



The Perceptron

Instructor: Jessica Wu -- Harvey Mudd College

The instructor gratefully acknowledges Andrew Ng (Stanford), Eric Eaton (UPenn), David Kauchak (Pomona), and the many others who made their course materials freely available online.

Robot Image Credit: Viktoriya Sukhanova © 123RF.com

Perceptron Basics

Learning Goals

- Describe the perceptron model
- Describe the perceptron algorithm
- Describe why the perceptron update works
- Describe the perceptron cost function
- Describe how a bias term affects the perceptron

The Perceptron

Assume there is a linear classifier

Start with simple learner and analyze with it does

Let $y \in \{-1, +1\}$

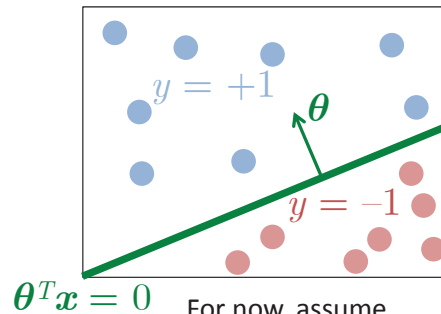
$$h_{\theta}(\mathbf{x}) = g(\theta^T \mathbf{x})$$

where

$$g(z) = \text{sgn}(z) \begin{cases} +1, & z \geq 0 \\ -1, & z < 0 \end{cases}$$

So if $\theta^T \mathbf{x} \geq 0$, then $y = +1$

if $\theta^T \mathbf{x} < 0$, then $y = -1$



For now, assume hyperplane through origin (no bias term).

Perceptron Algorithm

start with guess for θ (typically $\theta = \mathbf{0}$)

repeat until convergence

for $i = 1$ to n

if $y^{(i)} \theta^T \mathbf{x}^{(i)} \leq 0$ ← if mistake is made (why ≤ 0 rather than < 0 ?)

$\theta \leftarrow \theta + y^{(i)} \mathbf{x}^{(i)}$ ← update step

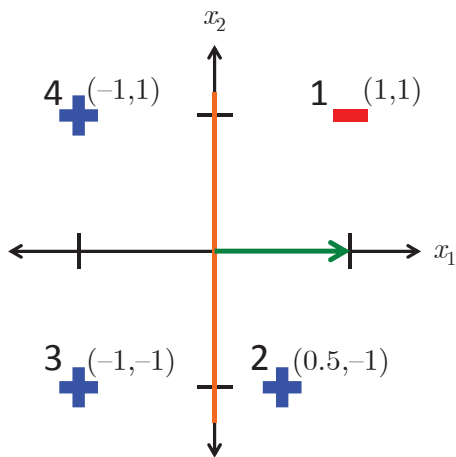
possible criteria:

- all training examples correctly classified
- if average update $(\|\theta_{\text{new}} - \theta_{\text{old}}\|_2) < \epsilon$
- fixed # of iterations
- single pass over data

Notes

- online learning algorithm
- guaranteed to find separating hyperplane if data is linearly separable (theorem later this lecture)

Perceptron Example



$$\theta_0 = [1, 0]^T$$

Repeat until convergence

Process points in order 1,2,3,4

Keep track of θ as it changes

Redraw the hyperplane after each step

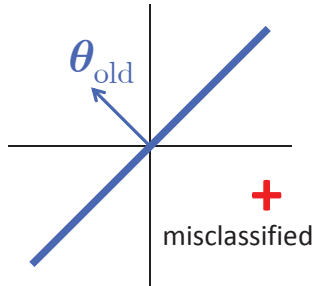
Based on slide by David Kauchak [originally by Piyush Rai]



(This slide intentionally left blank.)

Why the Perceptron Update Works

Geometric Interpretation



Based on slide by Eric Eaton [originally by Piyush Rai]



Why the Perceptron Update Works

Mathematic Proof

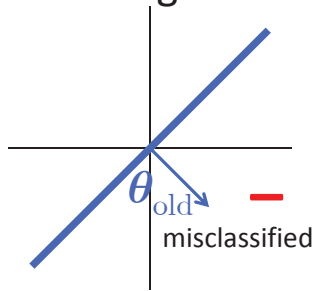
- Consider the misclassified example $y = +1$
 - Perceptron wrongly thinks $\theta_{old}^T \mathbf{x} < 0$

Based on slide by Eric Eaton [originally by Piyush Rai]



Why the Perceptron Update Works

similar arguments for misclassified negative example



- Consider the misclassified example $y = -1$
 - Perceptron wrongly things $\theta_{old}^T \mathbf{x} > 0$



Why the Perceptron Update Works

- Consider the misclassified example $y = -1$
 - Perceptron wrongly things $\theta_{old}^T \mathbf{x} > 0$



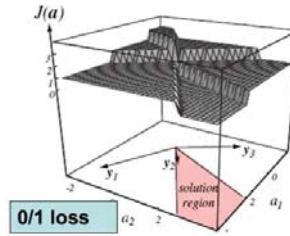
Cost Function

What if we tried to minimize 0/1 loss?

$$J(\theta) = \sum_{i=1}^n \mathbb{I}[[y^{(i)}\theta^T \mathbf{x}^{(i)} \leq 0]]$$

indicator function
0 if prediction is correct
1 otherwise

remember, prediction is correct if $y^{(i)}\theta^T \mathbf{x}^{(i)} > 0$



0/1 loss

Challenges

- NP-hard!
- small changes in θ can induce large changes in $J(\theta)$, and change is not continuous
- lots of local minima
- not useful gradient: at any point, no information to direct us towards any minima

Idea: Use surrogate loss function

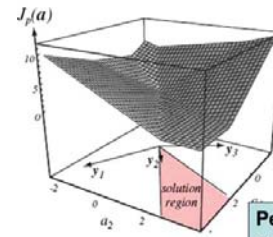
$$J(\theta) = \sum_{i=1}^n \max(0, -y^{(i)}\theta^T \mathbf{x}^{(i)})$$

If prediction is correct

- $\max(0, -y^{(i)}\theta^T \mathbf{x}^{(i)}) = 0$.

Otherwise

- it is confidence of mis-prediction
- walk in direction of negative gradient (GD!)
(for single example, $-\nabla_{\theta} J(\theta) = y^{(i)} \mathbf{x}^{(i)}$)



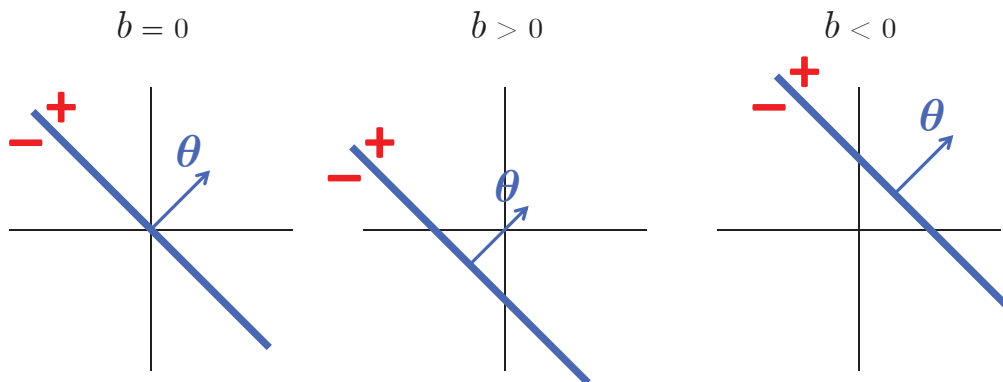
Perceptron criterion

Based on slide by Eric Eaton [originally by Alan Fern]

Bias Term

$$h_{\theta}(\mathbf{x}) = \text{sgn}(\theta^T \mathbf{x} + b)$$

Bias **shifts** hyperplane $-b$ units in direction of θ .



Positive bias means more examples classified as positive.

Shift away from θ means more space for positive classification.

Perceptron Convergence

Learning Goals

Does the perceptron converge?
If so, how long does it take?
(How many mistakes / updates?)

- State the perceptron convergence theorem
- State the implications of the theorem

Margins

For training set $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$ and parameter vector θ :

- **functional margin** of i^{th} example

$$\hat{\gamma}^{(i)} = y^{(i)} (\theta^T \mathbf{x} + b)$$

- this is positive if θ classifies $\mathbf{x}^{(i)}$ correctly
- absolute value = “confidence” in predicted label (or “mis-confidence”)

- **geometric margin** of i^{th} example

$$\gamma^{(i)} = \frac{\hat{\gamma}^{(i)}}{\|\theta\|} = \frac{y^{(i)} (\theta^T \mathbf{x} + b)}{\|\theta\|}$$

- signed distance of example to hyperplane
(positive if example classified correctly)

- **margin of training set** = minimum geometric margin

$$\gamma = \min_i \gamma^{(i)}$$

Perceptron Convergence Theorem

Theorem (Block & Novikoff) [$b = 0$]

Assume

- $\exists \theta^*$ s.t. $\|\theta^*\| = 1$ and $\forall i, \gamma^{(i)} \geq \gamma^* > 0$
 - data is linearly separable with margin γ^* by unit-norm hyperplane θ^*
- $\forall i, \|\mathbf{x}^{(i)}\| \leq R$
 - examples are not “too big”

Then

- perceptron converges after at most $\frac{R^2}{(\gamma^*)^2}$ updates.

Convergence Theorem Guarantees

- Convergence rate
 - depends on margin γ^* and “size” of data R
 - but not on number of training examples n or data dimensionality d
- If perceptron is given data that is linearly separable with margin γ^* , it will converge to a solution that separates data and converge quickly if γ^* is large

Proof

- We assumed that θ^* (a hyperplane consistent with data) exists
 - the classification problem is “easy” if γ^* is large
 - optimal case (maximum γ^*): θ^* is the **maximum-margin separator**
- But we are not guaranteed to find θ^* using perceptron algorithm
 - we will show how to find maximum-margin separator θ^* using SVMs

Order of Examples

Does the order in which we traverse examples matter?



The Name “Perceptron”

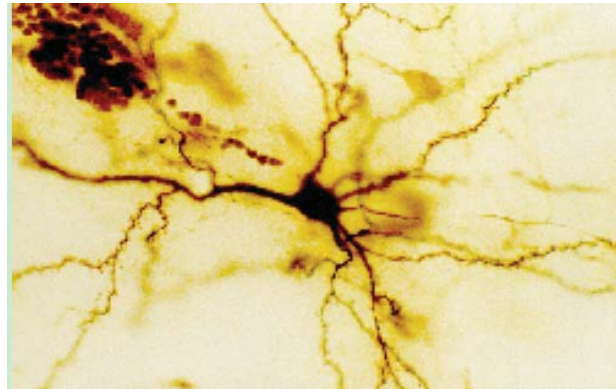
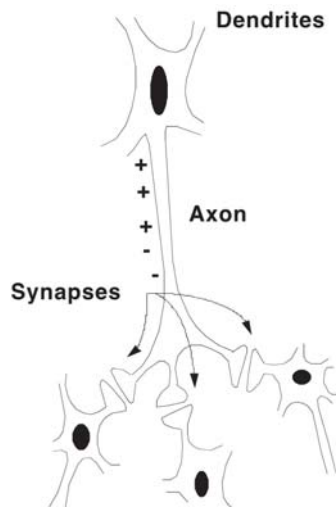
Learning Goals

Why is it called the “perceptron” learning algorithm if it learns a line?

Why not “line learning” algorithm?

The Name “Perceptron”

Comes from our nervous system



neuron

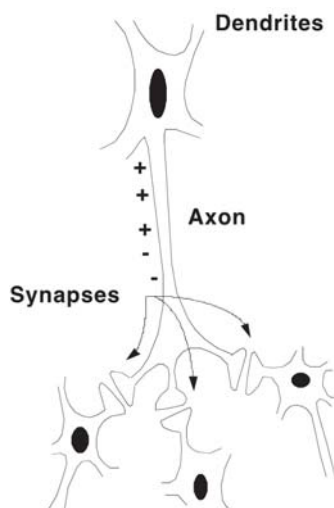
Based on slide by David Kauchak

Our Nervous System:

the human brain is a large collection of interconnected neurons

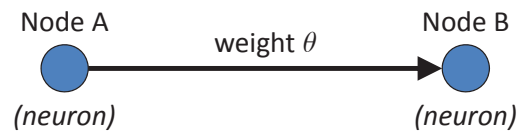
a **neuron** is a brain cell

- collects, processes, and disseminates electrical signals
- connected via synapses
- **fire** depending on conditions of neighboring neurons



Based on slide by David Kauchak

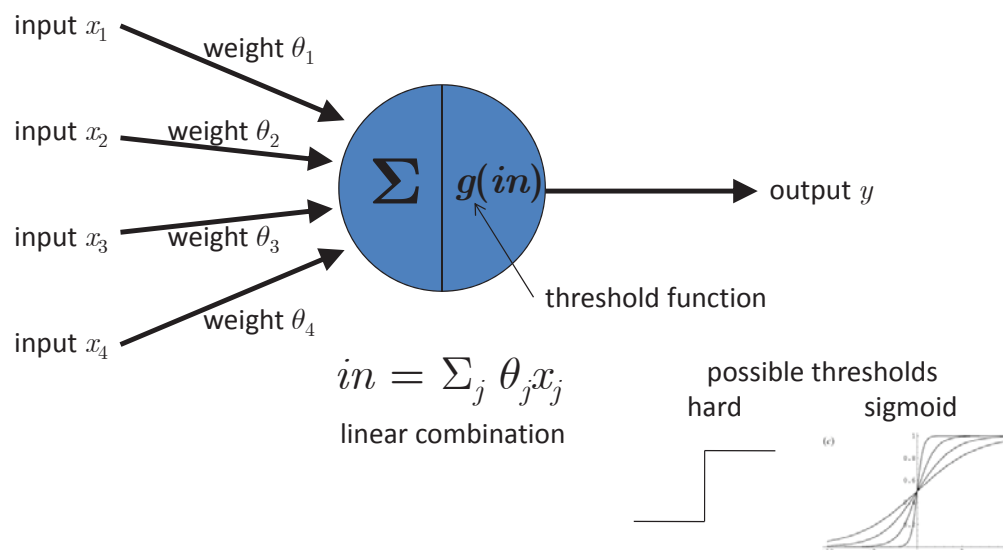
A Neuron



- θ is the strength of signal sent from A to B
- if A fires and θ is positive, then A stimulates B
- if A fires and θ is negative, then A inhibits B
- if B is stimulated enough, then it also fires
- amount of stimulation required is determined by its threshold

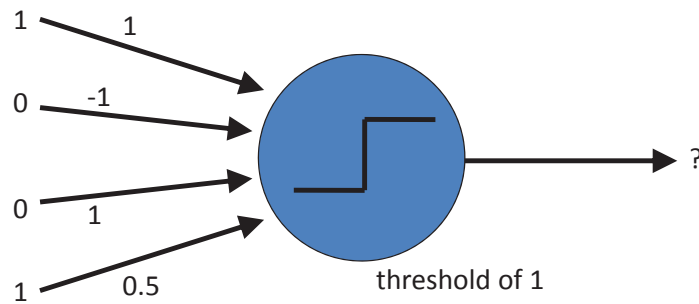
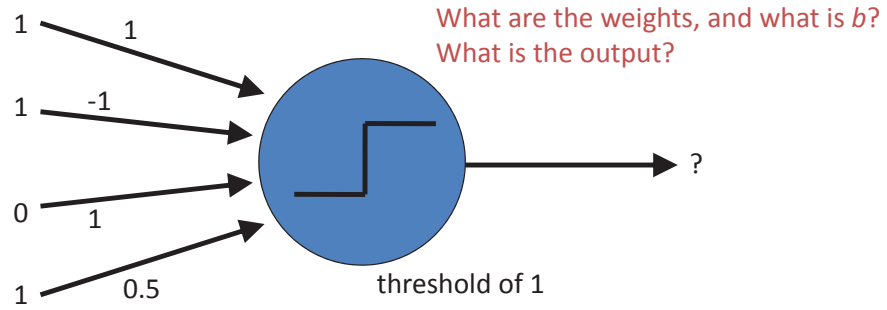
Based on slide by David Kauchak

A Single Neuron / Perceptron



Based on slide by David Kauchak

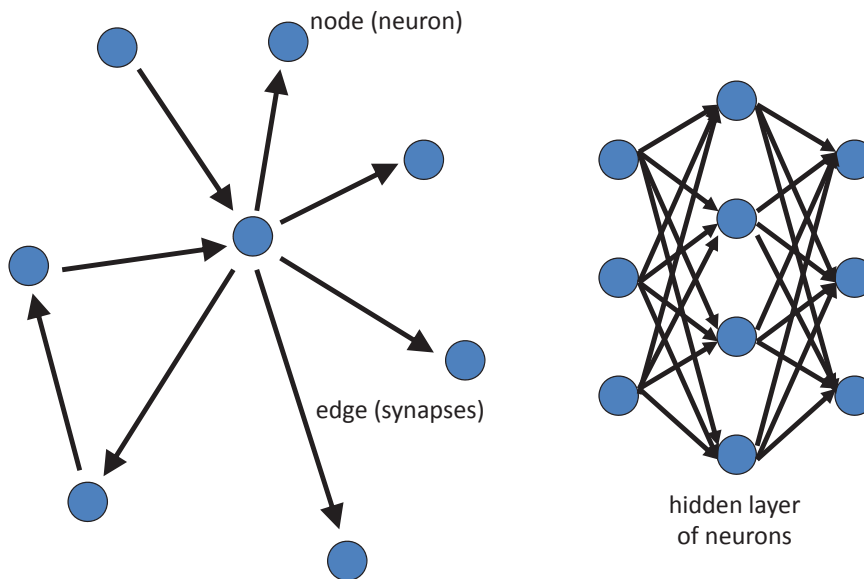
Neuron Example



Based on slide by David Kauchak



Neural Networks



Based on slide by David Kauchak

(extra slides)

Proof of Perceptron Convergence Theorem

Learning Goals

- Prove the Perceptron Convergence Theorem
=D

Proof Overview

- $\exists \theta^*$ s.t. data is linearly separable with margin γ^*
 - (we do not know θ^* but we know that it exists)
- perceptron algorithm tries to find θ that points roughly in same direction as θ^*
 - for large γ^* , “roughly” is very rough
 - for small γ^* , “roughly” is very precise
- every update, angle between θ and θ^* changes

recall $\cos \alpha = \frac{\mathbf{u}^T \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$, we will prove that

- $\mathbf{u}^T \mathbf{v}$ increases a lot
- $\|\mathbf{u}\|$ and $\|\mathbf{v}\|$ do not increase very much
- so angle α decreases each update

Proof by Induction



(This slide intentionally left blank.)