## Machine Learning: Perceptrons! CS 445/545





## Oops!

We're very sorry, but we're having trouble doing what you just asked us to do. Please give us another chance--click the Back button on your browser and try your request again. Or start from the beginning on our <u>homepage</u>.



500 Service Unavailable Error - Windows Internet Explorer





## Oops!

We're very sorry, but we're h asked us to do. Please give button on your browser and from the beginning on our ho





Information flow through neurons

#### Dendrites Collect electrical signals

Cell body Contains nucleus and organelles

#### Axon

Passes electrical signals on to dendrites of another cell or to an effector cell



## Hebb's Postulate

"When an axon of cell A is near enough to excite a cell B and repeatedly or persistently takes part in firing it, some growth process or metabolic change takes place in one or both cells such that A's efficiency, as one of the cells firing B, is increased."

 In other words: if two neurons fire "close in time" then strength of synaptic connection between them increases.

$$\Delta w_{ij}(t) = \eta v_i v_j g(t_{v_i}, t_{v_j})$$





Weights reflect correlation between firing events.

- Human brain contains  $\sim 10^{11}$  neurons
- Each individiaul neuron connects to  $\sim 10^4$  neuron
- $\sim 10^{14}$  total synapses!

	Brain	Computer
Number of Processing Units	$\approx 10^{11}$	$pprox 10^9$
Type of Processing Units	Neurons	Transistors
Form of Calculation	Massively Parallel	Generally Serial
Data Storage	Associative	Address-based
Response Time	$\approx 10^{-3} {\rm s}$	$\approx 10^{-9} {\rm s}$
Processing Speed	Very Variable	Fixed
Potential Processing Speed	$\approx 10^{13}$ FLOPS $^{14}$	$\approx 10^{18} \; \mathrm{FLOPS}$
Real Processing Speed	$\approx 10^{12} \; {\rm FLOPS}$	$\approx 10^{10} \; \mathrm{FLOPS}$
Resilience	Very High	Almost None
Power Consumption per Day	20W	$300W$ $^{15}$

### – Is the singularity near?







TYPE I CIVILIZATION harnesses all the resources of a planet. Carl Sagan estimated that Earth rates about 0.7 on the scale.



TYPE II CIVILIZATION harnesses all the radiation of a star. Humans might reach Type II in a few thousand years.



#### **TYPE III CIVILIZATION**

harnesses all the resources of a galaxy. Humans might reach Type III in a few hundred thousand to a million years.

### The Kardashev Scale -

## McCulloch & Pitts Neuron Model (1943)







(3) Components:

(1) Set of weighted inputs  $\{w_i\}$  that correspond to synapses

(2) An "adder" that sums the input signals (equivalent to membrane of the cell that collects the electrical charge)

(3) An activation function (initially a threshold function) that decides whether the neuron fires ("spikes") for the current inputs.

## McCulloch & Pitts Neuron Model (1943)

– Limitations & Deviations of the M-P Neuron Model:

(\*) Summing is linear.

(\*) No explicit model of "spike trains" (sequence of pulses that encodes information in biological neuron).

(\*) Threshold value is usually fixed.

(\*) Sequential updating implicit (biological neurons usually update themselves asynchronously)

(\*) Weights can be positive (excitatory) or negative (inhibitory); biological neurons do not change in this way.

(\*) Real neurons can have synapses that link back to themselves (e.g. feedback loop) – see RNNs (recurrent neural networks).

(\*) Other biological aspects ignored: chemical concentrations, refractory periods, etc.

input

 $x_2$ 

 $X_n$ 

 $x_1$ 

w<sub>n</sub>

 $W_1$ 

output

v

Input is  $(x_1, x_2, ..., x_n)$ 

Weights are  $(w_1, w_2, \dots, w_n)$ 

Output *y* is 1 ("the neuron fires") if the sum of the inputs times the weights is greater or equal to the threshold:

Input is  $(x_1, x_2, ..., x_n)$ 

Weights are  $(w_1, w_2, \dots, w_n)$ 

Output *y* is 1 ("the neuron fires") if the sum of the inputs times the weights is greater or equal to the threshold:

If  $w_1x_1 + w_2x_2 + \dots + w_nx_n > threshold$ then y = 1, else y = 0



 $x_2$ 

 $X_{n}$ 

 $x_1$ 

w<sub>n</sub>

W<sub>1</sub>

output

v

Input is  $(x_1, x_2, ..., x_n)$ 

Weights are  $(w_1, w_2, \dots, w_n)$ 

Output *y* is 1 ("the neuron fires") if the sum of the inputs times the weights is greater or equal to the threshold:

*If*  $w_1x_1 + w_2x_2 + ... + w_nx_n > threshold$ *then* y = 1, *else* y = 0

#### input



 $w_0$  is called the "bias"

 $-w_0$  is called the "**threshold**"

#### input



Input is  $(x_1, x_2, ..., x_n)$ 

Weights are  $(w_1, w_2, \dots, w_n)$ 

Output *y* is 1 ("the neuron fires") if the sum of the inputs times the weights is greater or equal to the threshold:

*If*  $w_1x_1 + w_2x_2 + ... + w_nx_n > threshold$ *then* y = 1, *else* y = 0

 $w_0$  is called the "**bias**"

 $-w_0$  is called the "**threshold**"

#### input



Input is  $(x_1, x_2, ..., x_n)$ 

Weights are  $(w_1, w_2, \dots, w_n)$ 

Output *y* is 1 ("the neuron fires") if the sum of the inputs times the weights is greater or equal to the threshold:

*If*  $w_1x_1 + w_2x_2 + ... + w_nx_n > threshold$ *then* y = 1, *else* y = 0*If*  $w_1x_1 + w_2x_2 + ... + w_nx_n > -w_0$ 

*then* y = 1, *else* y = 0

 $w_0$  is called the "bias"

 $-w_0$  is called the "**threshold**"

#### input



Input is  $(x_1, x_2, ..., x_n)$ 

Weights are  $(w_1, w_2, \dots, w_n)$ 

Output *y* is 1 ("the neuron fires") if the sum of the inputs times the weights is greater or equal to the threshold:

*If*  $w_1x_1 + w_2x_2 + ... + w_nx_n > threshold$  *then* y = 1, *else* y = 0 *If*  $w_1x_1 + w_2x_2 + ... + w_nx_n > -w_0$ *then* y = 1, *else* y = 0

*If*  $w_0 + w_1 x_1 + w_2 x_2 + ... + w_n x_n > 0$ *then* y = 1, *else* y = 0

Q: Why do we introduce a bias term?

 $w_0$  is called the "bias"

 $-w_0$  is called the "**threshold**"

#### input



Input is  $(x_1, x_2, ..., x_n)$ 

Weights are  $(w_1, w_2, \dots, w_n)$ 

Output *y* is 1 ("the neuron fires") if the sum of the inputs times the weights is greater or equal to the threshold:

*If*  $w_1x_1 + w_2x_2 + ... + w_nx_n > threshold$  *then* y = 1, *else* y = 0 *If*  $w_1x_1 + w_2x_2 + ... + w_nx_n > -w_0$ *then* y = 1, *else* y = 0

*If*  $w_0 + w_1 x_1 + w_2 x_2 + ... + w_n x_n > 0$ *then* y = 1, *else* y = 0

 $w_0$  is called the "bias"

 $-w_0$  is called the "**threshold**"

#### input



Input is  $(x_1, x_2, ..., x_n)$ 

Weights are  $(w_1, w_2, \dots, w_n)$ 

Output *y* is 1 ("the neuron fires") if the sum of the inputs times the weights is greater or equal to the threshold:

*If*  $w_1x_1 + w_2x_2 + ... + w_nx_n > threshold$  *then* y = 1, *else* y = 0 *If*  $w_1x_1 + w_2x_2 + ... + w_nx_n > -w_0$ *then* y = 1, *else* y = 0

*If*  $w_0 + w_1 x_1 + w_2 x_2 + ... + w_n x_n > 0$ *then* y = 1, *else* y = 0

*If*  $w_0 x_0 + w_1 x_1 + w_2 x_2 + ... + w_n x_n > 0$ *then* y = 1, *else* y = 0

 $w_0$  is called the "bias"

 $-w_0$  is called the "**threshold**"

#### input



Input is  $(x_1, x_2, ..., x_n)$ 

Weights are  $(w_1, w_2, \dots, w_n)$ 

Output *y* is 1 ("the neuron fires") if the sum of the inputs times the weights is greater or equal to the threshold:

*output* =  $y(\mathbf{x}) = a(w_0x_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n)$ 

where  $a(z) = \begin{cases} 0 & \text{if } z \le 0 \\ 1 & \text{if } z > 0 \end{cases}$ 

 $w_0$  is called the "bias"

 $-w_0$  is called the "**threshold**"

#### input



Input is  $(x_1, x_2, ..., x_n)$ 

Weights are  $(w_1, w_2, \dots, w_n)$ 

Output *y* is 1 ("the neuron fires") if the sum of the inputs times the weights is greater or equal to the threshold:

*output* =  $y(\mathbf{x}) = a(w_0x_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n)$ 

where 
$$a(z) = \begin{cases} 0 & if \ z \le 0 \\ 1 & if \ z > 0 \end{cases}$$

Let  $\mathbf{x} = (x_0, x_1, x_2, \dots, x_n)$  $\mathbf{w} = (w_0, w_1, w_2, \dots, w_n)$ 

 $w_0$  is called the "bias"

 $-w_0$  is called the "**threshold**"

 $x_1$ 

 $w_n$ 

w<sub>1</sub>

 $x_0$ 

+1

Wo

v

output

#### input

 $x_2$ 



Weights are  $(w_1, w_2, \dots, w_n)$ 

Output y is +1 ("the neuron fires") if the sum of the inputs times the weights is greater or equal to the threshold:

*output* =  $y(\mathbf{x}) = a(w_0x_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n)$ 

where  $a(z) = \begin{cases} 0 & \text{if } z \le 0 \\ 1 & \text{if } z > 0 \end{cases}$ 

Let  $\mathbf{x} = (x_0, x_1, x_2, \dots x_n)$  $\mathbf{w} = (w_0, w_1, w_2, \dots w_n)$ 

Then:

 $y = a(\mathbf{W} \times \mathbf{X})$ 

## **Decision Surfaces**

- Assume data can be separated into two classes, positive and negative, by a linear *decision surface*.
- Assuming data is *n*-dimensional, a perception represents a (*n*-1)-dimensional hyperplane that separates the data into two classes, *positive* (1) and *negative* (0).



Feature 2





Feature 2

# Example where line won't work?

 $\bigcirc$ 

Feature 1

Feature 2

# Example

What is the predicted class y?



# Example

• What is the predicted class y?



 $y = a(\vec{w} \cdot \vec{x}) = a((-.1)(1) + .4(1) + (-.4)(-1)) = a(0.7) = 1$ 

# **Geometry of the perceptron**



input



input



Input instance:  $\mathbf{x}^k = (x_1, x_2, \dots, x_n)$ , with target class  $t^k \mid \{0, 1\}$ 

 $x_0$  $x_1$ +1  $w_1$ Wo  $x_2$ output y w<sub>n</sub>

input

Goal is to use the training data to learn a set of weights that will:

(1) correctly classify the training data

(2) generalize to unseen data

Input instance:  $\mathbf{x}^k = (x_1, x_2, \dots, x_n)$ , with target class  $t^k \mid \{0, 1\}$ 

 $X_n$ 

Learning is often framed as an optimization problem:

• Find w that minimizes average "loss":

$$J(\mathbf{w}) = \frac{1}{M} \mathop{\overset{M}{\stackrel{}_{a=1}}}\limits^{M} L(\mathbf{w}, \mathbf{x}^{k}, t^{k})$$

where M is number of training examples and L is a "loss" function.

One part of the "art" of ML is to define a good loss function.

• Here, define the loss function as follows:

Let  $y = a(\mathbf{w} \times \mathbf{x})$  $L(\mathbf{w}, \mathbf{x}^{k}, t^{k}) = \frac{1}{2}(t^{k} - y)^{2}$  "squared loss"

How to solve this minimization problem? Gradient descent.

# Aside: Convex Regions

 Convex: for any pair of points x<sub>a</sub> and x<sub>b</sub> within a region, every point x<sub>c</sub> on a line between x<sub>a</sub> and x<sub>b</sub> is in the region



# Aside: Convex Functions

 Convex: for any pair of points x<sub>a</sub> and x<sub>b</sub> within a region, every point x<sub>c</sub> on a line between x<sub>a</sub> and x<sub>b</sub> is in the region


#### **Aside: Convex Functions**

 Convex: for any pair of points x<sub>a</sub> and x<sub>b</sub> within a region, every point x<sub>c</sub> on a line between x<sub>a</sub> and x<sub>b</sub> is in the region



#### **Aside: Convex Functions**

- Convex functions have a single maximum and minimum!
- How does this help us?
- (nearly) Guaranteed optimality of Gradient
   Descent

 $x_c$ 

 $x_a$ 

- The Gradient is defined (though we can't solve **directly)**  $\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 0$ • Points in the direction of fastest increase

 $\nabla_{\theta} R$ 

- Gradient points in the direction of fastest increase  $\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 0$ • To minimize R, move in the opposite direction



- Gradient points in the direction of fastest increase  $\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 0$ • To minimize R, move in the opposite direction

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- (nearly) guaranteed to converge to the minimum

$$-\nabla_{\theta}R$$

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 20$$

 $-\nabla_{\theta} R$ 

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- (nearly) guaranteed to converge to the minimum

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 30$$

Initialize Randomly

$$\theta_0 = random$$

• Update with small steps  $\theta_{t+1} = \theta_t - \eta \nabla_{\theta} R|_{\theta_t}$ 

 $-\nabla_{\theta}R$ 

 (nearly) guaranteed to converge to the minimum

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 40$$

Initialize Randomly

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- (nearly) guaranteed to converge to the minimum

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 50$$

Initialize Randomly

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- (nearly) guaranteed to converge to the minimum

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 0$$

Initialize Randomly

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- (nearly) guaranteed to converge to the minimum

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 0$$

Initialize Randomly

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- (nearly) guaranteed to converge to the minimum

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 0$$

 $-\nabla_{\theta}R$ 

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- Can oscillate if η is too large

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 0$$

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- Can oscillate if η is too large



 $\nabla_{\theta} R$ 

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- Can oscillate if η is too large



 $-\nabla_{\theta}R$ 

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- Can oscillate if η is too large

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 20$$

 $-\nabla_{\theta}R$ 

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- Can oscillate if η is too large

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 30$$

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- Can oscillate if η is too large



 $\nabla_{\theta} R$ 

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- Can oscillate if η is too large

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 50$$

 $-\nabla_{\theta}R$ 

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- Can oscillate if η is too large

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 0$$

 $-\nabla_{\theta}R$ 

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- Can oscillate if η is too large

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 0$$

Initialize Randomly

$$\theta_0 = random$$

- Update with small steps  $\theta_{t+1} = \theta_t \eta \nabla_{\theta} R|_{\theta_t}$
- Can stall if  $-\nabla_{\theta}R$  is ever 0 not at the minimum

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 0$$

• Find w that minimizes average "loss":

$$J(\mathbf{w}) = \frac{1}{M} \mathop{\text{a}}\limits_{k=1}^{M} L(\mathbf{w}, \mathbf{x}^{k}, t^{k})$$

where M is number of training examples and L is a "loss" function.

• Here, define the loss function as follows:

Let 
$$y = a(\mathbf{w} \times \mathbf{x})$$
  
 $L(\mathbf{w}, \mathbf{x}^k, t^k) = \frac{1}{2} (t^k - y)^2$  "squared loss"

How to solve this minimization problem? Gradient descent.

• To find direction of steepest descent, take the derivative of *J*(**w**) with respect to **w**.

• A vector derivative is called a "gradient":  $\nabla J(\mathbf{w}) = \left[\frac{\partial J}{\partial w_0}, \frac{\partial J}{\partial w_1}, \dots, \frac{\partial J}{\partial w_n}\right]$  • Here is how we change each weight:

For 
$$i = 0$$
 to  $n$ :  
 $w_i \leftarrow w_i + Dw_i$ 

where

$$\mathsf{D}w_i = -h \frac{\partial J}{\partial w_i}$$

# "True" (or "batch") gradient descent

- One *epoch* = one iteration through the training data.
- After each epoch, compute average loss over the training set:

$$J(\mathbf{w}) = \frac{1}{M} \mathop{\bigotimes}\limits_{k=1}^{M} L(\mathbf{w}, \mathbf{x}^{k}, t^{k}) = \frac{1}{M} \mathop{\bigotimes}\limits_{k=1}^{M} \frac{1}{2} \left(t^{k} - y\right)^{2}$$

 Change the weights to move in direction of steepest descent in the average-loss surface:



From T. M. Mitchell, *Machine Learning* 

• Problem with *true gradient descent*:

Training process is slow.

Training process can land in local optimum.

- Common approach to this: use *stochastic gradient descent*:
  - Instead of doing weight update after all training examples have been processed, do weight update after each training example has been processed (i.e., perceptron output has been calculated).
  - Stochastic gradient descent approximates true gradient descent increasingly well as  $\eta \rightarrow 1/\infty$ .

#### Derivation of perceptron learning rule (stochastic gradient descent)

We defined 
$$J = \frac{1}{2} (t - y)^2$$

Here, use

$$J = \frac{1}{2} \left( t - (\mathbf{w} \times \mathbf{x}) \right)^2$$

$$= \frac{1}{2} \left( t - (\mathbf{w}_1 x_1 + \mathbf{w}_2 x_2 + \mathbf{v}_n x_n) \right)^2$$

 $y = a(\mathbf{w} \times \mathbf{x}^k)$ 1 0 0

$$y = \mathbf{W} \times \mathbf{x}^k$$

Then,

 $\frac{\P J}{\P w_i} = -\left(t - (\mathbf{W} \times \mathbf{X})\right) x_i^k$ 

#### But we'll use

 $\frac{\P J}{\P w_i} = -(t - y) x_i^k$ 

$$\Delta \mathbf{w}_i = -\boldsymbol{\eta} \frac{\partial J}{\partial w_i} = \boldsymbol{\eta} (t^k - y^k) x_i^k$$

This is called the "perceptron learning rule"

#### **Perceptron Learning Algorithm**

Start with small random weights,  $\mathbf{w} = (w_0, w_1, w_2, \dots, w_n)$ , where  $w_i \hat{1} [-.05, .05]$ 

Repeat until accuracy on the training data stops increasing or for a maximum number of *epochs* (iterations through training set):

For k = 1 to M (total number in training set):

- 1. Select next training example  $(\mathbf{x}^k, t^k)$ .
- 2. Run the perceptron with input  $\mathbf{x}^k$  and weights  $\mathbf{w}$  to obtain y.
- 3. If y ≠ t<sup>k</sup>, update weights:
  for i = 0, ..., n: ; Note that bias weight w<sub>0</sub> is changed just like all other weights!
  - $w_i \leftarrow w_i + \Delta w_i$

where

 $\Delta w_i = \eta (t^k - y^k) x_i^k \quad (or \, equivalently : \Delta w_i = -\eta (y^k - t^k) x_i^k)$ 

1. Go to 1.

#### Example: logical OR

Training set: ((0, 0), 0) (0, 1), 1) (1, 0), 1) (1,1),1))





Initial weights:  $\{w_0, w_1, w_2\} = \{-.05, -0.2, 0.2\}$   $w_{i} \leftarrow w_{i} + \Delta w_{i}$ where  $\Delta w_{i} = -\eta (y^{k} - t^{k}) x_{i}$ 

Apply perceptron learning rule for one epoch with  $\eta = 0.25$ 

(\*) Please note that the text uses "-1" for the extra input weight that corresponds with the bias  $(w_0)$ . It is slightly more conventional to use "+1" so please note we will also use this convention from time to time. (Either value generates a model of comparable expressive power)

#### Example: logical OR

Training set: ((0, 0), 0) (0, 1), 1) (1, 0), 1) (1,1),1))





Initial weights:

 $\{w_0, w_1, w_2\} = \{-.05, -0.2, 0.2\}$ 

Input: (0,0)Input: (0,1) $w_0: -0.05 - 0.25 \times (1-0) \times -1 = 0.2$  $w_0: 0.2 - 0.25 \times (0-1) \times -1 = -.05$  $w_1: -0.02 - 0.25 \times (1-0) \times 0 = -0.02$  $w_1: -0.02 - 0.25 \times (0-1) \times 0 = -0.02$  $w_2: 0.02 - 0.25 \times (1-0) \times 0 = 0.02$  $w_2: 0.02 - 0.25 \times (0-1) \times 1 = 0.27$ 

 $w_i \leftarrow w_i + \Delta w_i$ where  $\Delta w_i = -\eta (y^k - t^k) x_i$ 

No weight updates are needed for inputs: (1,0) and (1,1); why?

Does the Perceptron Learning Algorithm (PLA) always work?

Q: Will the PLA yield a solution to the logical XOR problem?

Why/why not?



Does the Perceptron Learning Algorithm (PLA) always work?

Q: Will the PLA yield a solution to the logical XOR problem?

Why/why not?

The PLA will not yield a solution, since the data is not linearly separable.



Does the Perceptron Learning Algorithm (PLA) always work?

Q: Will the PLA yield a solution to the logical XOR problem?

Why/why not?

The PLA will <u>not yield a solution</u>, since the data is <u>not linearly separable</u>.

Q: Does this mean the classification problem is hopeless?



Q: Does this mean the classification problem is hopeless? No! There are at least (3) remedies:

(1) Project the XOR problem into a higher dimensional space (e.g. 3-space) where it is linearly separable!

- (2) Use a kernel/polynomial decision boundary in 2-space.
- (3) Use a multi-layer perceptron (i.e. a neural network).

(\*) Minsky and Papert (1969) published an influential (viz. notorious) <sup>A</sup> text, "Perceptrons", that identified the learning capabilities and limitations of Perceptrons.

(\*) The major effect of this text, unfortunately, was to <u>set back NN research for</u> <u>two decades</u> (see: "AI Winter")

(\*) What brought it back? LeCun, et al. 1990s MNIST results, etc.



B

В

## Recognizing Handwritten

28 pixels

- MNIST dataset Digits
  - 60,000 training examples
  - 10,000 test examples

Each example is a 28×28pixel image, where each pixel is a grayscale value in [0,255].

See csv files.

First value in each row is the target class.

0123456789 0123456799 0123456789 0/23456789 0123456789



Label: "2"

28 pixels
#### Perceptron architecture for handwritten digits classification



#### Preprocessing (To keep weights from growing very large)

# Scale each feature to a fraction between 0 and 1:

$$x_i' = \frac{x_i}{255}$$

#### **Processing an input**

1 **'**0' **'**1' **'**2' **'**3' **'**4' **'**5' **'**6' **'**7' **'**8' **'**9'

 $x_0$ 

 $x_1$ 

 $x_{n-1}$ 

 $x_{n}$ 

At each output, compute:



 $\overset{n}{a} w_i x_i$ i=0The output with highest value is the prediction.

If correct, do nothing. If not correct, adjust weights according to perceptron learning rule.

See assignment for details.

## Example

For each training example k, input  $\mathbf{x}^k$  to the perceptron.

 $x_0$ 

 $x_1$ 

 $x_{n-1}$ 

1)

Suppose  $\mathbf{x}^k$  is a '2'.

 $t^k = 1$  for the '2' output.

 $x_{n}$  (  $t^{k} = 0$  for all other outputs.

$\bigcirc$	<b>'</b> 0'
$\bigcirc$	<b>'</b> 1'
$\bigcirc$	<b>'</b> 2'
$\bigcirc$	'3'
$\bigcirc$	'4'
$\bigcirc$	<b>'</b> 5'
$\bigcirc$	<b>'</b> 6'
$\bigcirc$	'7'
	<b>'</b> 8'
$\overline{\bigcirc}$	<b>'</b> 9'

For each output, calculate



Set y = 1 for outputs that are greater than 0. Set y = 0otherwise.

Then for all weights coming into each output:

 $w_i \leftarrow w_i + \Delta w_i$ where  $\Delta w_i = \eta (t^k - y^k) x_i^k$ 

## Homework 1

- Implement perceptron (785 inputs, 10 outputs) and perceptron learning algorithm in any programming language you would like.
- Download MNIST data from website
- Preprocess MNIST data: Scale each feature to a fraction between 0 and 1.

## Homework 1 Experiments

Train perceptrons with three different learning rates:

 $\eta = 0.001, 0.01, and 0.1$ 

For each learning rate:

1. Choose small random weights  $w_i \hat{1} [-.05,.05]$ 

2. Repeat cycling through the training data until the accuracy on the training data has essentially stopped improving (i.e., the difference between training accuracy from one epoch to the next is less than some small number, like 0.01.)

After the initialization, and after each epoch (one cycle through training data), compute accuracy on training and test set (for plot), without changing weights.

### Homework 1: Presenting results For each experiment, plot training and test accuracy over epochs:



Accuracy (%)

Epoch

# For each experiment, give confusion matrix:

Predicted class

Actual class

	0	1	2	3	4	5	6	7	8	9
0		6							1	
1	X						1			
2	1	/	24							
3	7	\ \	/			$\langle \rangle$	11	1		
4				$\times$		1				
5	/			1.		<u></u>	/	+		
6	X.	$\swarrow$			X		/			
7			1					¥		~
8	$\geq$		1							
9			10		0					

In each cell (*i*,*j*). put number of test examples that were predicted (classified) as class *i*, and whose actual class is class *j*.

#### Homework 1 FAQ

Q: Can I use code for perceptrons from another source?A: No, you need to write your own code.

Q: How long will it take to train the perceptrons?A: Depends on your computer and your code, but probably an hour or more.

Q: What accuracy should I expect on the test set?A: It depends on initial weights, but probably over 80%.

Q: Should I wait until the last minute?A: No!!! Start as soon as possible.

(\*) Rosenblatt (1962) proved that:

Given a linearly separable dataset, the Perceptron will converge to a solution that separates the classes, and that it will do it <u>after a finite number of</u> <u>iterations</u>.

Pf. Let:  $\gamma$  := the distance between separating hyperplane and Closest data point; let w\* be a unit weight vector that separates the data (known to exist by assumption of linear separability).

WLOG, assume  $||x|| \le 1$  for all input data.

Recall to check "similarity" of two vectors (in particular, "perfect alignment" connotes parallel).



Recall: to check "similarity" of two vectors we take their dot product (in particular, "perfect alignment" connotes parallel); when two vectors are parallel, their inner product is maximal.

If we therefore show that at each weight update,  $w^* \cdot w$ Increases, then we have nearly show that the algorithm will converge.

However, we do need to check that the length of w does not increase too much as well, since  $\vec{w}^* \cdot \vec{w} = \|\vec{w}^*\| \|\vec{w}\| \cos\theta$ .

In summary, we need (2) checks: (1)  $\vec{w}^* \cdot \vec{w}$ (2)  $\|\vec{w}\|$ 

Suppose that at the *t*th iteration of the algorithm, the network sees a particular x that should output y, and that it gets this input wrong, so:

 $y\vec{w}^{(t-1)}\cdot\vec{x}<0$ 

(\*) This means that the weights need to be updated. (\*) The weight update will be:  $\vec{w}^{(t)} = \vec{w}^{(t-1)} + y\vec{x}$  (where  $\eta = 1$ , WLOG)

Consider:  $\vec{w} * \cdot \vec{w}^{(t)} = \vec{w} * \cdot \left( \vec{w}^{(t-1)} + y\vec{x} \right)$ 

Suppose that at the *t*th iteration of the algorithm, the network sees a particular x that should output y, and that it gets this input wrong, so:

 $y\vec{w}^{(t-1)}\cdot\vec{x}<0$ 

(\*) This means that the weights need to be updated. (\*) The weight update will be:  $\vec{w}^{(t)} = \vec{w}^{(t-1)} + y\vec{x}$  (where  $\eta = 1$ , WLOG)

Consider:  $\vec{w}^* \cdot \vec{w}^{(t)} = \vec{w}^* \cdot \left(\vec{w}^{(t-1)} + y\vec{x}\right)$ =  $\vec{w}^* \cdot \vec{w}^{(t-1)} + y\vec{w}^* \cdot \vec{x}$  Why?

Suppose that at the *t*th iteration of the algorithm, the network sees a particular x that should output y, and that it gets this input wrong, so:

 $y\vec{w}^{(t-1)}\cdot\vec{x}<0$ 

(\*) This means that the weights need to be updated. (\*) The weight update will be:  $\vec{w}^{(t)} = \vec{w}^{(t-1)} + y\vec{x}$  (where  $\eta = 1$ , WLOG)

Consider:  $\vec{w} * \cdot \vec{w}^{(t)} = \vec{w} * \cdot \left(\vec{w}^{(t-1)} + y\vec{x}\right)$  $= \vec{w} * \cdot \vec{w}^{(t-1)} + y\vec{w} * \cdot \vec{x}$   $\geq \vec{w} * \cdot \vec{w}^{(t-1)} + \gamma$ Why?

### Perceptron Convergence Theorem $\vec{w}^* \cdot \vec{w}^{(t)} = \vec{w}^* \cdot \left(\vec{w}^{(t-1)} + y\vec{x}\right)$ $= \vec{w}^* \cdot \vec{w}^{(t-1)} + y\vec{w}^* \cdot \vec{x}$ $\geq \vec{w}^* \cdot \vec{w}^{(t-1)} + \gamma$

This means that at each update of the weights, this inner product increases by at least  $\gamma$ , and so after t updates of the weights:

$$\vec{w}^* \cdot \vec{w}^{(t)} \ge t\gamma$$

(\*) We can use this to put a bound on the length of ||w(t)|| using the Cauchy-Schwartz inequality, which tells us:

$$\vec{w}^* \cdot \vec{w}^{(t)} \leq \left\| \vec{w}^* \right\| \left\| \vec{w}^{(t)} \right\| \text{ and so, } \left\| \vec{w}^{(t)} \right\| \geq t\gamma$$

Perceptron Convergence Theorem $\left\|\vec{w}^{(t)}\right\| \ge t\gamma$ 

The length of the weight vector after *t* steps is:

$$\begin{aligned} \left\| \vec{w}^{(t)} \right\|^{2} &= \left\| \vec{w}^{(t-1)} + y\vec{x} \right\|^{2} \\ &= \left\| \vec{w}^{(t-1)} \right\|^{2} + y^{2} \left\| \vec{x} \right\|^{2} + 2y\vec{w}^{(t-1)} \cdot \vec{x} \end{aligned}$$
 Why? 
$$\leq \left\| \vec{w}^{(t-1)} \right\|^{2} + 1$$

(\*) The last inequality holds: since  $y^2=1$ ,  $\|\vec{x}\| \le 1$ , and the network made an error so  $\vec{w}^{(t-1)}$  and  $\vec{x}$  yield a negative value (when we assume y=1, WLOG).

(\*) The inequality above shows that:  $\|\vec{w}^{(t)}\|^2 \leq t$ 

(\*) If we put these two inequalities together, we get:

 $t\gamma \le \left\|\vec{w}^{(t-1)}\right\| \le \sqrt{t}$ 

$$t\gamma \le \left\|\vec{w}^{(t-1)}\right\| \le \sqrt{t}$$

(\*) So  $t \leq 1/\gamma^2$ 

(\*) Thus after this many updates the algorithm must have converged. QED (\*) We showed that the quantity:  $\vec{w}^* \cdot \vec{w} = \|\vec{w}^*\|\|\vec{w}\|\cos\theta$  increases monotonically and it is upper-bounded.

Summary: We demonstrated that if the data are linearly separable, then the algorithm will converge and that the time this takes is a function of the distance between the separating hyperplane and the nearest data point (this distance is called the margin).

If the data are not linearly separable, then the algorithm is not guaranteed to converge and may oscillate infinitely.

## Neuron inspires Regression

- Edges multiply the signal  $(x_i)$  by some weight  $(\theta_i)$ .
- Nodes sum inputs
- Equivalent to Linear Regression

$$f(x,\theta) = \sum_{d=1}^{D} \theta_d x^d + \theta_0$$



