2^{(i)} + \cdots + x_d^{(i)} \cdot w_d$
- Can be rewritten: $y = \vec{x} \cdot \vec{w}$

# Option 1: Closed Form Solution

Goal: Minimize *Loss* function

*Process:*

- Find derivative (or gradient) of loss function

- Set derivative to 0

- Solve for parameters

# Option 1: Closed Form Solution

Goal: Minimize *Loss* function

MSE (Mean Squared Error)

*Process:*

- Find derivative (or gradient) of loss function

- Set derivative to 0

- Solve for parameters

# Option 1: Closed Form Solution

Goal: Minimize *Loss* function

MSE (Mean Squared Error)

*Process:*

- Find derivative (or gradient) of loss function

- Set derivative to 0

- Solve for parameters

Generalization of derivatives to functions with multiple inputs

# Option 1: Closed Form Solution

Goal: Minimize *Loss* function

MSE (Mean Squared Error)

*Process:*

- Find derivative (or gradient) of loss function

Generalization of derivatives to functions with multiple inputs

- Set derivative to 0

- Solve for parameters

Is this guaranteed to find the global best parameter settings?

# Option 1: Closed Form Solution

Goal: Minimize *Loss* function

*Process:*

- Find derivative (or gradient) of loss function

- Set derivative to 0

- Solve for parameters

MSE (Mean Squared Error)

Generalization of derivatives to functions with multiple inputs

Is this guaranteed to find the global best parameter settings?

weight vector $w \in \mathbb{R}^d$

# Gradients

The gradient of a function *f* is a vector of partial derivatives:

$$\nabla f_\theta = [\frac{\partial f}{\theta_1}, \frac{\partial f}{\theta_2}, \frac{\partial f}{\theta_3}, \dots, \frac{\partial f}{\theta_d}]$$

# Gradients

The gradient of a function *f* is a vector of partial derivatives:

$$\nabla f_\theta = [\frac{\partial f}{\theta_1}, \frac{\partial f}{\theta_2}, \frac{\partial f}{\theta_3}, \ldots, \frac{\partial f}{\theta_d}]$$

For a linear regression model with one input variable what dimension is $\nabla f_\theta$ in?

$$\nabla f_\theta \in \mathbb{R}^?$$

# Gradients

The gradient of a function *f* is a vector of partial derivatives:

$$\nabla f_\theta = [\frac{\partial f}{\theta_1}, \frac{\partial f}{\theta_2}, \frac{\partial f}{\theta_3}, \ldots, \frac{\partial f}{\theta_d}]$$

For a linear regression model with one input variable what dimension is $\nabla f_\theta$ in?

$$\nabla f_\theta \in \mathbb{R}^?$$

$\nabla f_w$ tells us what happens to *f* with small adjustments to each parameter *w*

# Option 1: Closed Form Solution

# Option 1: Closed Form Solution

$$\mathcal{L} = MSE = \frac{\sum_{i}^{n}(y_i - \vec{w}^T \vec{x})^\wedge 2}{n}$$

# Option 1: Closed Form Solution

Matrix notation will make our lives easy!
$$\mathbb{X} \in \mathbb{R}^{n \times d}, \mathrm{y} \in \mathbb{R}^n, \vec{\mathrm{w}} \in \mathbb{R}^d$$
Vectors are assumed to be column vectors, i.e., $\mathrm{y} \in \mathbb{R}^{n \times 1}$

$$\mathcal{L} = MSE = \frac{\sum_i^n (y_i - \vec{w}^T \vec{x})\^2}{n}$$

# Option 1: Closed Form Solution

Matrix notation will make our lives easy!
$$\mathbb{X} \in \mathbb{R}^{n \times d}, \mathbb{y} \in \mathbb{R}^{n}, \vec{\mathbb{w}} \in \mathbb{R}^{d}$$
Vectors are assumed to be column vectors, i.e., $\mathbb{y} \in \mathbb{R}^{n \times 1}$

Is this a legal operation: $\mathbb{y} - \vec{w}\mathbb{X}$?

$$\mathcal{L} = MSE = \frac{\sum_{i}^{n}(y_i - \vec{w}^T \vec{x})^{\wedge}2}{n}$$

# Option 1: Closed Form Solution

Matrix notation will make our lives easy!
$$\mathbb{X} \in \mathbb{R}^{n \times d}, \mathbb{y} \in \mathbb{R}^n, \vec{w} \in \mathbb{R}^d$$
Vectors are assumed to be column vectors, i.e., $\mathbb{y} \in \mathbb{R}^{n \times 1}$

$$\mathcal{L} = MSE = \frac{\sum_i^n (y_i - \vec{w}^T \vec{x})^{\wedge}2}{n}$$

Is this a legal operation: $\mathbb{y} - \vec{w}\mathbb{X}$?

Is this a legal operation: $(\mathbb{y} - \mathbb{X}\vec{w})(\mathbb{y} - \mathbb{X}\vec{w})$?

# Option 1: Closed Form Solution

Matrix notation will make our lives easy!

$\mathbb{X} \in \mathbb{R}^{n \times d}, \mathrm{y} \in \mathbb{R}^n, \vec{\mathrm{w}} \in \mathbb{R}^d$

Vectors are assumed to be column vectors, i.e., $\mathrm{y} \in \mathbb{R}^{n \times 1}$

$$\mathcal{L} = MSE = \frac{\sum_i^n (y_i - \vec{w}^T \vec{x})^\wedge 2}{n}$$

Is this a legal operation: $\mathrm{y} - \vec{w}\mathbb{X}$?

Is this a legal operation: $(\mathrm{y} - \mathbb{X}\vec{w})(\mathrm{y} - \mathbb{X}\vec{w})$?

Shape errors are the most common errors you will face when starting deep learning

# Option 1: Closed Form Solution

Matrix notation will make our lives easy!
$$\mathbb{X} \in \mathbb{R}^{n \times d}, \mathbb{y} \in \mathbb{R}^{n}, \vec{\mathbb{w}} \in \mathbb{R}^{d}$$
Vectors are assumed to be column vectors, i.e., $\mathbb{y} \in \mathbb{R}^{n \times 1}$

Is this a legal operation: $\mathbb{y} - \vec{w}\mathbb{X}$?

Is this a legal operation: $(\mathbb{y} - \mathbb{X}\vec{w})(\mathbb{y} - \mathbb{X}\vec{w})$?

Shape errors are the most common errors you will face when starting deep learning

$$\mathcal{L} = MSE = \frac{\sum_{i}^{n}(y_i - \vec{w}^T \vec{x})\char`\^2}{n}$$

$$\mathcal{L} = \frac{(\mathbb{y} - \mathbb{X}\vec{w})^T(\mathbb{y} - \mathbb{X}\vec{w})}{n}$$

# Option 1: Closed Form Solution

Matrix notation will make our lives easy!

$$\mathbb{X} \in \mathbb{R}^{n \times d}, \mathbb{y} \in \mathbb{R}^n, \vec{\mathbb{w}} \in \mathbb{R}^d$$

Vectors are assumed to be column vectors, i.e., $\mathbb{y} \in \mathbb{R}^{n \times 1}$

Is this a legal operation: $\mathbb{y} - \vec{w}\mathbb{X}$?

Is this a legal operation: $(\mathbb{y} - \mathbb{X}\vec{w})(\mathbb{y} - \mathbb{X}\vec{w})$?

Shape errors are the most common errors you will face when starting deep learning

$$\mathcal{L} = MSE = \frac{\sum_i^n (y_i - \vec{w}^T \vec{x})^\wedge 2}{n}$$

$$\mathcal{L} = \frac{(\mathbb{y} - \mathbb{X}\vec{w})^T (\mathbb{y} - \mathbb{X}\vec{w})}{n}$$

$$\mathcal{L} = \frac{\mathbb{y}^T \mathbb{y} - \mathbb{y}^T \mathbb{X}\vec{w} - \vec{w^T}\mathbb{X}^T \mathbb{y} + \vec{w}^T \mathbb{X}^T \mathbb{X}\vec{w}}{n}$$

# Option 1: Closed Form Solution

Matrix notation will make our lives easy!
$$\mathbb{X} \in \mathbb{R}^{n \times d}, \mathbb{y} \in \mathbb{R}^{n}, \vec{\mathbb{w}} \in \mathbb{R}^{d}$$
Vectors are assumed to be column vectors, i.e., $\mathbb{y} \in \mathbb{R}^{n \times 1}$

Is this a legal operation: $\mathbb{y} - \vec{w}\mathbb{X}$?

Is this a legal operation: $(\mathbb{y} - \mathbb{X}\vec{w})(\mathbb{y} - \mathbb{X}\vec{w})$?

Shape errors are the most common errors you will face when starting deep learning

$$\mathcal{L} = MSE = \frac{\sum_{i}^{n}(y_i - \vec{w}^T \vec{x})^{\wedge}2}{n}$$

$$\mathcal{L} = \frac{(\mathbb{y} - \mathbb{X}\vec{w})^T(\mathbb{y} - \mathbb{X}\vec{w})}{n}$$

$$\mathcal{L} = \frac{\mathbb{y}^T\mathbb{y} - \mathbb{y}^T\mathbb{X}\vec{w} - \vec{w}^T\mathbb{X}^T\mathbb{y} + \vec{w}^T\mathbb{X}^T\mathbb{X}\vec{w}}{n}$$

$$\nabla\mathcal{L}_w = \frac{-2\mathbb{X}^T\mathbb{y} + 2\mathbb{X}^T\mathbb{X}\vec{w}}{n}$$

# Option 1: Closed Form Solution

Matrix notation will make our lives easy!
$$\mathbb{X} \in \mathbb{R}^{n \times d}, \mathbb{y} \in \mathbb{R}^n, \vec{\mathbb{w}} \in \mathbb{R}^d$$
Vectors are assumed to be column vectors, i.e., $\mathbb{y} \in \mathbb{R}^{n \times 1}$

Is this a legal operation: $\mathbb{y} - \vec{w}\mathbb{X}$?

Is this a legal operation: $(\mathbb{y} - \mathbb{X}\vec{w})(\mathbb{y} - \mathbb{X}\vec{w})$?

Shape errors are the most common errors you will face when starting deep learning

$$\mathcal{L} = MSE = \frac{\sum_i^n (y_i - \vec{w}^T \vec{x})\wedge 2}{n}$$

$$\mathcal{L} = \frac{(\mathbb{y} - \mathbb{X}\vec{w})^T (\mathbb{y} - \mathbb{X}\vec{w})}{n}$$

$$\mathcal{L} = \frac{\mathbb{y}^T\mathbb{y} - \mathbb{y}^T\mathbb{X}\vec{w} - \vec{w}^T\mathbb{X}^T\mathbb{y} + \vec{w}^T\mathbb{X}^T\mathbb{X}\vec{w}}{n}$$

$$\nabla\mathcal{L}_w = \frac{-2\mathbb{X}^T\mathbb{y} + 2\mathbb{X}^T\mathbb{X}\vec{w}}{n}$$

$$0 = -\mathbb{X}^T\mathbb{y} + \mathbb{X}^T\mathbb{X}\vec{w}$$

# Option 1: Closed Form Solution

Matrix notation will make our lives easy!

$$\mathbb{X} \in \mathbb{R}^{n \times d}, \mathbb{y} \in \mathbb{R}^n, \vec{\mathbb{w}} \in \mathbb{R}^d$$

Vectors are assumed to be column vectors, i.e., $\mathbb{y} \in \mathbb{R}^{n \times 1}$

Is this a legal operation: $\mathbb{y} - \vec{w}\mathbb{X}$?

Is this a legal operation: $(\mathbb{y} - \mathbb{X}\vec{w})(\mathbb{y} - \mathbb{X}\vec{w})$?

Shape errors are the most common errors you will face when starting deep learning

$$\mathcal{L} = MSE = \frac{\sum_i^n (y_i - \vec{w}^T \vec{x}) \char`\^ 2}{n}$$

$$\mathcal{L} = \frac{(\mathbb{y} - \mathbb{X}\vec{w})^T (\mathbb{y} - \mathbb{X}\vec{w})}{n}$$

$$\mathcal{L} = \frac{\mathbb{y}^T \mathbb{y} - \mathbb{y}^T \mathbb{X}\vec{w} - \vec{w}^T \mathbb{X}^T \mathbb{y} + \vec{w}^T \mathbb{X}^T \mathbb{X}\vec{w}}{n}$$

$$\nabla \mathcal{L}_w = \frac{-2\mathbb{X}^T \mathbb{y} + 2\mathbb{X}^T \mathbb{X}\vec{w}}{n}$$

$$0 = -\mathbb{X}^T \mathbb{y} + \mathbb{X}^T \mathbb{X}\vec{w}$$

$$(\mathbb{X}^T \mathbb{X})^{-1} (\mathbb{X}^T \mathbb{y}) = \vec{w}$$

# Closed Form Solution

Advantages:

- Simple/fast to implement

Disadvantages:

- Need to invert: $(\mathbb{X}\mathbb{X}^T)^{-1}$

- Matrix inversion is $O(n^3)$

- $(\mathbb{X}\mathbb{X}^T)$  May not be invertible

- Doesn't necessarily exist for other models

# A Linear Classification Model

# A Linear Classification Model

Linear Regression is a linear model for *regression*.

What's a natural way to make a linear *classifier*?

# A Classifier

Everything above the line (or hyperplane in >2D) is classified as 1, everything below the line as 0



$$\hat{y} = w_1 \cdot x_1 + w_2 \cdot x_2 + b$$

# A Classifier

Everything above the line (or hyperplane in >2D) is classified as 1, everything below the line as 0

How can you tell if a point is above or below the line?

$$\hat{y} = w_1 \cdot x_1 + w_2 \cdot x_2 + b$$

# A Classifier

Everything above the line (or hyperplane in >2D) is classified as 1, everything below the line as 0

How can you tell if a point is above or below the line?

If $\hat{y} = 0$, the point is **on** the line,
If $\hat{y} > 0$, the point is **"above"** the line,
If $\hat{y} < 0$, the point is **"below"** the line

$x_2$

$x_1$

$$\hat{y} = w_1 \cdot x_1 + w_2 \cdot x_2 + b$$

# A Classifier

Everything above the line (or hyperplane in >2D) is classified as 1, everything below the line as 0

How can you tell if a point is above or below the line?

If $\hat{y} = 0$, the point is **on** the line,
If $\hat{y} > 0$, the point is **"above"** the line,
If $\hat{y} < 0$, the point is **"below"** the line

If $\hat{y} > 0$, predict 1.
If $\hat{y} \leq 0$, predict 0.

$x_2$

$x_1$

$$\hat{y} = w_1 \cdot x_1 + w_2 \cdot x_2 + b$$

# Perceptrons: A Linear Classifier

(Our first building block of Deep Learning)

# Biological Motivation

- Loosely inspired by neurons, basic working unit of the brain
- Serve to transmit information between cells

# The Perceptron



Biological Neuron

Artificial Neuron (Perceptron)

# Inputs

Inputs are $\vec{x} = [x_1, x_2, \ldots, x_d]$

Features of the data

# Predicting with a Perceptron

1. Take each of the inputs and multiply by corresponding weight
2. Sum the results, add bias term

# Predicting with a Perceptron

1. Take each of the inputs and multiply by corresponding weight
2. Sum the results, add bias term

Until here, a Perceptron and Linear Regression are equivalent

# Predicting with a Perceptron

1. Take each of the inputs and multiply by corresponding weight
2. Sum the results, add bias term
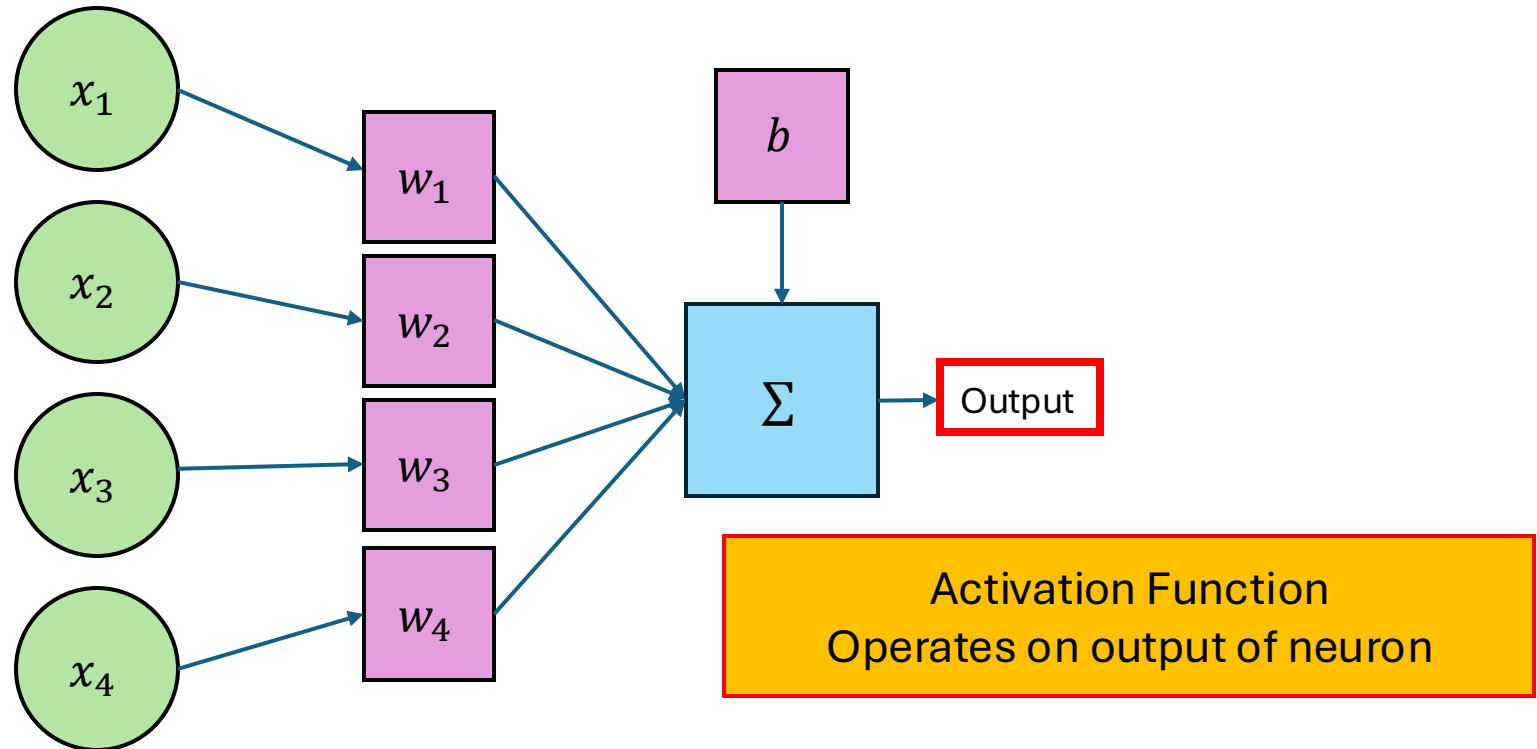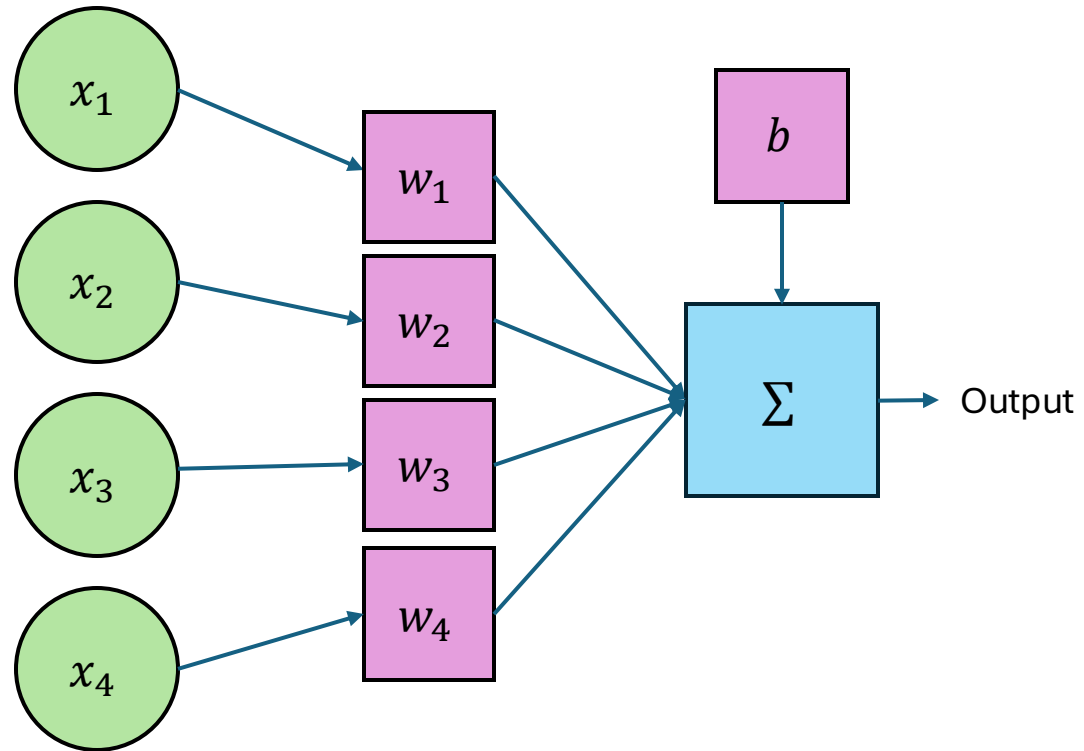3. If output is above 0, return 1, otherwise return 0

# Predicting with a Perceptron

1. Take each of the inputs and multiply by corresponding weight
2. Sum the results, add bias term
3. If output is above 0, return 1, otherwise return 0
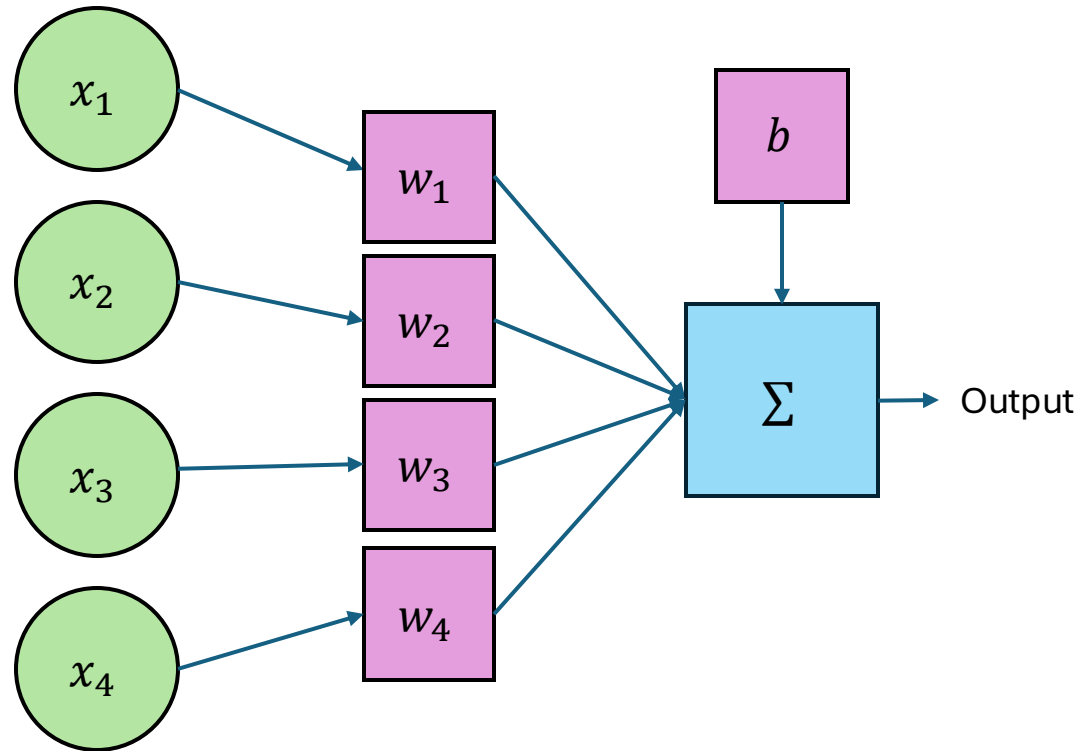
Activation Function
(many more to come)

$x_1$

$x_2$

$x_3$

$x_4$

$w_1$

$w_2$

$w_3$

$w_4$

$b$

$\Sigma$

Output

# Predicting with a Perceptron

1. Take each of the inputs and multiply by corresponding weight
2. Sum the results, add bias term
3. If output is above 0, return 1, otherwise return 0

Activation Function
(many more to come)

$x_1$

$x_2$

$x_3$

$x_4$

$w_1$

$w_2$

$w_3$

$w_4$

$b$

$\Sigma$

Output

Activation Function
Operates on output of neuron

# Understanding Weights

# Understanding Weights

What would it mean for a weight to be 0?

$x_1$

$x_2$

$x_3$

$x_4$

$w_1$

$w_2$

$w_3$

$w_4$
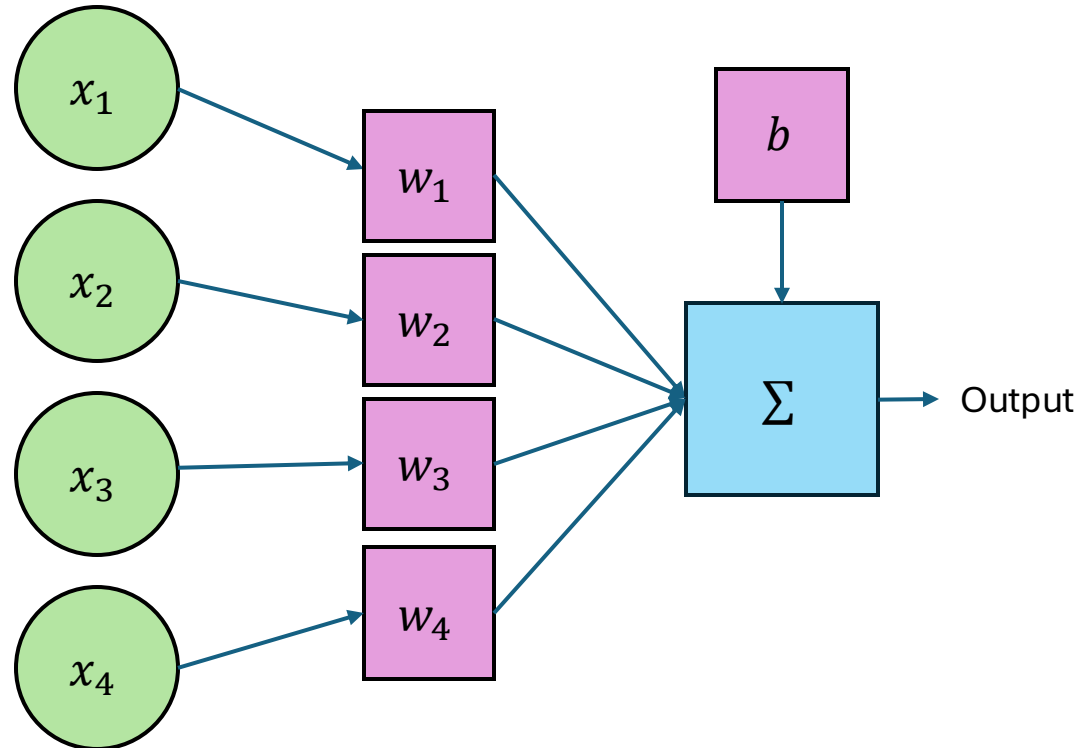
$b$

$\Sigma$

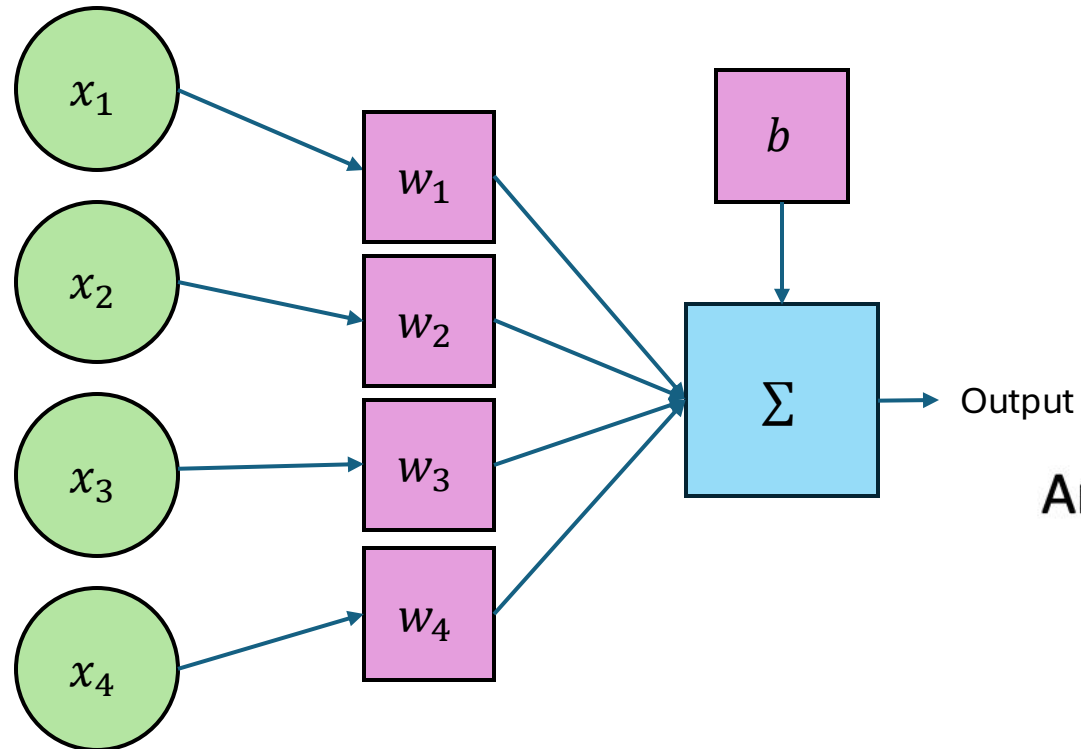Output

# Understanding Weights

# Understanding Weights



What would it mean for a weight to be 0?

What would it mean for a weight to be very positive?

What would it mean for a weight to be very negative?

$x_1$

$x_2$

$x_3$

$x_4$

$w_1$

$w_2$

$w_3$

$w_4$

$b$

$\Sigma$

Output

# Understanding Weights

What would it mean for a weight to be 0?

What would it mean for a weight to be very positive?

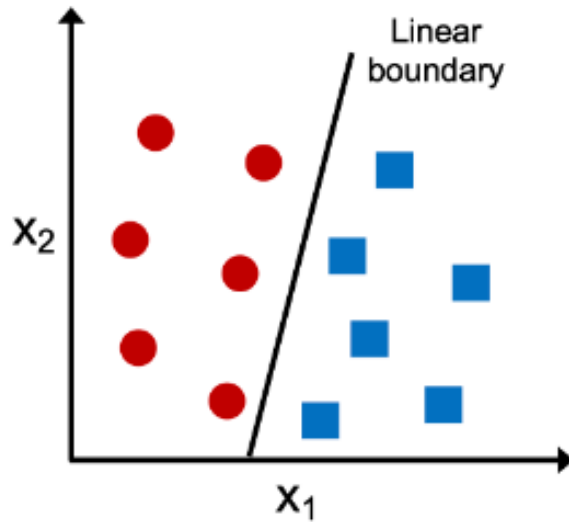What would it mean for a weight to be very negative?

$x_1$

$x_2$

$x_3$

$x_4$

$w_1$

$w_2$

$w_3$

$w_4$

$b$

$\Sigma$

Output

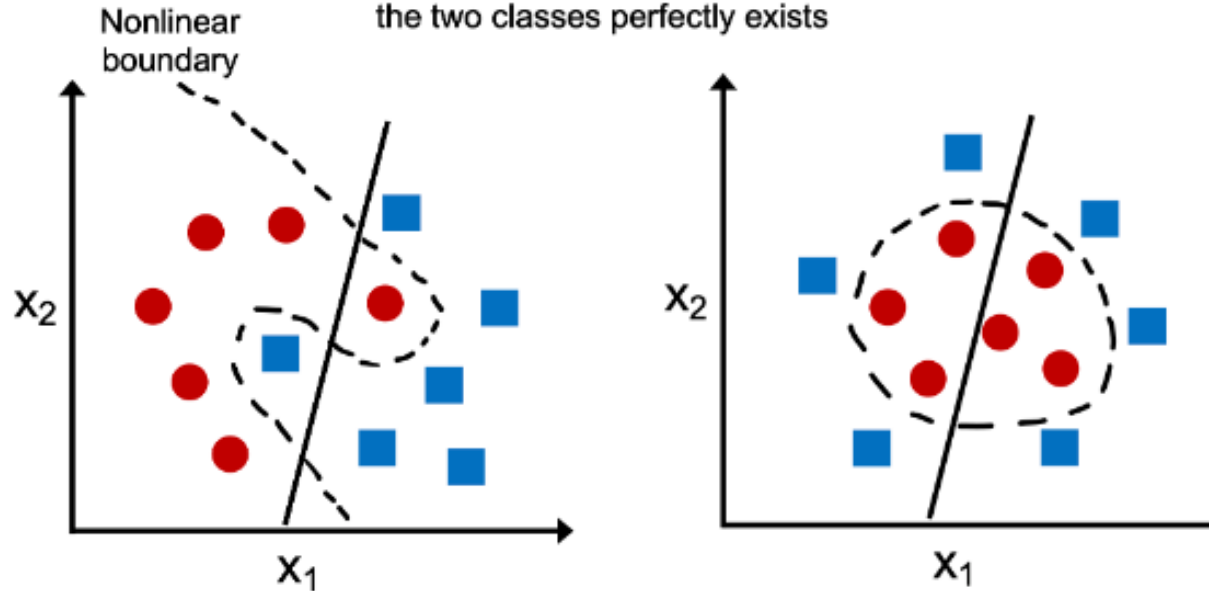Any questions?

???

# How Strong are Linear Separators?



Linearly separable
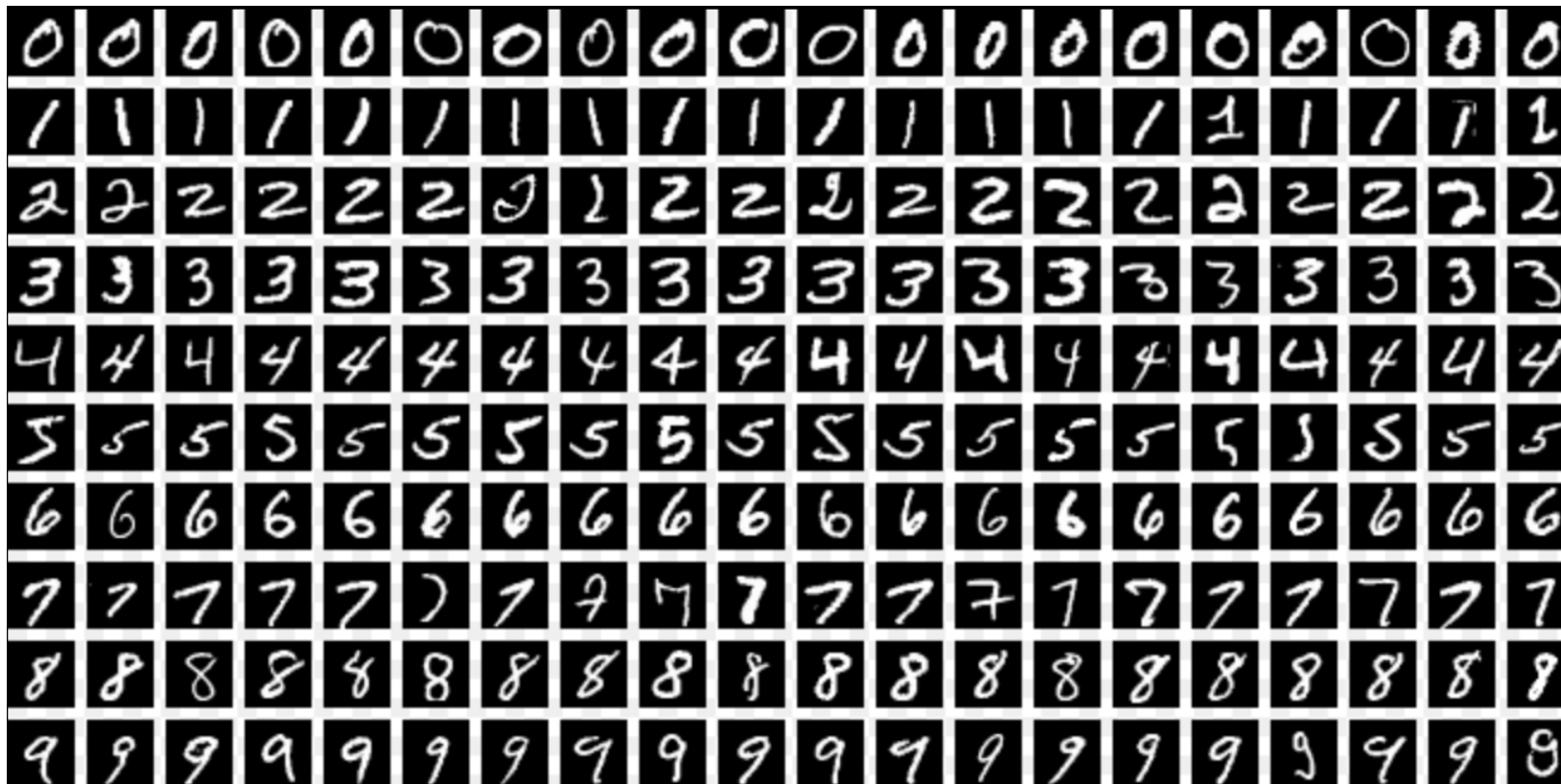A linear decision boundary that separates the two classes exists

Not linearly separable
No linear decision boundary that separates the two classes perfectly exists
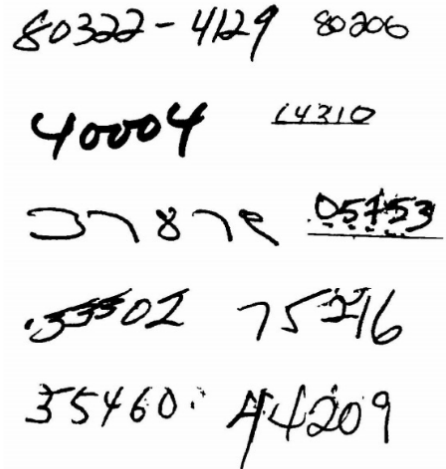
# MNIST

The most famous dataset in Deep Learning

**M**odified **N**ational **I**nstitute of **S**tandards and **T**echnology database



Image courtesy of Wikipedia

# Motivation: Zip Code Recognition

- In 1990s, great increase in documents on paper (mail, checks, books, etc.)

- Motivation for a ZIP code recognizer on real U.S. mail for the postal service!

# Our Problem:

Input: $\mathbb{X}$

Target: $\mathbb{Y}$

Which digit is it?

"3"

Function: f

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

"three"
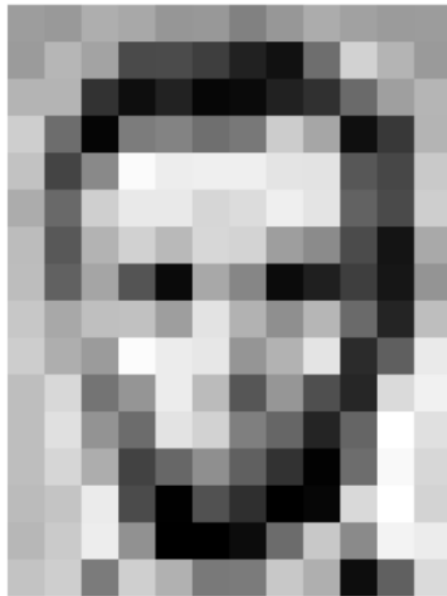
How Does a Computer know this is a three?

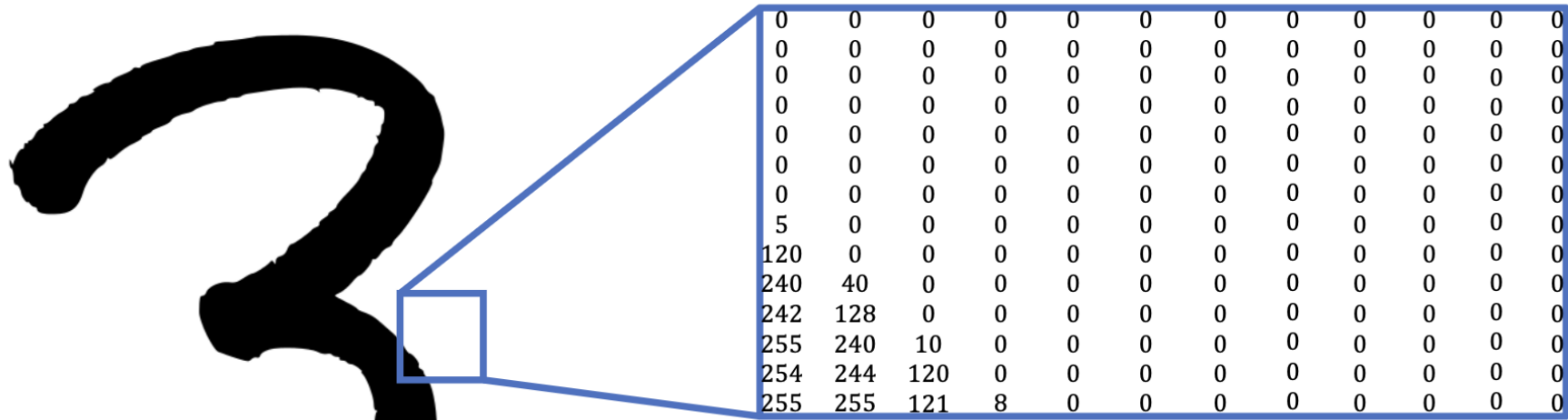# Representing digits in the computer

- Numbers known as *pixel values* (a grid of discrete values that make up an image)

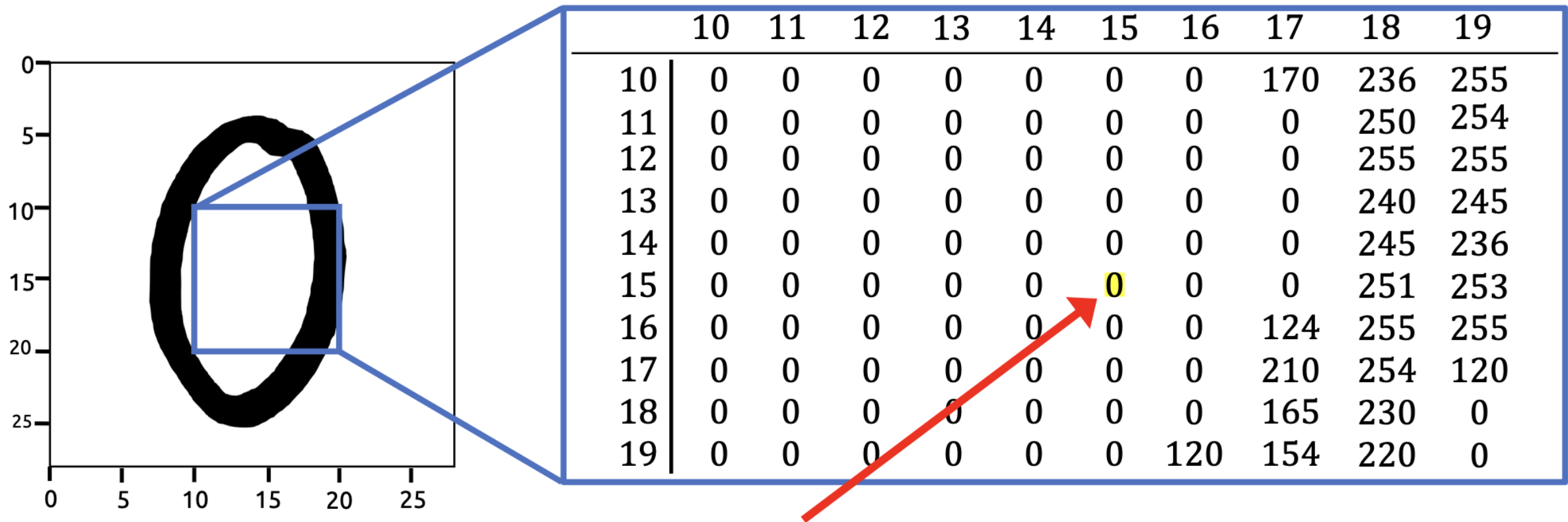0 is white, 255 is black, and numbers in between are shades of gray

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 120 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 240 | 40 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 242 | 128 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 255 | 240 | 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 254 | 244 | 120 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 255 | 255 | 121 | 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

what the
computer sees

|     | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
|-----|----|----|----|----|----|----|----|----|----|----|
| 10  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 170 | 236 | 255 |
| 11  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 250 | 254 |
| 12  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 255 | 255 |
| 13  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 240 | 245 |
| 14  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 245 | 236 |
| 15  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 251 | 253 |
| 16  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 124 | 255 | 255 |
| 17  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 210 | 254 | 120 |
| 18  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 165 | 230 | 0  |
| 19  | 0  | 0  | 0  | 0  | 0  | 0  | 120 | 154 | 220 | 0  |

what the computer sees

- Pixel in position [15, 15] is light.

Center is typically empty for 0's.
How does this compare with 3's?

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 255 | 255 | 255 | 255 | 255 | 253 | 254 | 245 | 255 |
| 255 | 255 | 251 | 255 | 255 | 255 | 254 | 235 | 252 |
| 255 | 252 | 255 | 250 | 255 | 245 | 255 | 253 | 234 |
| 253 | 255 | 255 | 255 | 251 | 254 | 255 | 255 | 235 |
| 255 | 255 | 252 | 255 | 249 | 255 | 239 | 243 | 255 |
| 255 | 250 | 255 | 245 | 255 | 255 | 254 | 244 | 254 |
| 255 | 255 | 255 | 255 | 249 | 255 | 255 | 255 | 244 |
| 249 | 255 | 253 | 255 | 233 | 255 | 249 | 245 | 239 |
| 255 | 255 | 255 | 250 | 255 | 254 | 251 | 243 | 251 |
| 245 | 240 | 244 | 240 | 239 | 244 | 255 | 244 | 248 |
| 242 | 128 | 140 | 150 | 130 | 128 | 110 | 245 | 246 |
| 240 | 240 | 4 | 5 | 4 | 3 | 2 | 118 | 120 |
| 240 | 5 | 4 | 2 | 0 | 0 | 0 | 4 | 2 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

| 255 | 255 | 255 | 255 | 255 | 253 | 254 | 245 | 255 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 255 | 255 | 251 | 255 | 255 | 255 | 254 | 235 | 252 |
| 255 | 252 | 255 | 250 | 255 | 245 | 255 | 253 | 234 |
| 253 | 255 | 255 | 255 | 251 | 254 | 255 | 255 | 235 |
| 255 | 255 | 252 | 255 | 249 | 255 | 239 | 243 | 255 |
| 255 | 250 | 255 | 245 | 255 | 255 | 254 | 244 | 254 |
| 255 | 255 | 255 | 255 | 249 | 255 | 255 | 255 | 244 |
| 249 | 255 | 253 | 255 | 233 | 255 | 249 | 245 | 239 |
| 255 | 255 | 255 | 250 | 255 | 254 | 251 | 243 | 251 |
| 245 | 240 | 244 | 240 | 239 | 244 | 255 | 244 | 248 |
| 242 | 128 | 140 | 150 | 130 | 128 | 110 | 245 | 246 |
| 240 | 240 | 4 | 5 | 4 | 3 | 2 | 118 | 120 |
| 240 | 5 | 4 | 2 | 0 | 0 | 0 | 4 | 2 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Can we define a set of *heuristics* (i.e. rules based on our intuition), to classify digits?

# Machine Learning Pipeline for Digit Recognition

# Train, validation, and test sets

- **Training Set**: Used to adjust parameters of model
- *Validation set* — used to test how well we're doing as we develop
  - Prevents *overfitting*
- *Test Set* — used to evaluate the model once the model is done

# MNIST

- 60,000 Images in training set

- 10,000 Images in test set

- No explicit validation set

# MNIST

- 60,000 Images in training set

- 10,000 Images in test set

- No explicit validation set

What do you suggest we do?

# Machine Learning Pipeline for Digit Recognition

# Our Problem:

Classifying MNIST digits requires predicting 1 of 10 possible values

Input: $\mathbb{X}$

Target: $\mathbb{Y}$

Pixel Grid

Which digit is it?

$x^{(1)} =$ 

2

Function: f

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

28x28 pixels

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

$x^{(2)} =$ 

0

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

# Our Problem:

Classifying MNIST digits requires predicting 1 of 10 possible values

Input: $\mathbb{X}$

What is our input space?

Target: $\mathbb{Y}$

**Pixel Grid**

**2**

$x^{(1)} =$

28x28 pixels

➡ Function: f ➡

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

**Which digit is it?**

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

**0**

$x^{(2)} =$

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

# Our Problem:

Classifying MNIST digits requires predicting 1 of 10 possible values

Input: $\mathbb{X}$

Target: $\mathbb{Y}$

What is our input space?

What is our output space?

**Pixel Grid**

Which digit is it?

$x^{(1)} =$

$2$

→ Function: f →

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

28x28 pixels

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

$x^{(2)} =$

$0$

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

# Our Problem:

Classifying MNIST digits requires predicting 1 of 10 possible values

Input: $\mathbb{X}$

Target: $\mathbb{Y}$

**Pixel Grid**

Which digit is it?

| What is our input space? |
| What is our output space? |
| What is our prediction task? |

$x^{(1)} =$ 2

**28x28 pixels**

Function: f

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

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

$x^{(2)} =$ 0

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

# Our simplified problem:

Input: $\mathbb{X}$

What is our input space?

What is our output space?

What is our prediction task?

Target: $\mathbb{Y}$

Pixel Grid

Is it digit 2?

$x^{(1)} =$

2

Function: f

$y^{(1)} = 1$  ✅

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

28x28 pixels
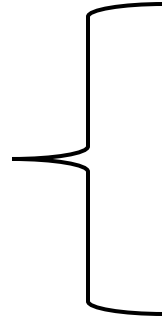
$x^{(2)} =$

0

$y^{(2)} = 0$  ❌

# A bit of a cliffhanger...

- How well do you think a perceptron will do on this task?

- Perceptrons are linear classifiers... what does it mean for images to be linearly separable?

- Perceptrons have a discontinuous activation function, which is not differentiable. How are we going to find good parameters without a nice closed-form solution?

# Recap

Linear Regression

Matrix and Vector Notation

Closed Form solution for finding optimal parameters

# Recap

Linear Regression

Matrix and Vector Notation

Closed Form solution for finding optimal parameters

Perceptrons

natural extension of linear models to binary classification tasks

Biological inspiration of activation threshold

Only differ from Linear Regression in terms of activation function

# Recap

**Linear Regression**
- Matrix and Vector Notation
- Closed Form solution for finding optimal parameters

**Perceptrons**
- natural extension of linear models to binary classification tasks
- Biological inspiration of activation threshold
- Only differ from Linear Regression in terms of activation function

**MNIST**
- Handwritten Digit Representation
- Handwriting Classification Task Framing