

Artificial Neurons and Gradient Descent

Spring 2018

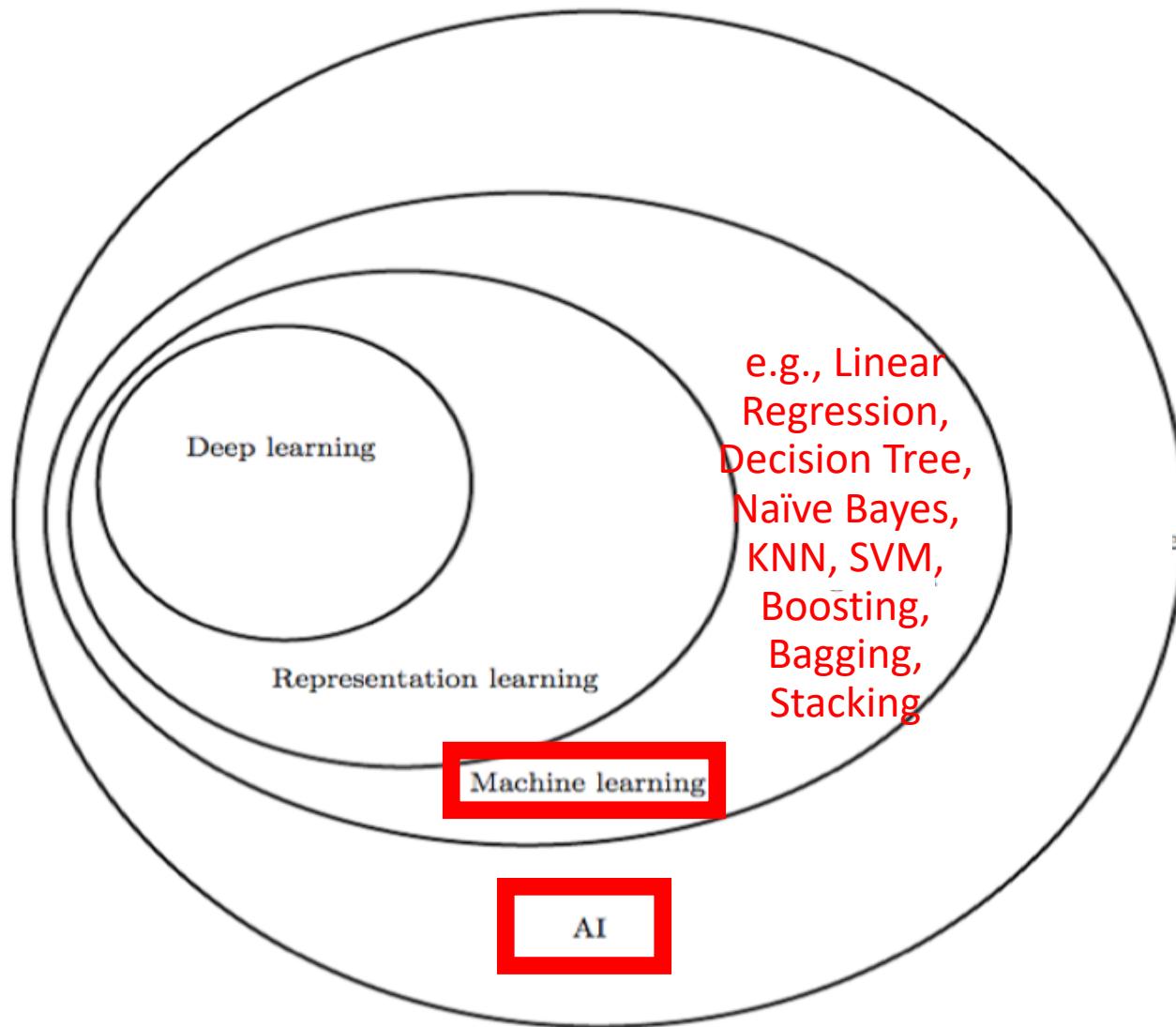
Today's Topics

- History of Modern Neural Networks and “Deep Learning”
- Single Layer Neural Network: Perceptron
- Single Layer Neural Network: Adaline
- Variants of Gradient Descent
- Lab

Today's Topics

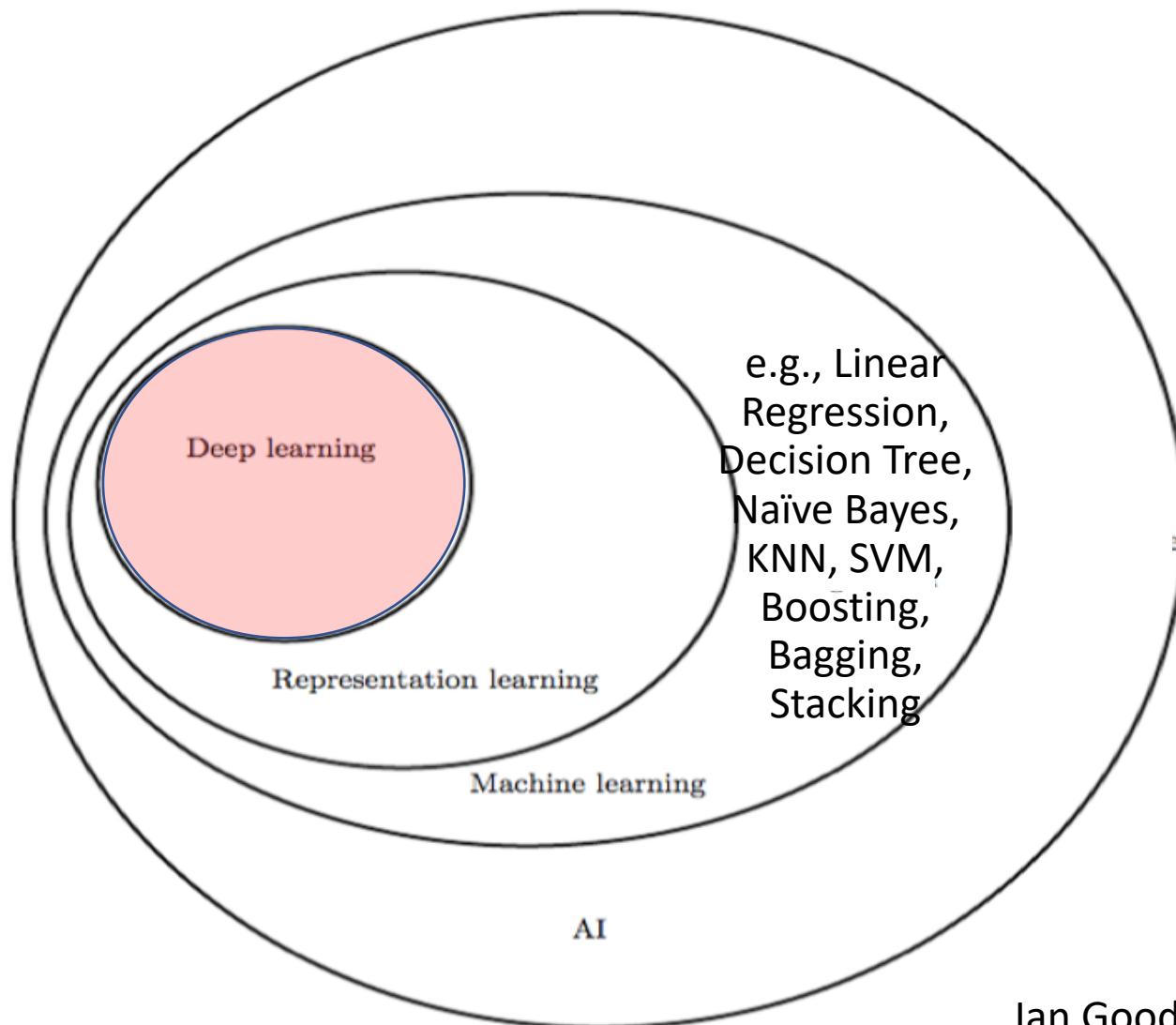
- History of Modern Neural Networks and “Deep Learning”
- Single Layer Neural Network: Perceptron
- Single Layer Neural Network: Adaline
- Variants of Gradient Descent
- Lab

Review: AI & Relationship to Different Disciplines



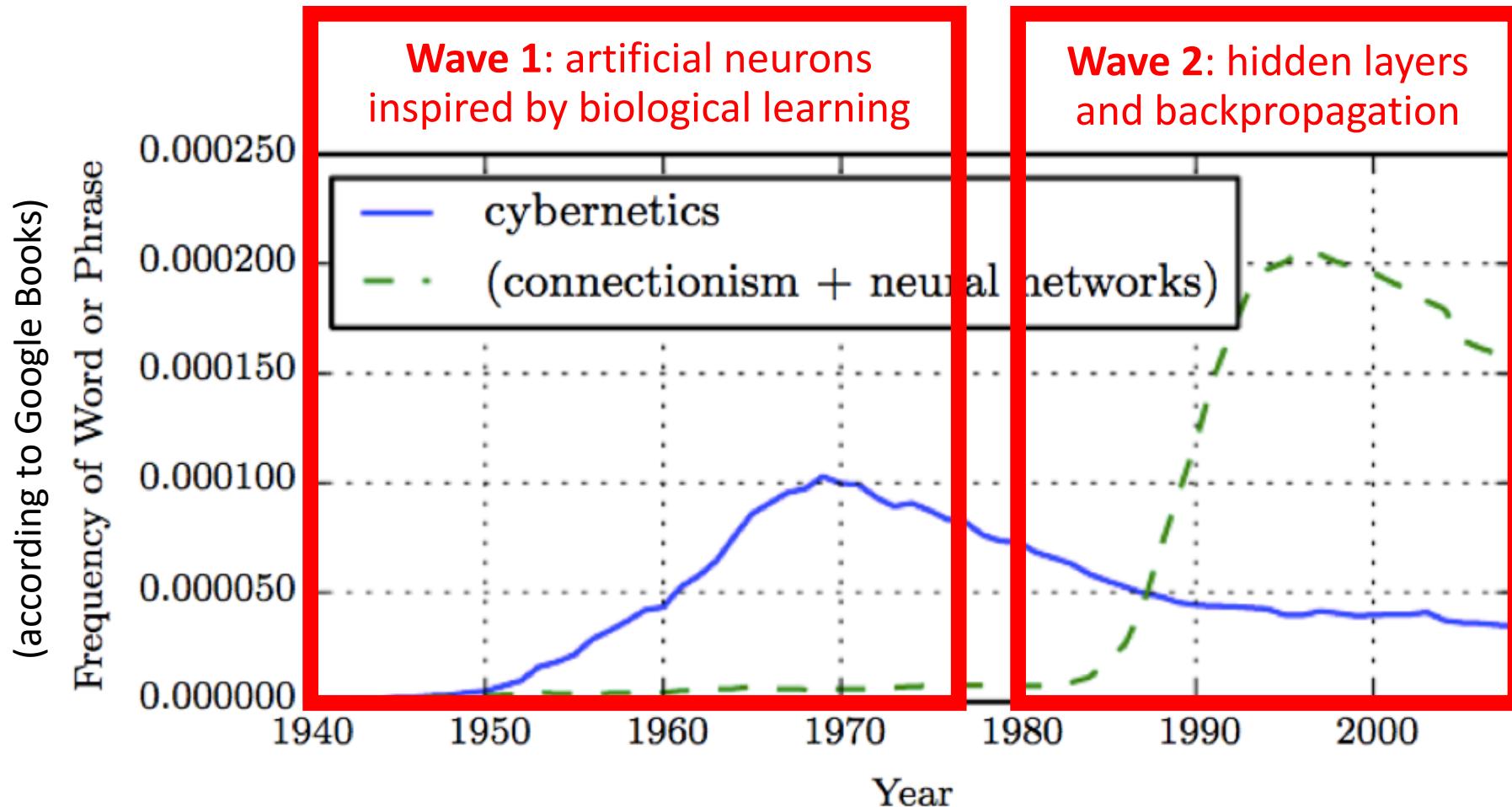
- What is artificial intelligence?
 - 1956: machines that do intelligent things
- What is machine learning?
 - 1956: algorithms that learn on their own
- What machine learning algorithms have we studied so far?

Focus in Next Few Weeks: Deep Learning



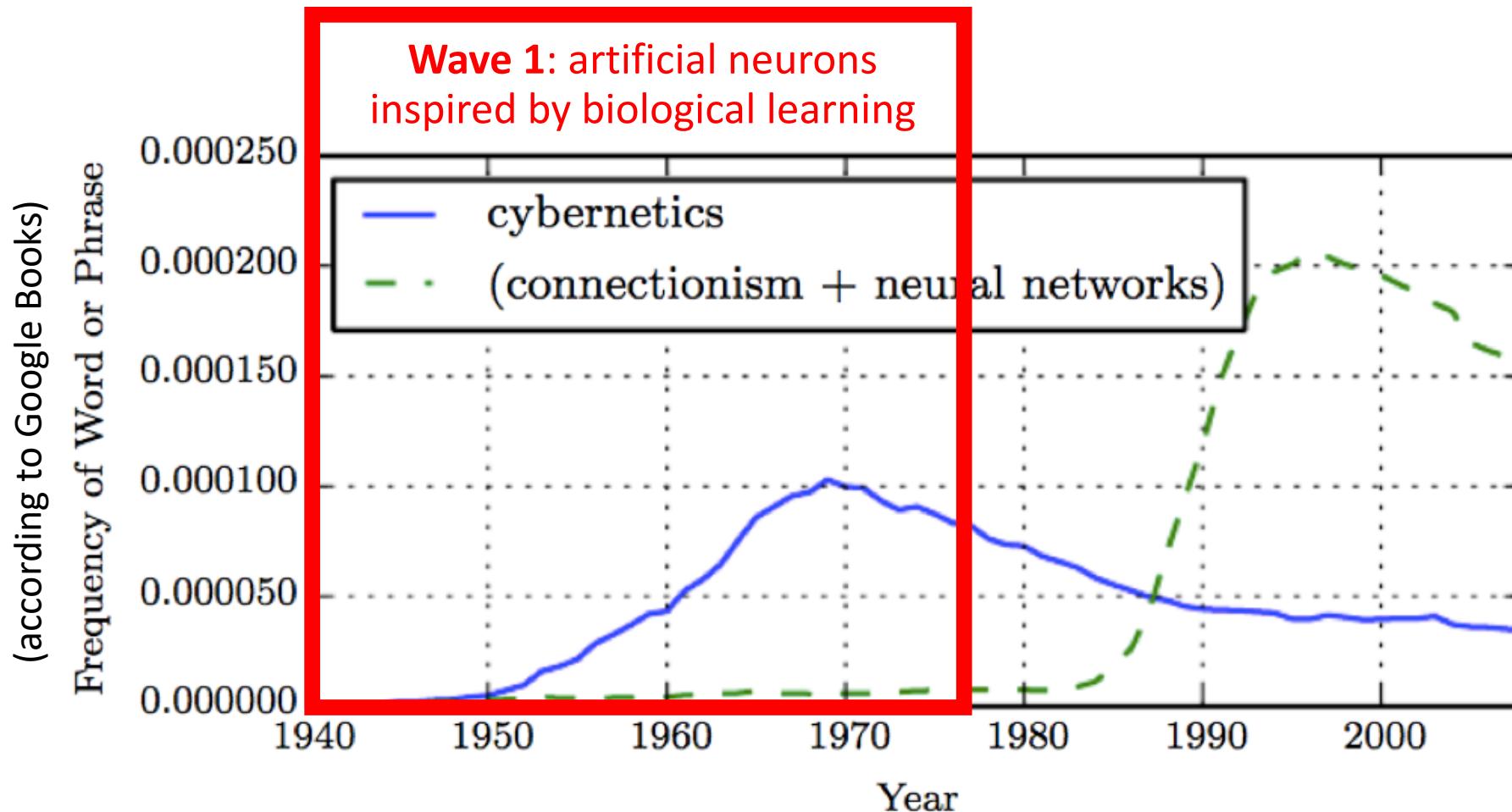
- What is artificial intelligence?
 - 1956: machines that do intelligent things
- What is machine learning?
 - 1956: algorithms that learn on their own

The History of “Deep Learning”

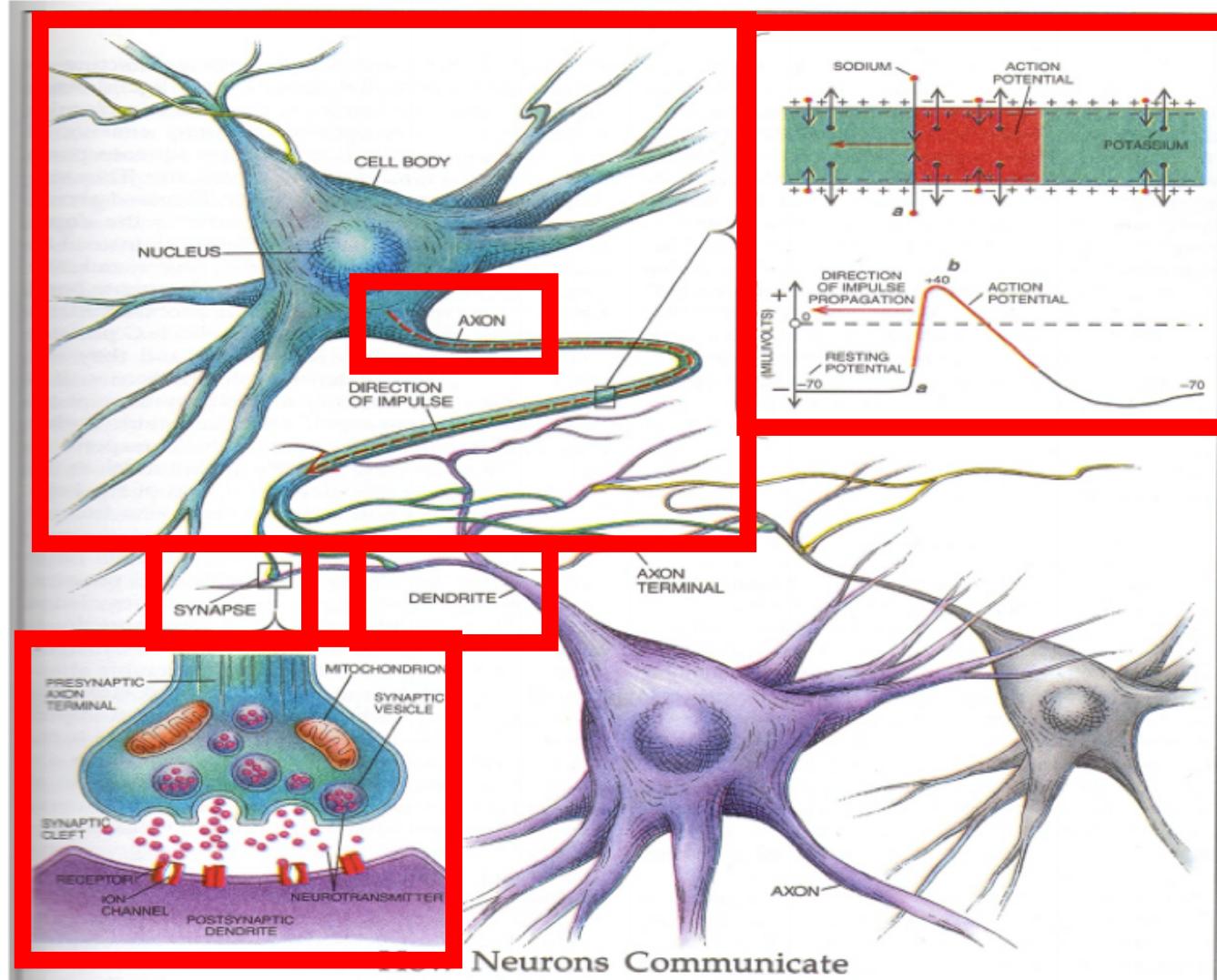


Wave 3: field emerged again as “deep learning” in 2006

Today's Topic: Artificial Neurons

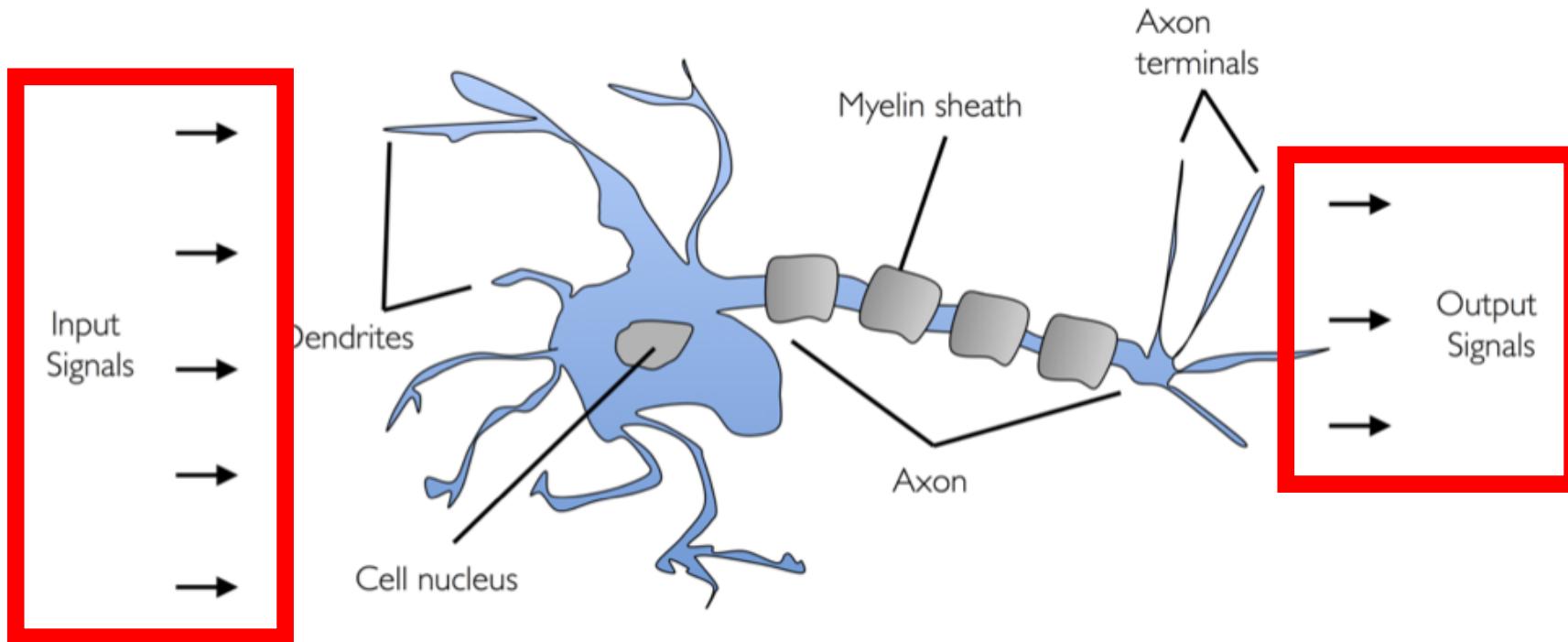


Biological Neuron: Basic Logic Unit in Brain



- What is a neuron?
 - Primary component of central nervous system for receiving, processing, and transmitting information
- How does a neuron communicate?
 - Synapse connects neurons
 - It receives information from its dendrites
 - It transmits information through its axon
 - It transmits information from its axon to another neuron's dendrites
 - It “Fires” (transmits information) ONLY when triggered above a certain voltage
- How many neurons does a human have?
 - Approximately 10^{11}

Artificial Neuron: Model of Biological Neuron



- Neuron “firing” (outputs signal) is an “all-or-none” process.
- A fixed number of input signals must be received within a short period of time for a neuron to “fire”

Artificial Neuron: McCulloch-Pitts Neuron

- 1943: First mathematical model of neuron proposed by...



Warren McCulloch: Neurophysiologist



Walter Pitts: Mathematician

- Recall: first computing machines created in early 1940s

https://en.wikipedia.org/wiki/Walter_Pitts

http://web.csulb.edu/~cwallis/artificialn/warren_mcculloch.html

Warren McCulloch and Walter Pitts, ^mA Logical Calculus of Ideas Immanent in Nervous Activity, 1943.

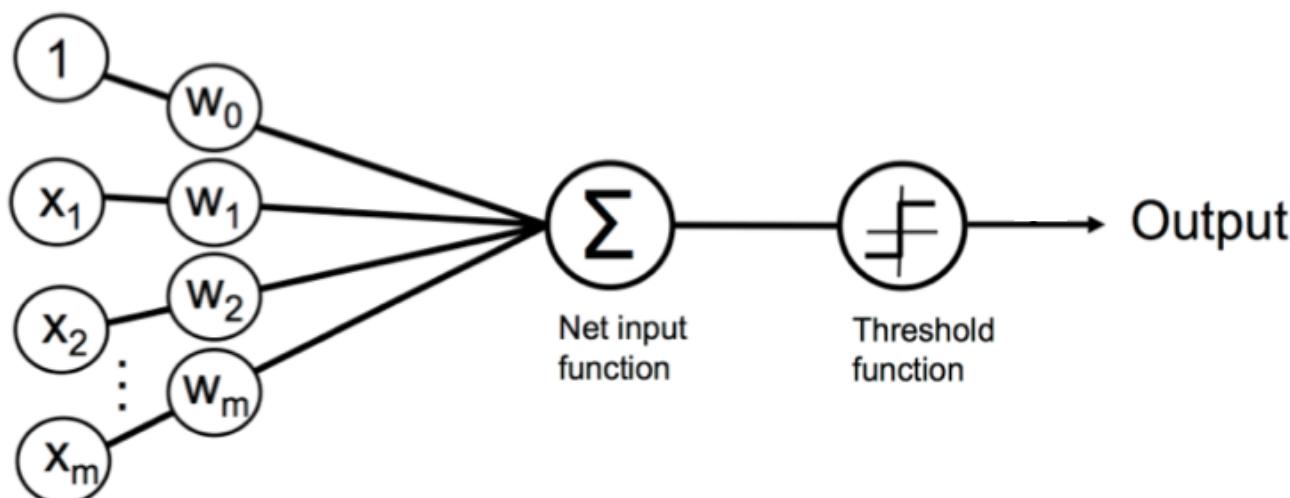
Today's Topics

- History of Modern Neural Networks and “Deep Learning”
- Single Layer Neural Network: Perceptron
- Single Layer Neural Network: Adaline
- Variants of Gradient Descent
- Lab

Artificial Neuron: Perceptron

- 1958: Linear Threshold Unit (LTU)

- Input features
- Model coefficients
- Outputs weighted sum of inputs

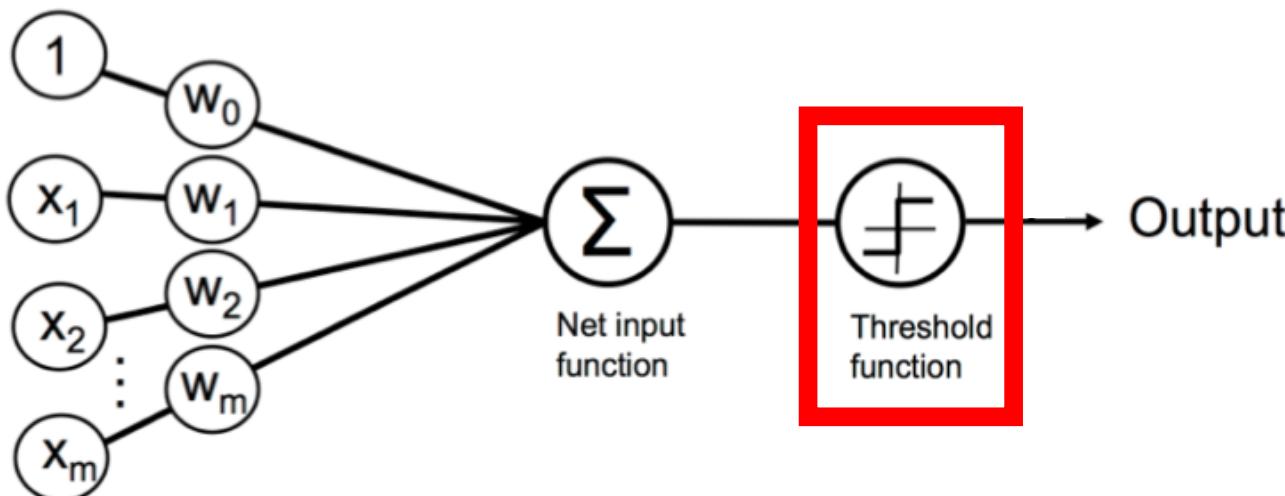


Python Machine Learning; Raschka & Mirjalili

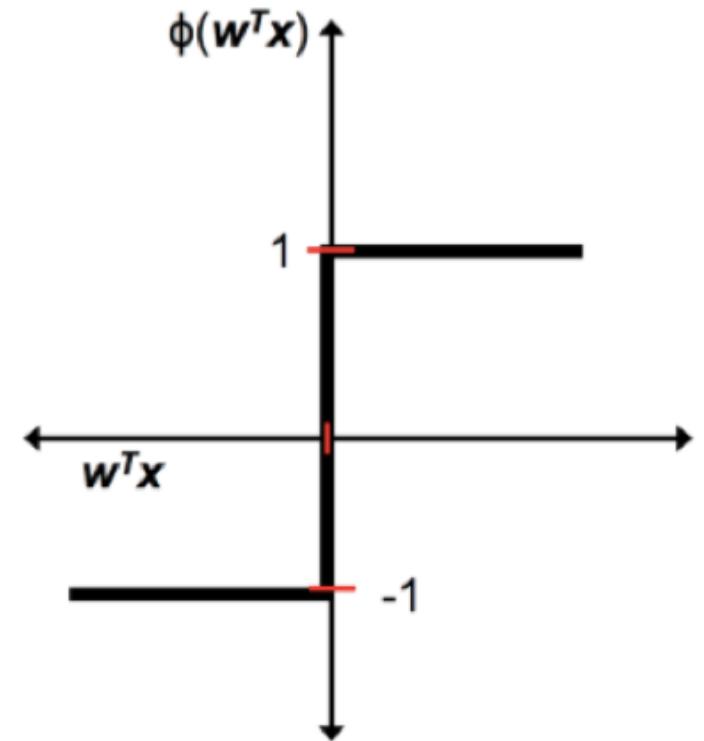
Frank Rosenblatt, The **perceptron**: a probabilistic model for information storage and organization in the brain, 1958.

Artificial Neuron: Perceptron

- Formal definition of artificial neuron: $\phi(\mathbf{w}^T \mathbf{x})$
- Decision function is a unit step function:
 - “fire” or “do not fire”
 - mimics human brain



$$\phi(\mathbf{w}^T \mathbf{x}) = \begin{cases} 1 & \text{if } \phi(\mathbf{w}^T \mathbf{x}) \geq \theta \\ -1 & \text{otherwise} \end{cases}$$

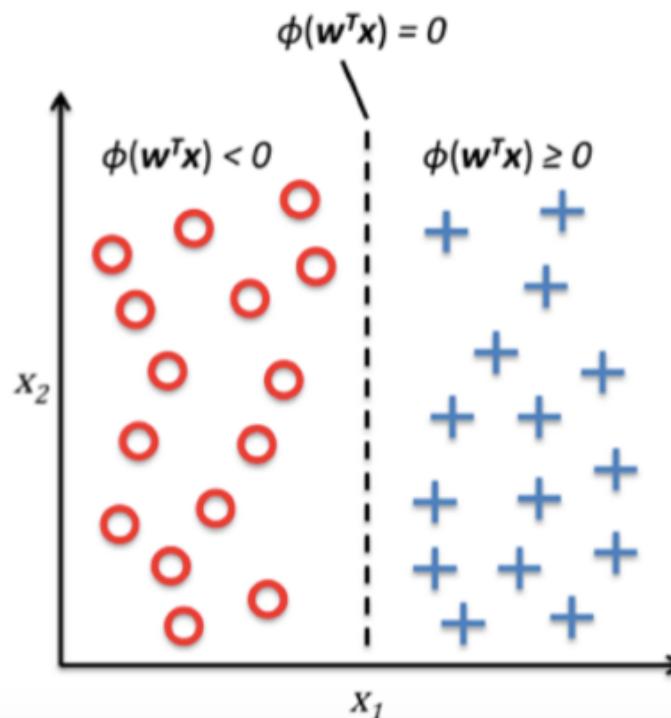


Artificial Neuron: Perceptron

- Formal definition of artificial neuron: $\phi(\mathbf{w}^T \mathbf{x}) = \begin{cases} 1 & \text{if } \phi(\mathbf{w}^T \mathbf{x}) \geq 0 \\ -1 & \text{otherwise} \end{cases}$
- Rewrite to: $w_0x_0 + w_1x_1 + \dots + w_mx_m = \sum_{j=0}^m \mathbf{x}_j \mathbf{w}_j = \mathbf{w}^T \mathbf{x}$

$$w_0x_0 + w_1x_1 + \dots + w_mx_m$$

A diagram showing a vector \mathbf{x} with components x_0, x_1, \dots, x_m . The component x_0 is labeled w_0 and has a red box around it. Arrows point from the labels w_0 and x_0 to the corresponding terms in the equation above. The other components x_1, \dots, x_m are labeled w_1, \dots, w_m respectively.

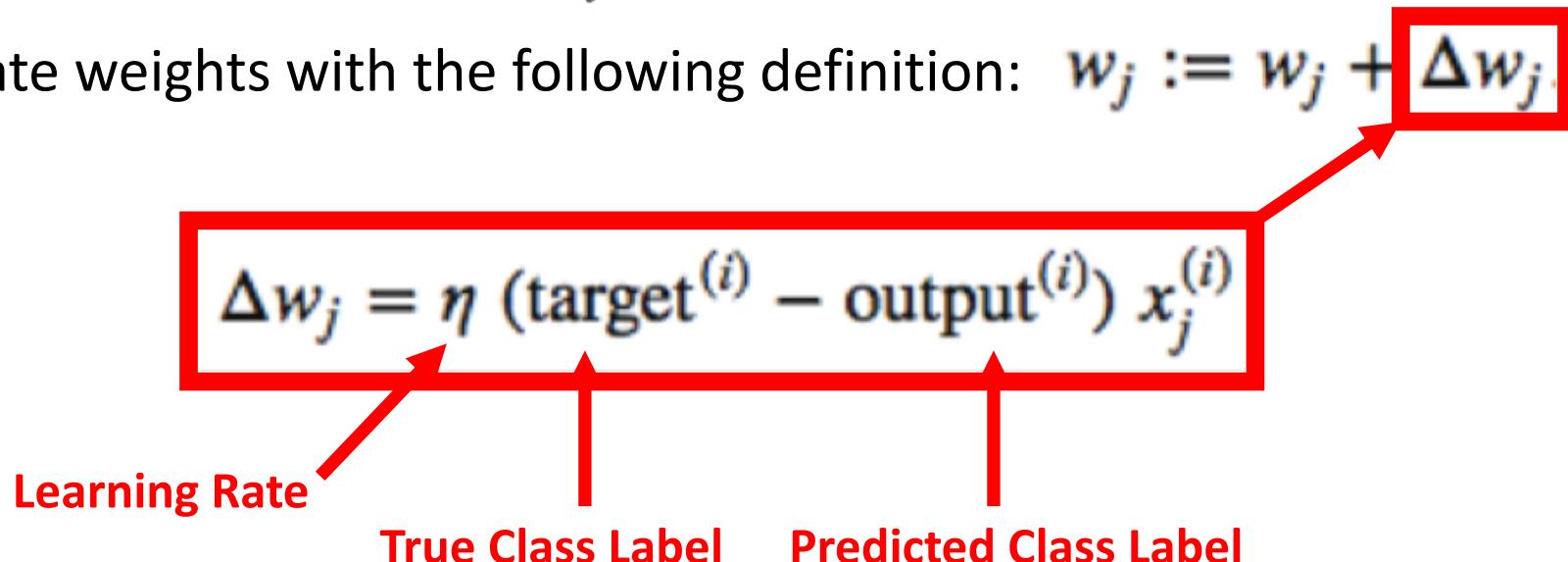


Perceptron Learning Algorithm

1. Initialize weights to 0 or small random numbers
2. For each training sample:

1. Compute output value: $\sum_{j=0}^m \mathbf{x}_j \mathbf{w}_j = \mathbf{w}^T \mathbf{x}$

2. Update weights with the following definition: $\mathbf{w}_j := \mathbf{w}_j + \Delta \mathbf{w}_j$



Perceptron Learning Algorithm: e.g., 2D dataset

1. Initialize weights to 0 or small random numbers
2. For each training sample:

1. Compute output value: $\sum_{j=0}^m \mathbf{x}_j \mathbf{w}_j = \mathbf{w}^T \mathbf{x}$

2. Update weights with the following definition: $\mathbf{w}_j := \mathbf{w}_j + \Delta \mathbf{w}_j$

$$\Delta \mathbf{w}_0 = \eta (\text{target}^{(i)} - \text{output}^{(i)})$$

$$\Delta \mathbf{w}_1 = \eta (\text{target}^{(i)} - \text{output}^{(i)}) \mathbf{x}_1^{(i)}$$

$$\Delta \mathbf{w}_2 = \eta (\text{target}^{(i)} - \text{output}^{(i)}) \mathbf{x}_2^{(i)}$$

All weights updated simultaneously

Learning: What Happens to Weights When Perceptron Predicts Correct Class Label?

1. Initialize weights to 0 or small random numbers
2. For each training sample:
 1. Compute output value: $\sum_{j=0}^m \mathbf{x}_j \mathbf{w}_j = \mathbf{w}^T \mathbf{x}$
 2. Update weights with the following definition: $\mathbf{w}_j := \mathbf{w}_j + \Delta \mathbf{w}_j$

$$1. \text{ Compute output value: } \sum_{j=0}^m \mathbf{x}_j \mathbf{w}_j = \mathbf{w}^T \mathbf{x}$$

$$2. \text{ Update weights with the following definition: } \mathbf{w}_j := \mathbf{w}_j + \Delta \mathbf{w}_j$$

equals 0, so no weight update

$$\Delta \mathbf{w}_j = \eta (\text{target}^{(i)} - \text{output}^{(i)}) \mathbf{x}_j^{(i)}$$

Learning Rate True Class Label Predicted Class Label

Learning: What Happens to Weights When Perceptron Predicts Wrong Class Label?

1. Initialize weights to 0 or small random numbers
2. For each training sample:

1. Compute output value: $\sum_{j=0}^m \mathbf{x}_j \mathbf{w}_j = \mathbf{w}^T \mathbf{x}$

2. Update weights with the following definition: $\mathbf{w}_j := \mathbf{w}_j + \Delta \mathbf{w}_j$

equals 2 or -2 so moves weights closer to positive or negative target class

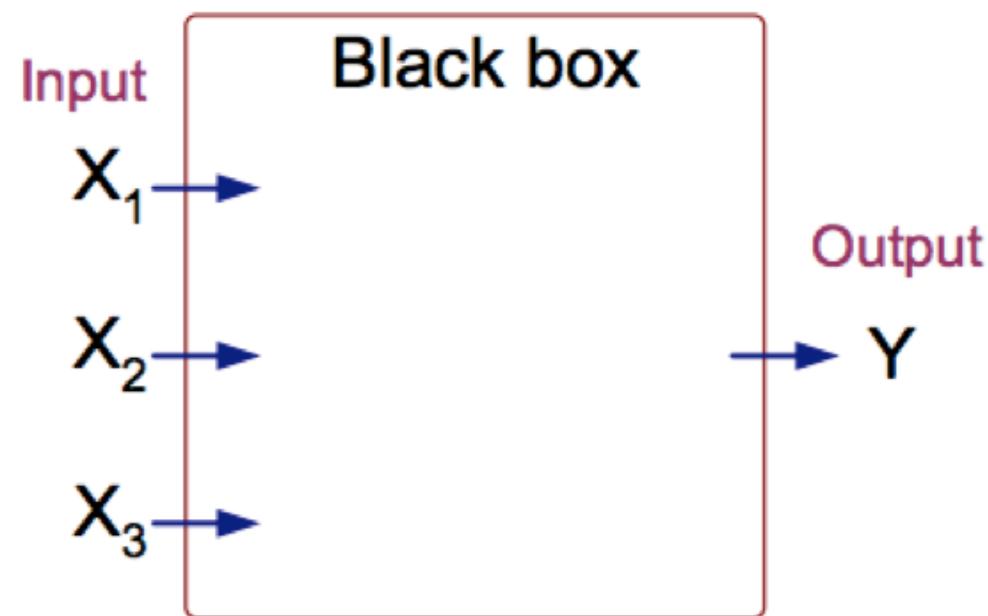
$$\Delta \mathbf{w}_j = \eta (\text{target}^{(i)} - \text{output}^{(i)}) \mathbf{x}_j^{(i)}$$

↑ ↑ ↑
Learning Rate True Class Label Predicted Class Label

Learning Example: Example Data

- True Model: Y is 1 if at least two of the three inputs are equal to 1.

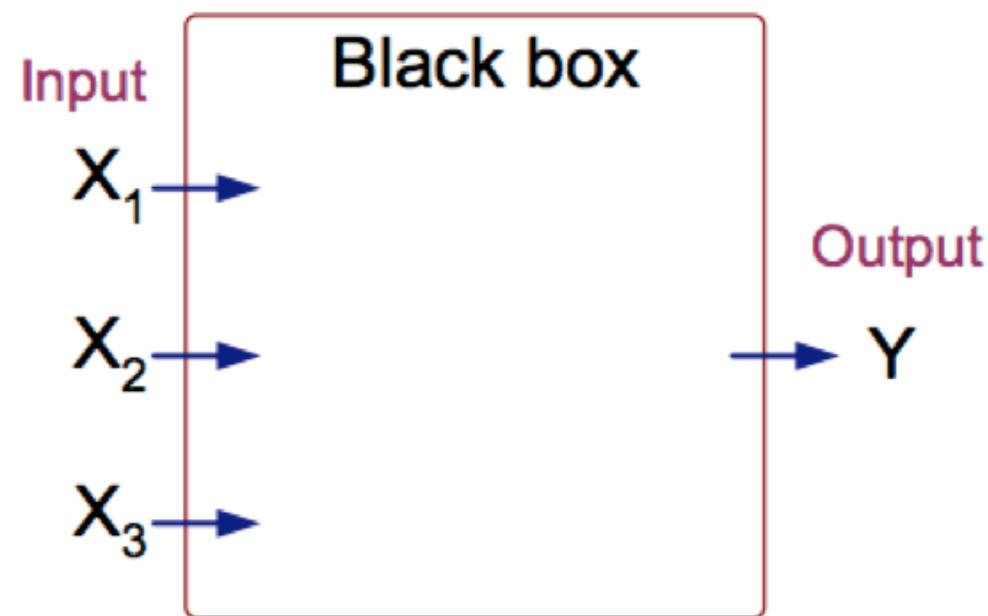
X_1	X_2	X_3	Y
1	0	0	
1	0	1	
1	1	0	
1	1	1	
0	0	1	
0	1	0	
0	1	1	
0	0	0	



Learning Example: Example Data

- True Model: Y is 1 if at least two of the three inputs are equal to 1.

X_1	X_2	X_3	Y
1	0	0	-1
1	0	1	1
1	1	0	1
1	1	1	1
0	0	1	-1
0	1	0	-1
0	1	1	1
0	0	0	-1



Learning Example: Learn with First Sample

- Compute output value: $\sum_{j=0}^m \mathbf{x}_j \mathbf{w}_j = \mathbf{w}^T \mathbf{x}$;

$$\phi(\mathbf{w}^T \mathbf{x}) = \begin{cases} 1 & \text{if } \phi(\mathbf{w}^T \mathbf{x}) \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

X ₁	X ₂	X ₃	Y	Predicted	
1	0	0	-1	?	

w ₀	w ₁	w ₂	w ₃
0	0	0	0

Learning Example: Learn with First Sample

- Compute output value: $\sum_{j=0}^m \mathbf{x}_j \mathbf{w}_j = \mathbf{w}^T \mathbf{x}$;

$$\phi(\mathbf{w}^T \mathbf{x}) = \begin{cases} 1 & \text{if } \phi(\mathbf{w}^T \mathbf{x}) \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

X ₁	X ₂	X ₃	Y	Predicted	w ₀	w ₁	w ₂	w ₃
1	0	0	-1	1	0	0	0	0

Learning Example: Learn with First Sample

- Update weights: $w_j = w_j + \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_j^{(i)}$; learning rate = 0.1

X ₁	X ₂	X ₃	Y	Predicted
1	0	0	-1	1

w ₀	w ₁	w ₂	w ₃
0	0	0	0
?	?	?	?

$$\Delta w_0 = \eta (\text{target}^{(i)} - \text{output}^{(i)})$$

$$\Delta w_1 = \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_1^{(i)}$$

$$\Delta w_2 = \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_2^{(i)}$$

$$\Delta w_3 = \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_3^{(i)}$$

Learning Example: Learn with First Sample

- Update weights: $w_j = w_j + \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_j^{(i)}$; learning rate = 0.1

X ₁	X ₂	X ₃	Y
Predicted			
1	0	0	-1

w ₀	w ₁	w ₂	w ₃
0	0	0	0
?	?	?	?

$$\Delta w_0 = 0.1(-1-1)*1 = -0.2$$

$$\Delta w_1 = 0.1(-1-1)*1 = -0.2$$

$$\Delta w_2 = 0.1(-1-1)*0 = 0$$

$$\Delta w_3 = 0.1(-1-1)*0 = 0$$

Learning Example: Learn with First Sample

- Update weights: $w_j = w_j + \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_j^{(i)}$; learning rate = 0.1

X ₁	X ₂	X ₃	Y	Predicted
1	0	0	-1	1

w ₀	w ₁	w ₂	w ₃
0	0	0	0
-0.2	-0.2	0	0

$$\Delta w_0 = 0.1(-1-1)*1 = -0.2$$

$$\Delta w_1 = 0.1(-1-1)*1 = -0.2$$

$$\Delta w_2 = 0.1(-1-1)*0 = 0$$

$$\Delta w_3 = 0.1(-1-1)*0 = 0$$

Learning Example: Learn with Second Sample

- Compute output value: $\sum_{j=0}^m \mathbf{x}_j \mathbf{w}_j = \mathbf{w}^T \mathbf{x}$;

$$\phi(\mathbf{w}^T \mathbf{x}) = \begin{cases} 1 & \text{if } \phi(\mathbf{w}^T \mathbf{x}) \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

X ₁	X ₂	X ₃	Y	Predicted	w ₀	w ₁	w ₂	w ₃
1	0	0	-1	1	0	0	0	0
1	0	1	1	?	-0.2	-0.2	0	0

Learning Example: Learn with Second Sample

- Compute output value: $\sum_{j=0}^m \mathbf{x}_j \mathbf{w}_j = \mathbf{w}^T \mathbf{x}$;

$$\phi(\mathbf{w}^T \mathbf{x}) = \begin{cases} 1 & \text{if } \phi(\mathbf{w}^T \mathbf{x}) \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

X ₁	X ₂	X ₃	Y	Predicted	w ₀	w ₁	w ₂	w ₃
1	0	0	-1	1	0	0	0	0
1	0	1	1	-1	-0.2	-0.2	0	0

Learning Example: Learn with Second Sample

- Update weights: $w_j = w_j + \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_j^{(i)}$; learning rate = 0.1

X ₁	X ₂	X ₃	Y	Predicted	w ₀	w ₁	w ₂	w ₃
1	0	0	-1		0	0	0	0
1	0	1	1	1	-0.2	-0.2	0	0

$$\Delta w_0 = \eta (\text{target}^{(i)} - \text{output}^{(i)})$$

$$\Delta w_1 = \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_1^{(i)}$$

$$\Delta w_2 = \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_2^{(i)}$$

$$\Delta w_3 = \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_3^{(i)}$$

Learning Example: Learn with Second Sample

- Update weights: $w_j = w_j + \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_j^{(i)}$; learning rate = 0.1

X ₁	X ₂	X ₃	Y	Predicted	w ₀	w ₁	w ₂	w ₃
1	0	0	-1	1	0	0	0	0
1	0	1	1	-1	-0.2	-0.2	0	0

$$\Delta w_0 = 0.1(1--1)*1 = 0.2$$

$$\Delta w_1 = 0.1(1--1)*1 = 0.2$$

$$\Delta w_2 = 0.1(1--1)*0 = 0$$

$$\Delta w_3 = 0.1(1--1)*1 = 0.2$$

Learning Example: Learn with Second Sample

- Update weights: $w_j = w_j + \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_j^{(i)}$; learning rate = 0.1

X ₁	X ₂	X ₃	Y	Predicted	w ₀	w ₁	w ₂	w ₃
1	0	0	-1	1	0	0	0	0
1	0	1	1	-1	-0.2	-0.2	0	0

$$\Delta w_0 = 0.1(1--1)*1 = 0.2$$

$$\Delta w_1 = 0.1(1--1)*1 = 0.2$$

$$\Delta w_2 = 0.1(1--1)*0 = 0$$

$$\Delta w_3 = 0.1(1--1)*1 = 0.2$$

Learning Example: One Epoch (All Examples)

- $w_j = w_j + \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_j^{(i)}$; learning rate = 0.1

X ₁	X ₂	X ₃	Y
1	0	0	-1
1	0	1	1
1	1	0	1
1	1	1	1
0	0	1	-1
0	1	0	-1
0	1	1	1
0	0	0	-1

	w ₀	w ₁	w ₂	w ₃
0	0	0	0	0
1	-0.2	-0.2	0	0
2	0	0	0	0.2
3	0	0	0	0.2
4	0	0	0	0.2
5	-0.2	0	0	0
6	-0.2	0	0	0
7	0	0	0.2	0.2
8	-0.2	0	0.2	0.2

Learning Example: 6 Epochs (All Examples)

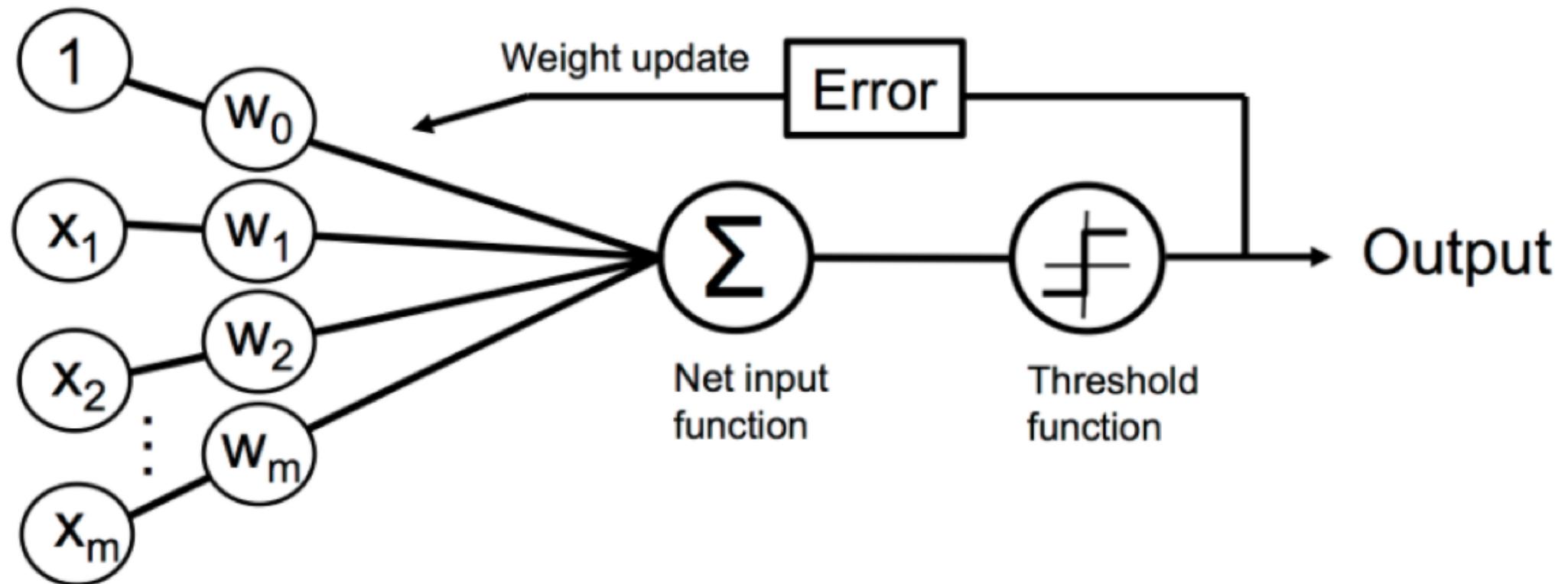
- $w_j = w_j + \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_j^{(i)}$; learning rate = 0.1

X ₁	X ₂	X ₃	Y
1	0	0	-1
1	0	1	1
1	1	0	1
1	1	1	1
0	0	1	-1
0	1	0	-1
0	1	1	1
0	0	0	-1

	w ₀	w ₁	w ₂	w ₃
0	0	0	0	0
1	-0.2	-0.2	0	0
2	0	0	0	0.2
3	0	0	0	0.2
4	0	0	0	0.2
5	-0.2	0	0	0
6	-0.2	0	0	0
7	0	0	0.2	0.2
8	-0.2	0	0.2	0.2

Epoch	w ₀	w ₁	w ₂	w ₃
0	0	0	0	0
1	-0.2	0	0.2	0.2
2	-0.2	0	0.4	0.2
3	-0.4	0	0.4	0.2
4	-0.4	0.2	0.4	0.4
5	-0.6	0.2	0.4	0.2
6	-0.6	0.4	0.4	0.2

Perceptron Learning Algorithm Pictorially

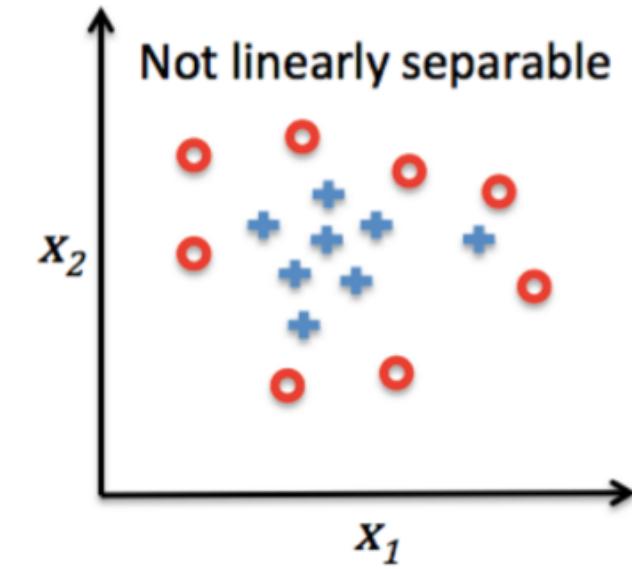
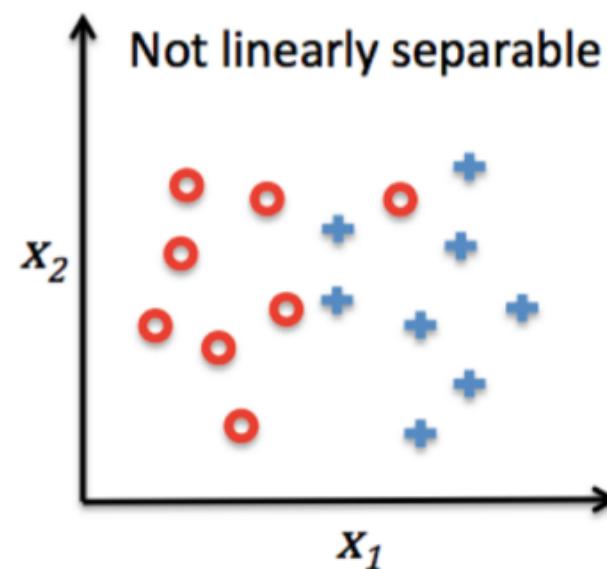
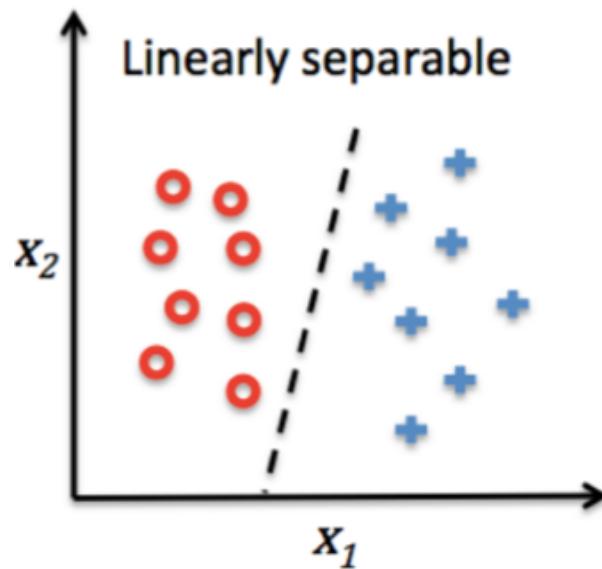


Perceptron Learning Algorithm: What are the Hyperparameters?

- Learning rate
- Number of epochs (passes over the dataset)

Limitations

1. Assumes Data is Linearly Separable

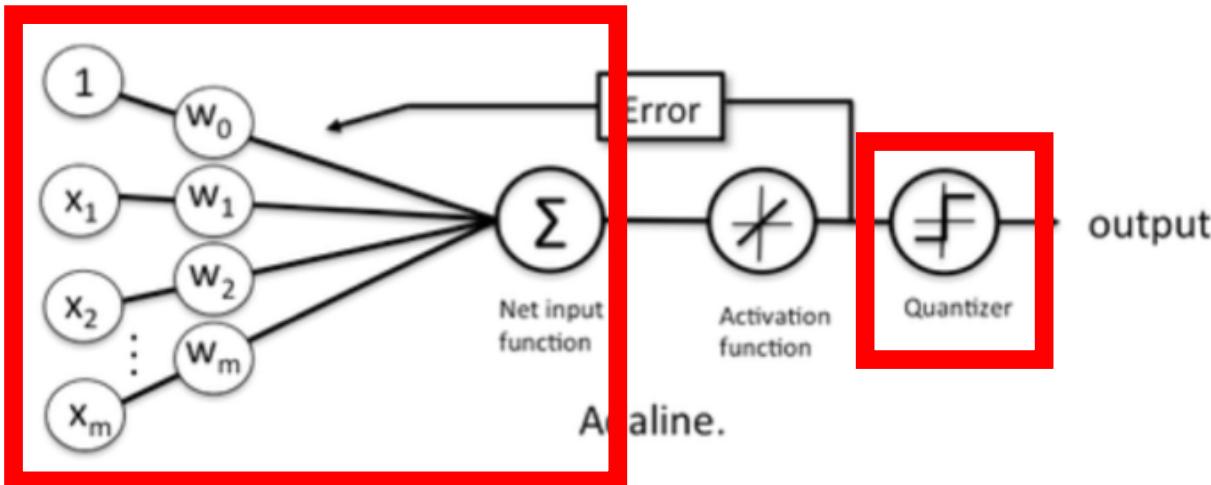
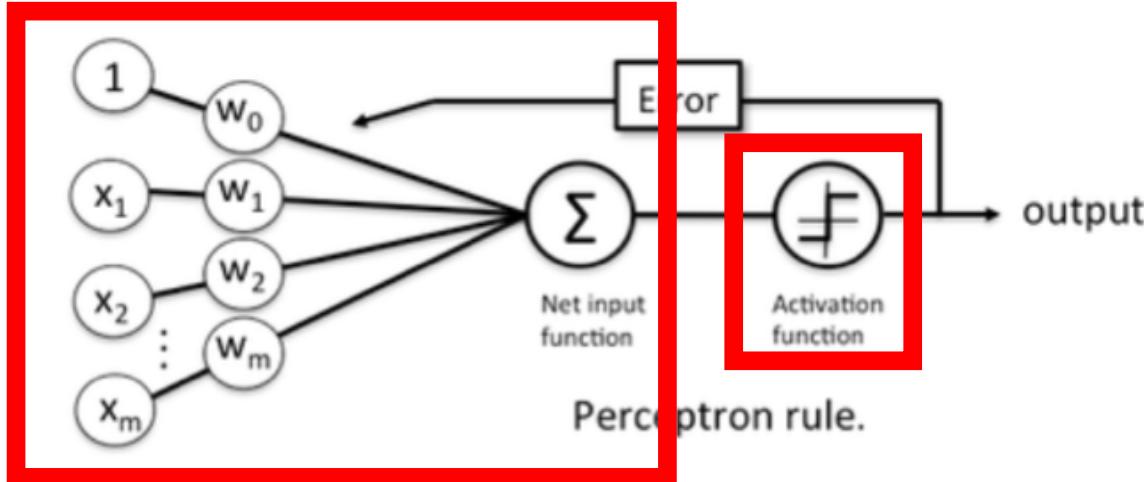


2. Results depend on initial values of weights

Today's Topics

- History of Modern Neural Networks and “Deep Learning”
- Single Layer Neural Network: Perceptron
- Single Layer Neural Network: Adaline
- Variants of Gradient Descent
- Lab

Adaline (ADAptive LInear NEuron): Similarity to Perceptron

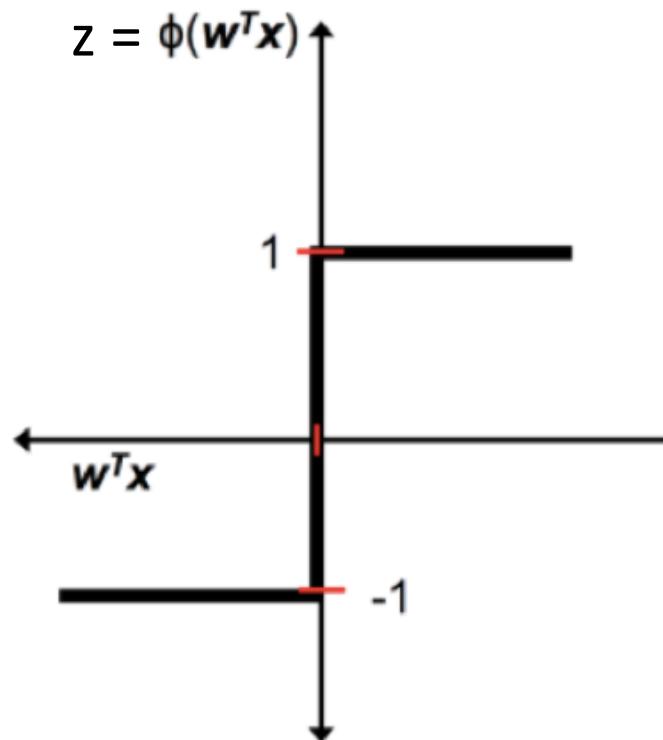


Python Machine Learning; Raschka & Mirjalili

Bernard Widrow and Ted Hoff, An Adaptive “Adaline” Neuron Using Chemical “Memistors”, 1960.

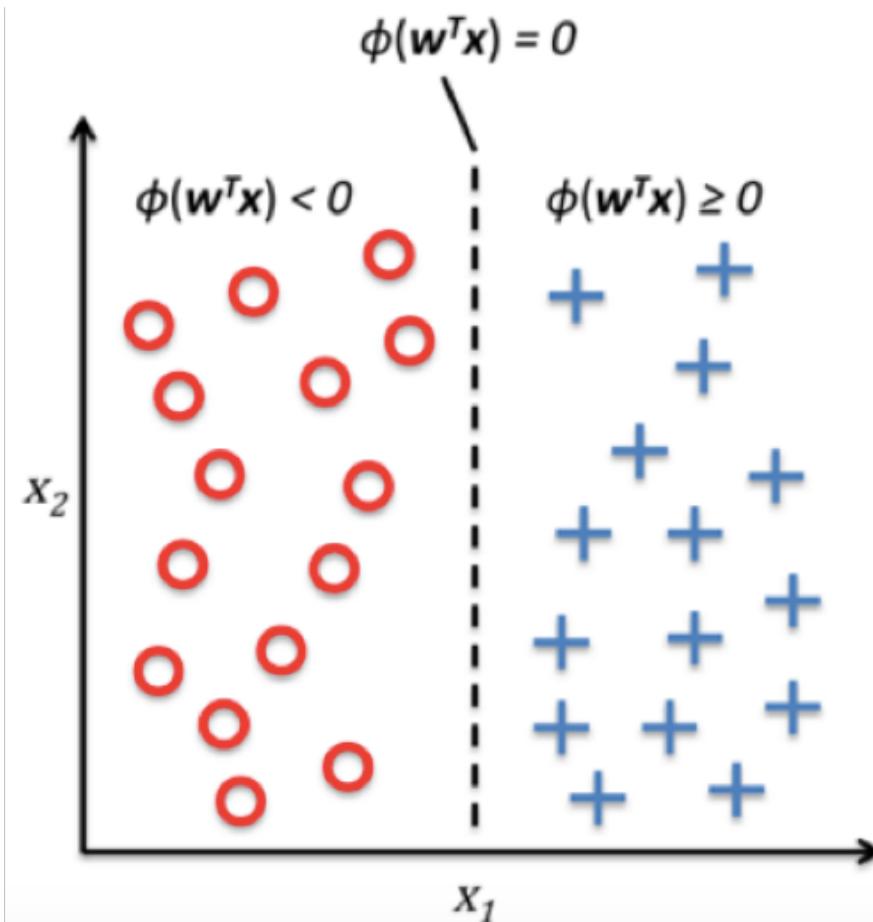
Adaline: Similarity to Perceptron

- Formal definition of artificial neuron: $z = \begin{cases} 1 & \text{if } z \geq \theta \\ -1 & \text{otherwise} \end{cases}$
- Decision function is a unit step function:
 - “fire” or “do not fire”
 - mimics human brain

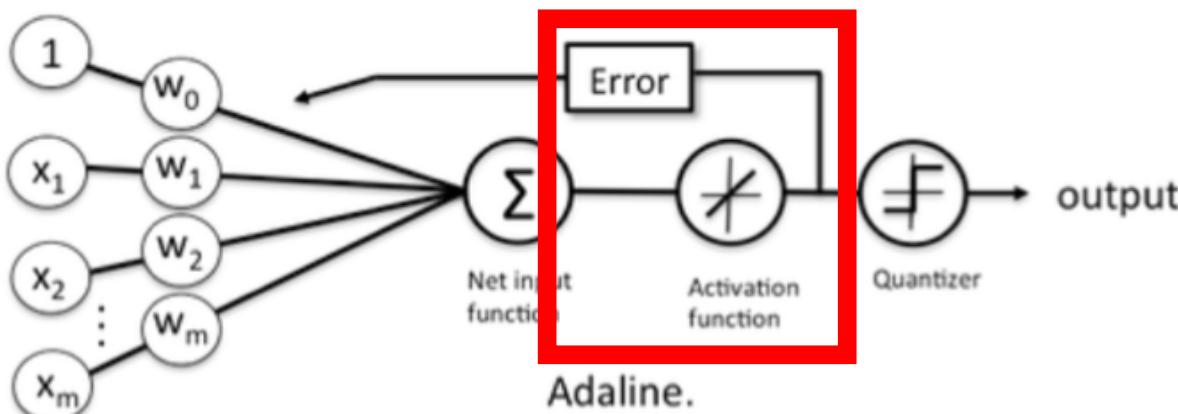
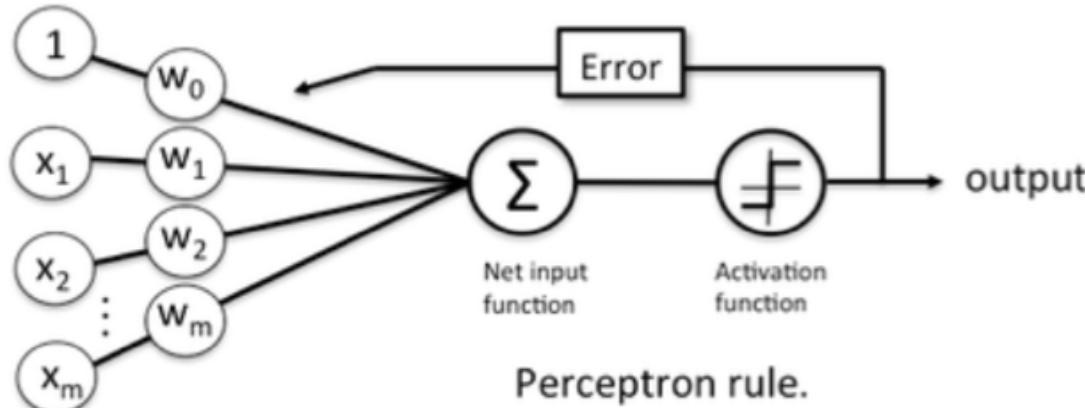


Adaline: Similarity to Perceptron

- Formal definition of artificial neuron: $z = \begin{cases} 1 & \text{if } z \geq \theta \\ -1 & \text{otherwise} \end{cases}$



Adaline: Difference to Perceptron



Adaline Learning Algorithm

1. Initialize the weights to 0 or small random numbers.
2. For k epochs (passes over the training set)
 1. For each training sample
 1. Compute the predicted output value y
 2. Compare predicted to actual output and compute "weight update" value
 3. Update the "weight update" value
 2. Update weights with accumulated "weight update" values

Unlike Perceptron, does not make updates per sample

Adaline Learning Algorithm

1. Initialize the weights to 0 or small random numbers.
2. For k epochs (passes over the training set)
 1. For each training sample
 1. Compute the predicted output value y
 2. Compare predicted to actual output and compute "weight update" value
 3. Update the "weight update" value
 2. Update weights with accumulated "weight update" values: $w_j := w_j + \Delta w_j$

Key Idea: this is differentiable!!!

$$J(\mathbf{w}) = \frac{1}{2} \sum_i (y^{(i)} - \phi(z)_A^{(i)})^2$$

Mathematical
Simplification

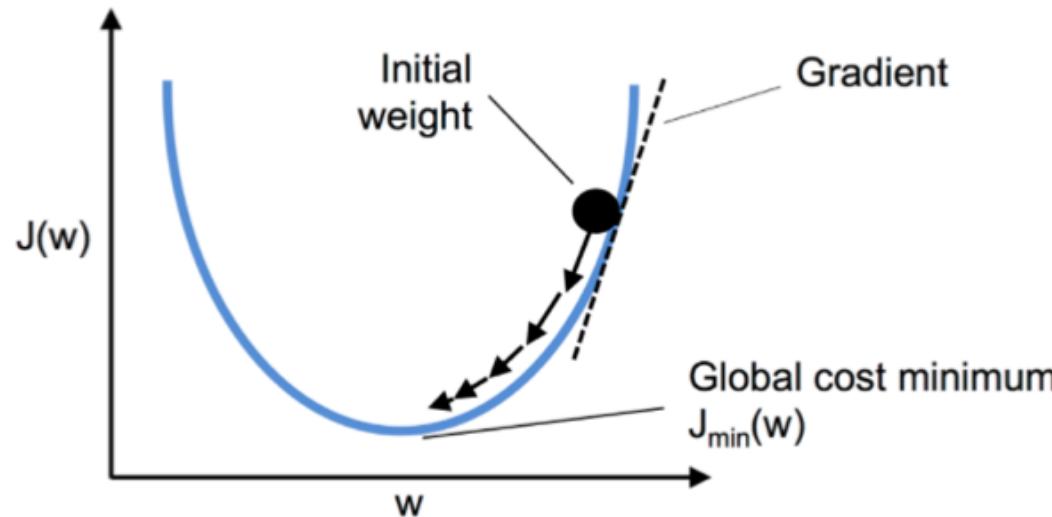
Sum of squared errors

$$\Delta \mathbf{w} = -\eta \nabla J(\mathbf{w})$$

Learning Rate

Take step away from gradient

Adaline Learning Algorithm



2. Update weights with accumulated "weight update" values: $w_j := w_j + \Delta w_j$

Key Idea: this is differentiable!!!

$$J(\mathbf{w}) = \frac{1}{2} \sum_i (y^{(i)} - \phi(z)_A^{(i)})^2$$

Mathematical Simplification

Sum of squared errors

$$\Delta \mathbf{w} = -\eta \nabla J(\mathbf{w})$$

Learning Rate

Take step away from gradient

Derivation of Equation to Update Weights

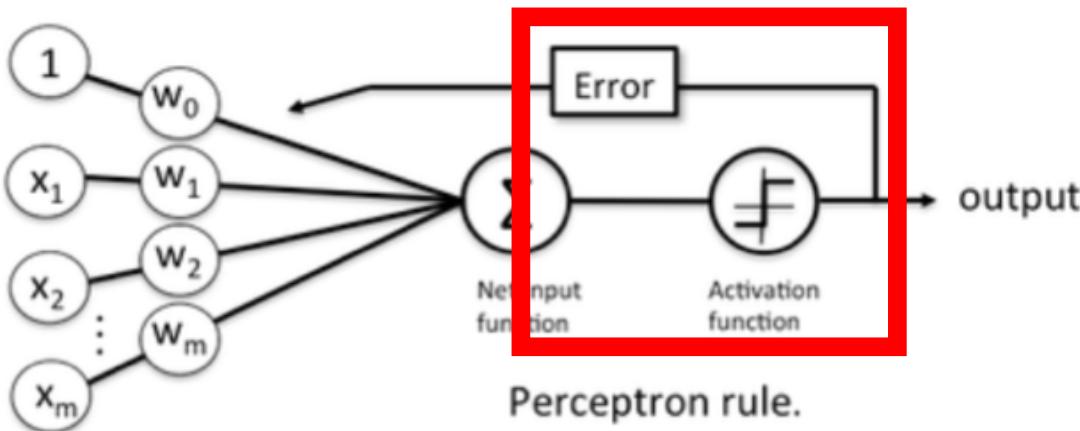
$$\begin{aligned}\frac{\partial J}{\partial w_j} &= \frac{\partial}{\partial w_j} \frac{1}{2} \sum_i (y^{(i)} - \phi(z)_A^{(i)})^2 \\ &= \frac{1}{2} \frac{\partial}{\partial w_j} \sum_i (y^{(i)} - \phi(z)_A^{(i)})^2 \\ &= \frac{1}{2} \sum_i (y^{(i)} - \phi(z)_A^{(i)}) \frac{\partial}{\partial w_j} (y^{(i)} - \phi(z)_A^{(i)}) \\ &= \sum_i (y^{(i)} - \phi(z)_A^{(i)}) \frac{\partial}{\partial w_j} \left(y^{(i)} - \sum_i (w_j^{(i)} x_j^{(i)}) \right) \\ &= \sum_i (y^{(i)} - \phi(z)_A^{(i)}) (-x_j^{(i)}) \\ &= - \sum_i (y^{(i)} - \phi(z)_A^{(i)}) x_j^{(i)}\end{aligned}$$

Adaline Learning Algorithm

$$\begin{aligned}\frac{\partial J}{\partial w_j} &= \frac{\partial}{\partial w_j} \frac{1}{2} \sum_i (y^{(i)} - \phi(z_A^{(i)}))^2 \\ &= \frac{1}{2} \frac{\partial}{\partial w_j} \sum_i (y^{(i)} - \phi(z_A^{(i)}))^2 \\ &= \frac{1}{2} \sum_i (y^{(i)} - \phi(z_A^{(i)})) \frac{\partial}{\partial w_j} (y^{(i)} - \phi(z_A^{(i)})) \\ &= \sum_i (y^{(i)} - \phi(z_A^{(i)})) \frac{\partial}{\partial w_j} \left(y^{(i)} - \sum_i (w_j^{(i)} x_j^{(i)}) \right) \\ &= \sum_i (y^{(i)} - \phi(z_A^{(i)})) (-x_j^{(i)}) \\ &= - \sum_i (y^{(i)} - \boxed{\phi(z_A^{(i)})}) x_j^{(i)}\end{aligned}$$

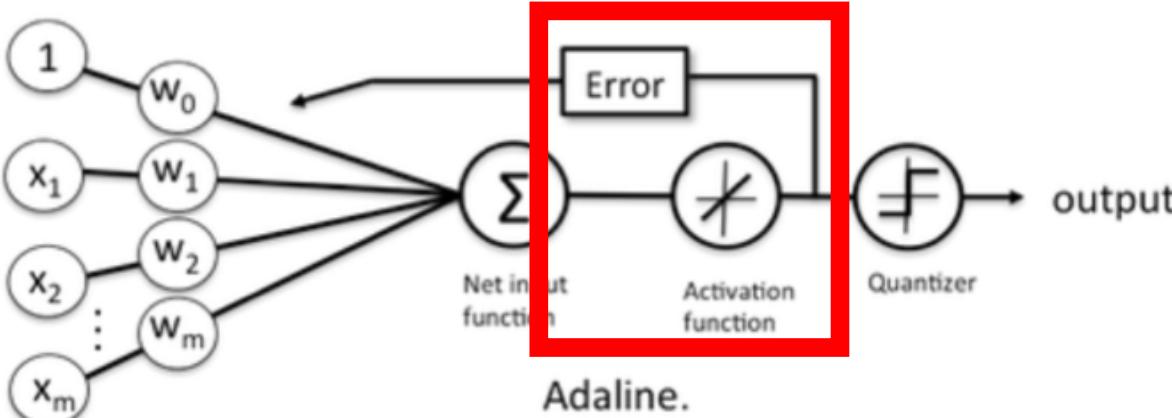
Updates based on continuous valued prediction!

Adaline: Difference to Perceptron



Perceptron rule.

Updates based on continuous valued prediction rather than integer predictions!

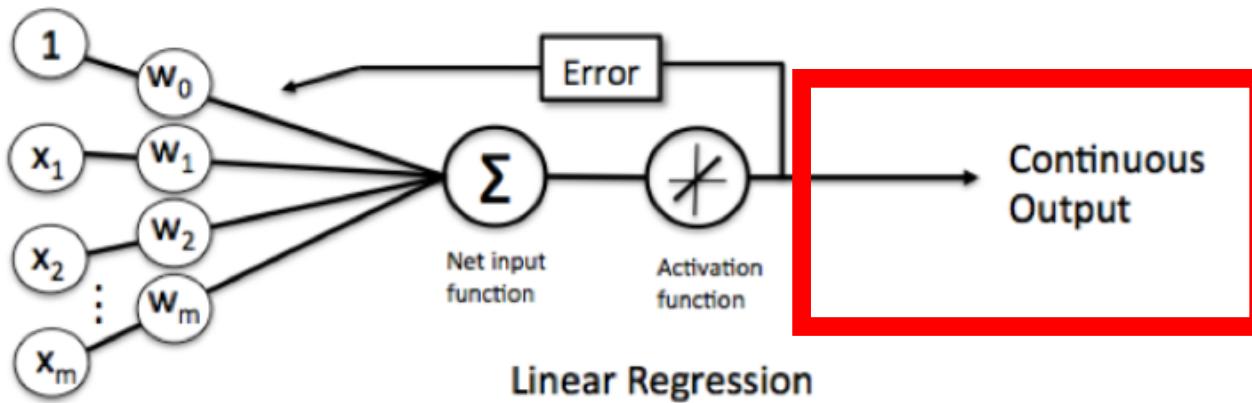


Adaline.

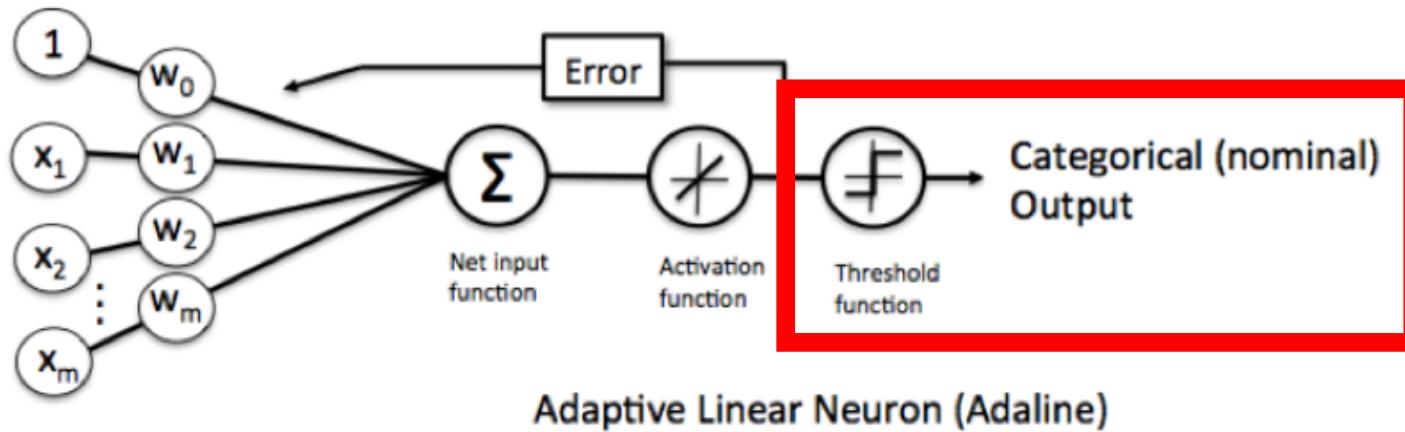
Python Machine Learning; Raschka & Mirjalili

Bernard Widrow and Ted Hoff, An Adaptive “Adaline” Neuron Using Chemical “Memistors”, 1960.

Adaline: Similarity and Difference to Linear Regression



Linear Regression



Adaptive Linear Neuron (Adaline)

Today's Topics

- History of Modern Neural Networks and “Deep Learning”
- Single Layer Neural Network: Perceptron
- Single Layer Neural Network: Adaline
- **Variants of Gradient Descent**
- Lab

Gradient Descent Intuition

- Repeat:
 1. Guess
 2. Calculate error
- e.g., learn linear model for converting kilometers to miles



Miles → **Kilometers = miles x constant** → Kilometers

Gradient Descent Intuition

- Repeat:
 1. Guess
 2. Calculate error
- e.g., learn constant multiplier to convert 1\$ to one yuan

\$10 → Dollar = dollars x **constant**

Gradient Descent Intuition

- Repeat:
 1. Guess
 2. Calculate error
- e.g., learn constant multiplier to convert 1\$ to one yuan

\$10 → Dollar = dollars x **constant** → Error = Guess - Correct

Gradient Descent Intuition

- Repeat:
 1. Guess
 2. Calculate error
- e.g., learn constant multiplier to convert 1\$ to one yuan

\$10 → Dollar = dollars x **constant**

Gradient Descent Intuition

- Repeat:
 1. Guess
 2. Calculate error
- e.g., learn constant multiplier to convert 1\$ to one yuan

\$10 → Dollar = dollars x **constant** → Error = Guess - Correct

Gradient Descent Intuition

- Repeat:
 1. Guess
 2. Calculate error
- e.g., learn constant multiplier to convert 1\$ to one yuan

\$10 → Dollar = dollars x **constant**

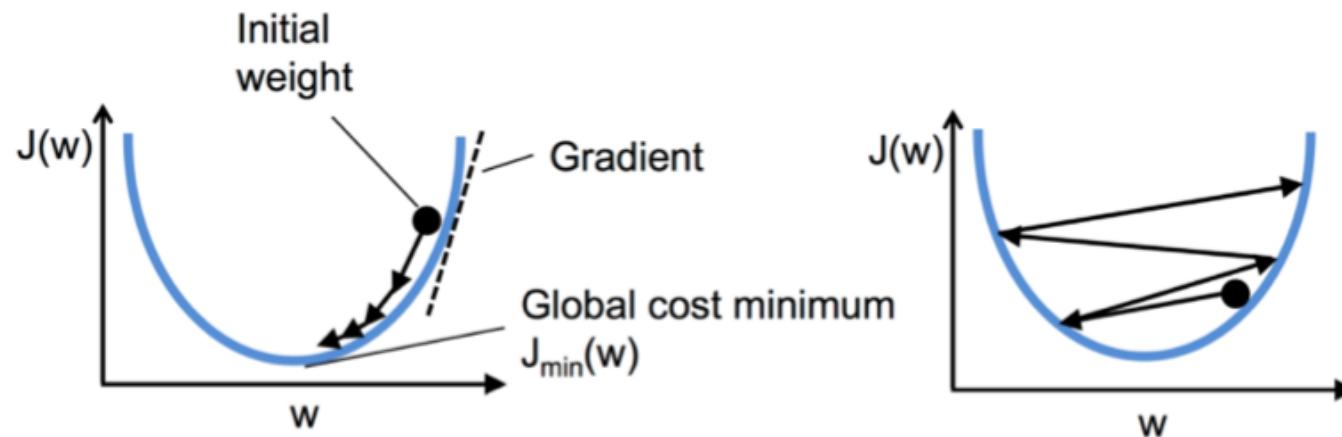
Gradient Descent Intuition

- Repeat:
 1. Guess
 2. Calculate error
- e.g., learn constant multiplier to convert 1\$ to one yuan



- Idea: iteratively adjust **constant** (i.e., model parameter) to try to make error get smaller

Gradient Descent: Influence of Learning Rate



- Learning Rate: amount new evidence is prioritized when updating weights
- What happens when learning rate is too small?
 - Convergence to good solution will be slow!
- What happens when learning rate is too large?
 - May not be able to converge to a good solution
- How to address the cons of different learning rates?
 - Gradually reduce learning rate over time

Stochastic Gradient Descent (SGD)

- For each step (update), use calculations from ***one training example***
- What are strengths of this approach?
 - Each iteration is fast to compute
 - Can train using huge datasets (stores one instance in memory at each iteration)
- What are weaknesses of this approach?
 - Updates will bounce a lot
- Which algorithm uses this?
 - Perceptron

Batch Gradient Descent (BGD)

- For each step (update), use calculations over ***all training examples***
- What are strengths of this approach?
 - Does not bounce too much
- What are weaknesses of this approach?
 - Very slow or infeasible when dataset is large
- Which algorithm uses this?
 - Adaline

Mini-batch Gradient Descent

- For each step (update), use calculations over *subset of training examples*
- What are strengths of this approach?
 - Bounces less erratically when finding model parameters than SGD
 - Can train using huge datasets (store some instances in memory at each iteration)
- What are weaknesses of this approach?
 - Very slow or infeasible when dataset is large
- Which algorithm uses this?
 - To be explored in future classes

Today's Topics

- History of Modern Neural Networks and “Deep Learning”
- Single Layer Neural Network: Perceptron
- Single Layer Neural Network: Adaline
- Variants of Gradient Descent
- Lab