

9/17/90

①

Outline

- linear threshold func.
- training using lin. prog.
- " " " Perceptron alg.
- Perceptron Conv. alg.
- Lower bounds
- Discussion

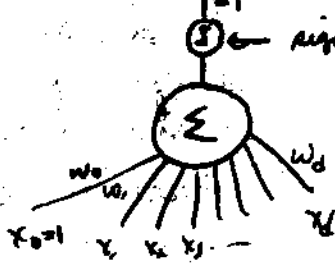
- input  $\leftrightarrow (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$

- w.l.o.g.  $x_0 = 1$  (could be an extra value that is added, etc.)

- learning from examples (positive & negative - i.e. labelled examples)

prior knowledge: examples can be classified by a linear threshold function

LTF:



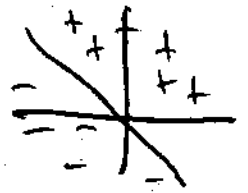
$w = (w_0, \dots, w_d)$

$x = (x_0, \dots, x_d)$

$w \cdot x = 0$  - hyperplane in  $d$  dimension

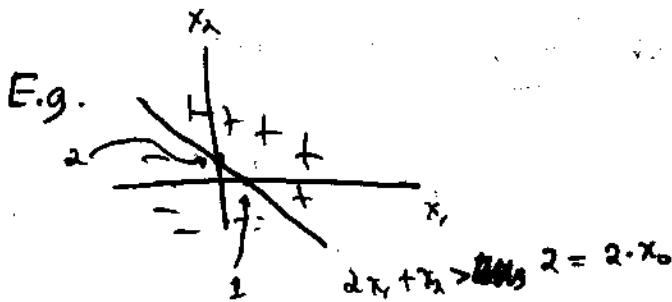
↑  
graphically

↑  
algebraically



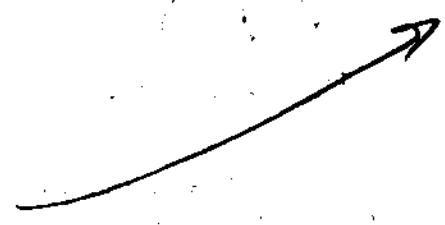
$w \cdot x > 0$  positive example

$w \cdot x < 0$  negative exam.



$x_0 = 1$

or  $(2, 1, 2) \cdot \vec{x} > 0$   
 $(-2, 2, 1) \cdot \vec{x} > 0$

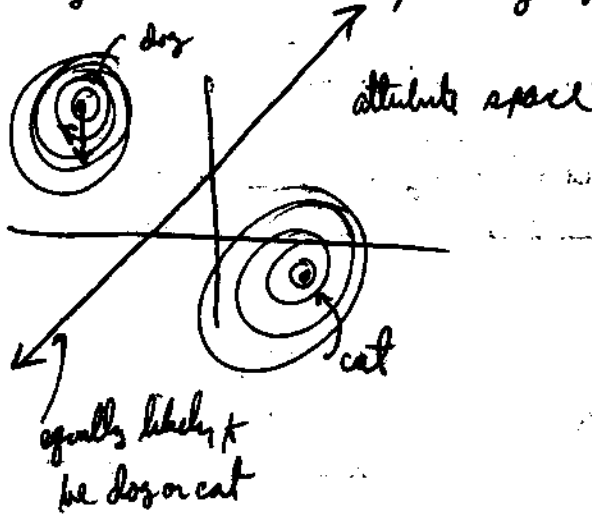


$X_+$  - positive training examples  
 $X_-$  - negative training examples

Def  $(X_+, X_-)$  are linearly separable if  $\exists w: w \cdot x > 0 \forall x \in X_+$   
 $w \cdot x < 0 \forall x \in X_-$

Why is this stuff interesting?

- 1) neuron is something like this ... (abstract model of a neuron)
- 2) ??
- 3) they are sometimes provably optimal

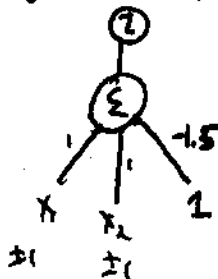


- examples are given by multivariate distribution from prototypical dog or cat - i.e. likelihood of dog (or cat) is  $e^{-r^2/2}$  where  $r$  is distance from prototypical dog (or cat).

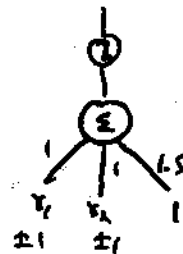
↳ this linear separator however will as you can do.

4) can represent great many things: AND, NOT, OR gates, etc.

AND



OR



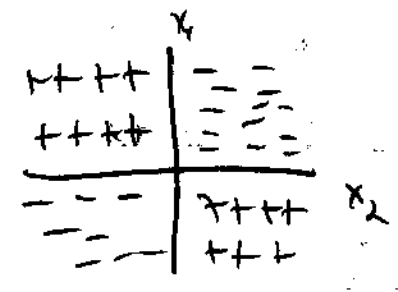
NOT



Note: can't do XOR or parity (can easily do majority)

XOR

	$x_1$	$x_2$
	+	-
$x_1$	-	+
-	+	-



can draw linear separator

How many boolean fun. on  $d$  variables (inputs) are threshold fun.?

$\rightarrow 2^{\Theta(d^2)}$

How many total?  $\rightarrow 2^{2^d}$

note  $\frac{2^{\Theta(d^2)}}{2^{2^d}} \rightarrow 0 \quad d \rightarrow \infty$  (as ratio many asymptotically)

Learning Algorithm

① Batch Alg. -  $x_+, x_-$  lin. sep.

- find  $w$  s.t.  $w \cdot x = 0$  separates  $x_+$  from  $x_-$

now, w.l.o.g. assume  $x_- = \emptyset$  (add  $-x$  to  $x_+$   $\forall x \in x_-$ )

now, find  $w \cdot x > 0 \quad \forall x \in x_+$

- this is just linear programming

$\therefore$  can solve very quickly (P algorithm) (only which can be done in poly time)

if  $w \cdot \tilde{x} < 0 \quad \forall \tilde{x} \in x_-$   
 then  
 $w \cdot (-\tilde{x}) > 0 \quad \forall \tilde{x} \in x_-$   
 $\Rightarrow$  add  $-\tilde{x}$  to  $x_+$   
 $\forall \tilde{x} \in x_-$

② Incremental (on-line) alg.

Perceptron alg. (Rosenblatt) (assumes again that all examples are in  $X_+$ )

$w=0$   
 repeat  
 for each training example  
 if  $w \cdot x > 0$  do nothing  
 if  $w \cdot x \leq 0$   
 $w \leftarrow w + x$   
 until  $w \cdot x > 0 \forall x \in X_+$

Why does this work?

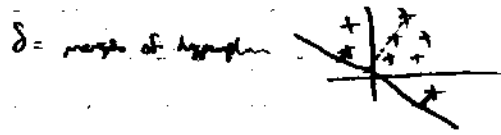
$w' \leftarrow w + x$

so,  $w' \cdot x = w \cdot x + \underbrace{x \cdot x}_{> 0}$

$L$  will be dec to 0 or positive

Correctness proof

Thm: Assume  $\exists v \in V, v \cdot x > \delta > 0 \forall x \in X_+$   
 $\exists M \in \mathbb{R} \forall x \in X_+ |x|^2 \leq M$



then Perceptron Alg. converges after finite # of updates.

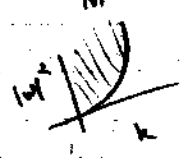
Pf: ① After  $k$  updates,  $w = x^{(1)} + \dots + x^{(k)}$   $x^{(i)} = i^{th}$  misclassified example

$v \cdot w = v \cdot x^{(1)} + \dots + v \cdot x^{(k)}$   
 $> k \cdot \delta$

$(\sum_{i=1}^k x^{(i)})^2 = (\sum_{i=1}^k |x^{(i)}|)^2$   
 $(v \cdot w)^2 \leq |v|^2 |w|^2$

Cauchy-Schwarz

$\Rightarrow |w|^2 \geq \frac{(v \cdot w)^2}{|v|^2} \geq \frac{(k\delta)^2}{|v|^2} = k^2 \frac{\delta^2}{|v|^2}$  - growing quadratically w/ updates



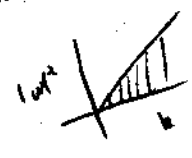
$\frac{v \cdot w}{|v||w|} = \cos \theta$   
 $\frac{(v \cdot w)^2}{|v|^2 |w|^2} = \cos^2 \theta \leq 1$

②  $w \cdot w = |w|^2 = |w|^2 + 2w \cdot x + |x|^2$   
 $\text{each update } \underbrace{2w \cdot x}_{< 0} + |x|^2$   
 $\Rightarrow |w|^2 = |w|^2 \leq |x|^2$

$\Rightarrow$  after  $k$  updates,  $|w|^2 \leq k \cdot M$

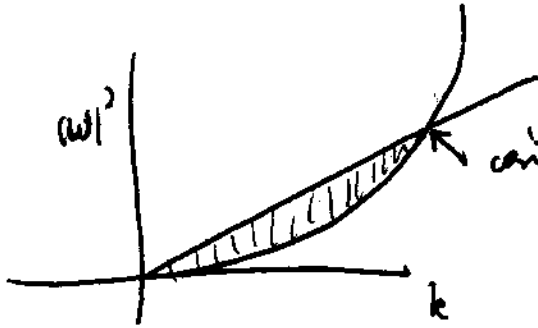
- growing at most linear

When does  $k^2 \frac{\delta^2}{|v|^2} > kM$ ?  $\Rightarrow k > \frac{M|v|^2}{\delta^2} \neq \text{so, } k < \frac{M|v|^2}{\delta^2}$



3

So,



$$kM \geq |w|^2 \geq \frac{k^2 \delta^4}{|w|^2}$$

$$kM = (k\delta)^4 / |w|^2$$

$$\Rightarrow k \leq \frac{M \cdot |w|^2}{\delta^2}$$

In fact this alg is exponential in  $d$  in worst case:

Then (Monoga)

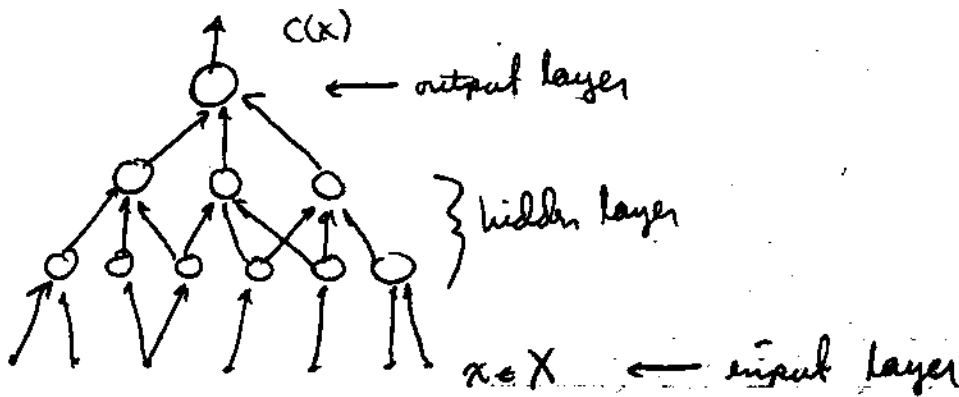
$$(-((x_1 \wedge x_2) \vee x_3) \wedge x_4) \vee x_5) \wedge \dots x_d$$

Claim:  $\uparrow$  is a threshold function, but weights grow exponentially in  $d$ .  
(turn out to be Fibonacci #'s)

$\Rightarrow$  Perceptron alg. can take exp. time.

11/4/92

Neural nets



neuron:



$$w \cdot x \geq \theta ?$$



sharp threshold

or



sigmoid

Network w/ sigmoid is a continuous function of inputs & weights.

Train using gradient descent... (back propagation)

Error we'd like to minimize: M.S.E.  $E = \frac{1}{n} \sum_{x \in S} (f(x) - c(x))^2$

$\nabla E = \left( \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_k} \right)$  (change in error as we tweak weights)

$w_0 \in \mathbb{R}^k$  initial weights (may be random)

$$\vec{w}_{i+1} = \vec{w}_i - \underset{\substack{\uparrow \\ \text{learning rate}}}{\epsilon} \nabla E$$

sometimes we momentum term...

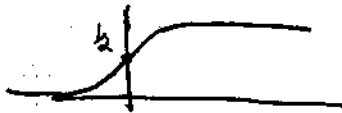
~~$\vec{w}_{i+1} = \vec{w}_i + \epsilon \nabla E$~~

now  $\Delta w_{i+1} = (\alpha \Delta w_i + (1-\alpha) \epsilon (-\nabla E))$

$\&$   $w_{i+1} = w_i + \Delta w_{i+1}$


---

types of sigmoids:

1) logistic function  $z = \frac{1}{1+e^{-x}}$  

$$\frac{dz}{dx} = z(1-z) = \frac{e^{-x}}{(1+e^{-x})^2}$$

2) hyperbolic tangent

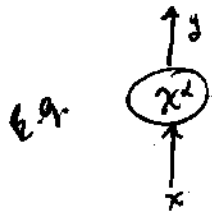
$$z = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$
 

$$\frac{dz}{dx} = \text{sech}^2(x) = \frac{4}{(e^x + e^{-x})^2}$$

Chain rule

(back-propagation)

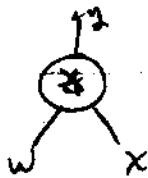
Let  $\delta_x = \frac{\partial E}{\partial x}$  (notation)



$\delta_x = \delta_y \cdot 2x$  (e.g.  $\frac{\partial E}{\partial x} = \frac{\partial E}{\partial y} \cdot \frac{\partial y}{\partial x}$ )



$\delta_x = \delta_y \cdot y(1-y)$  (logistic)

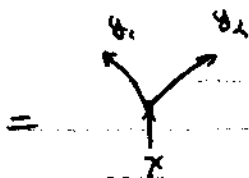


$\delta_w = \delta_y \cdot x$

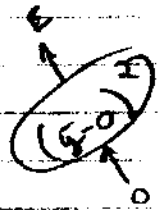
$\delta_x = \delta_y \cdot w$



$\delta_{x1} = \delta_y$   
 $\delta_{x2} = -\delta_y$



$\delta_x = \delta_{y1} + \delta_{y2}$



$\frac{\partial E}{\partial 0} = -\delta_y(-1)$   
 $= \delta_y$

$\delta_{x1} = \delta_{y1}$   
 $\delta_{x2} = \delta_{y2}$

- add to get cumulative effect

Net Talk (text  $\Rightarrow$  speech)

actually, text  $\Rightarrow$  phonemes

- all paper



11/9/92

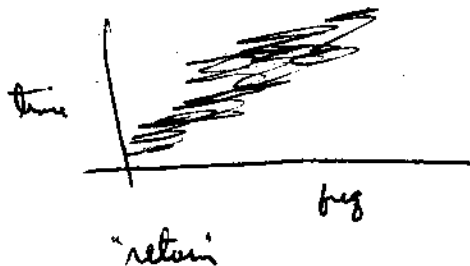
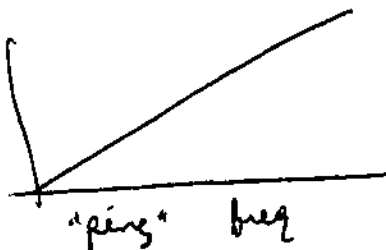
## Net-Talk

- Resistance to damage
  - added noise to weights
  - turns out that network performed well until fewer than 4 bits of weight left
  - $\therefore \approx 4$  bits info/weight suffice
- Vary # hidden units
  - hidden units (perception) - 82% accuracy
  - 120 " " " - 99% accuracy
- Vary # layers - similar performance w/ two layers as w/ one.
- Tea leaf reading
  - can you see things<sup>patterns, etc.</sup> in the trained network?
  - generally, no. Info seems to be dist. over nodes.

## Sonar example

(Neural Networks 1 (1988) 75-89  
Gormant Sejnowski)

- try to tell rocks from mines from sonar "pings".



9/26/90

Well done with training a single neuron, what about a collection?

Need new techniques & criteria...

### LMS criterion

target function:  $t: \mathbb{R}^n \rightarrow \mathbb{R}$

$(x, t(x))$  typical training example

$f$  = output of learner

label

prediction error on  $x = \begin{cases} 0 & \text{if } f(x) = t(x) \\ 1 & \text{if } f(x) \neq t(x) \end{cases}$

now

square error on  $x = (f(x) - t(x))^2$

mean squared error:

$$MSE = \frac{1}{M} \sum_{x \in S} (f(x) - t(x))^2$$

(try to minimize this)

(note: if  $f, t \in \{0, 1\}$  values this is identical to the above).

## Gradient Descent

$$f(x) = f_w(x) \quad w \in \mathbb{R}^k$$

↑  
weights, parameters, knobs, etc.

finding good  $f \equiv$  finding good  $w$

$$E(w) = \text{MSE} = \frac{1}{n} \sum_{x \in S} (f_w(x) - t(x))^2$$

assume  $E, f$  are differentiable w.r.t.  $w$ .

Idea: start somewhere, & run downhill until we get to a local minimum —

$$\nabla E = \left( \frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_n} \right)$$

if we go in direction of  $\nabla E$ , go uphill

so, go in direction of  $-\nabla E$



alg

$\vec{w}_0 =$  initial pt. =  $(0, 0, \dots, 0)$  (perhaps) (or maybe random)

$$\vec{w}_{i+1} \leftarrow \vec{w}_i + \underset{\substack{\uparrow \\ \text{learning rate}}}{\epsilon} \Delta \vec{w}_i \quad \Delta \vec{w}_i = (-\nabla E)$$

## Remarks

$$\textcircled{1} \quad \nabla E(S) = \frac{1}{m} \sum_{x \in S} \nabla E(S; x)$$

↑  
gradient  
of cost  
over entire set

↑  
sum of gradients of singletons

[updates are  
everywhere at a time]

there, can either look at all points & then take a step (update)  
or, look at each point & then take step.

## momentum



want to avoid bouncing back & forth across valley

let our move this time be an average of our proposed move  
this time & our move last time (will avoid weird bouncing back  
& forth)

$$\Delta \tilde{w}_i = \alpha \Delta \tilde{w}_{i-1} + (1-\alpha)(-\nabla E)$$

## Gradient Descent for Linear Function

$$f_w(x) = w \cdot \vec{x}$$

$$S = \{x_i\} \quad t(x_i) = y_i$$

$$E(w) = \frac{1}{m} \sum_i (w \cdot \vec{x}_i - y_i)^2$$

$$\nabla E(w) = \frac{2}{m} \sum_i (w \cdot \vec{x}_i - y_i) \cdot \vec{x}_i$$

in this case,

local minimum  $\Rightarrow$  global minimum

background / names:

Wrochow/Hoff, Adaline (1960's)

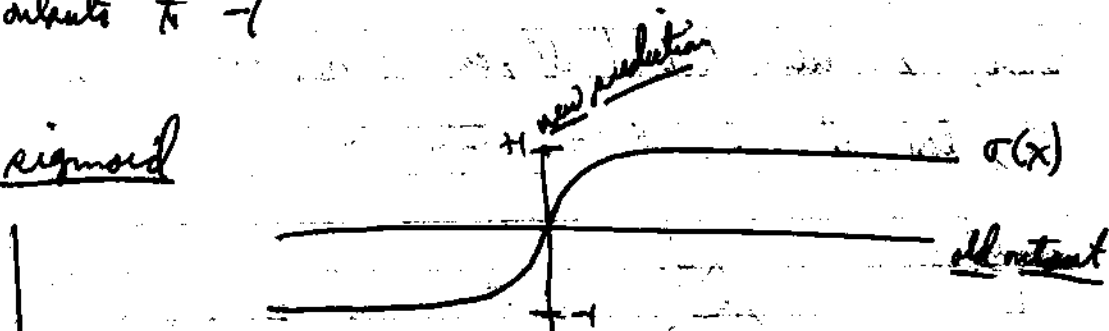
Now, think about threshold functions.

- a data point which is positive (if this we should output 1) may actually output 999 - MSE will be high. (if function will be dropped that way to compensate)

like to transform all positive ~~outputs~~ to 1 & all negative

outputs to -1

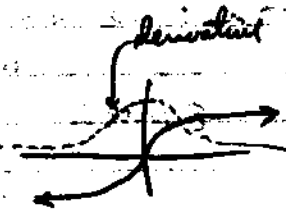
- use sigmoid



↳ differentiable approximation to a threshold

①  $z = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

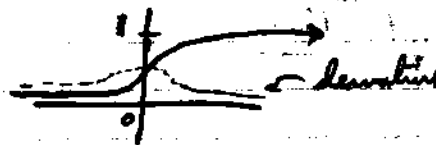
hyperbolic tangent



$$\frac{dz}{dx} = \text{sech}^2(x) = \frac{4}{(e^x + e^{-x})^2}$$

② logistic function

$$z = \frac{1}{1 + e^{-x}}$$



$$\frac{dz}{dx} = z(1-z)$$

note: looking at threshold function, so  $\in \{+1, -1\}$

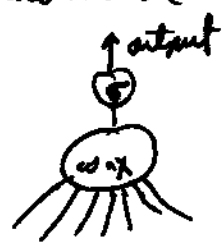
Now, let  $f_w(x) = \sigma(w \cdot x)$

$$E(w) \{x_i\} = (\sigma(w \cdot \vec{x}_i) - y_i)^2$$

$$\nabla E(w) \{x_i\} = \underbrace{2(\sigma(w \cdot \vec{x}_i) - y_i)}_{\text{error}} \cdot \text{sech}^2(w \cdot \vec{x}_i) \cdot \vec{x}_i$$

weight  $\rightarrow$  weight stuff  
closer to  $\pm 1$  more heavily  
don't weight others so much

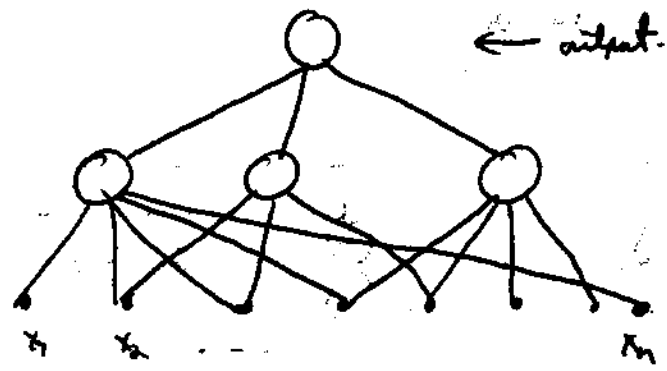
neuronid, pseudo-neuron



will represent  $\Rightarrow$



### "Neural Networks"



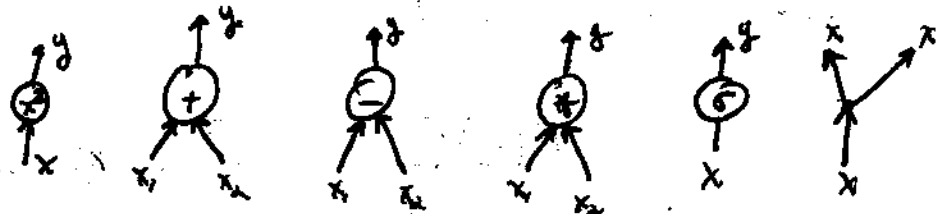
hidden layer

inputs

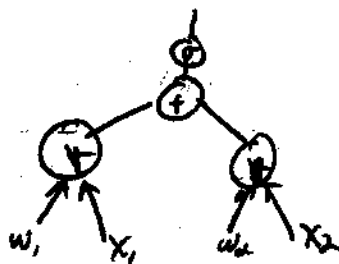
Now, gradient descent method still applies... back-prop

# Computing $\nabla E$

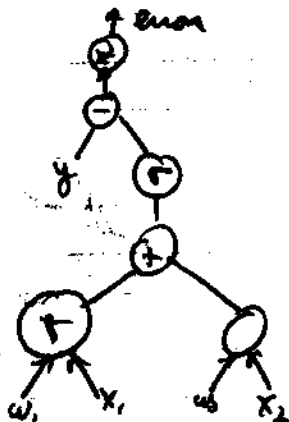
some primitive elements:



e.g., 2 input neuron



circuit errors:



idea: compute forward to compute  $y, \delta_y$  & backward to get  $x, \delta_x$

e.g.



$$\delta_x = \delta_y \cdot \Delta x$$



$$\delta_y = \delta x_1$$

$$= \delta x_2$$



$$\delta x_1 = \delta y$$

$$\delta x_2 = \delta y$$



$$\delta_{x_1} = \delta_y \cdot x_2$$

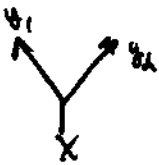
$$\delta_{x_2} = \delta_y \cdot x_1$$



$$\delta_x = \delta_y \cdot (y)(1-y)$$

$$\text{or}$$

$$\delta_x = \delta_y \cdot \text{sech}^2(x)$$

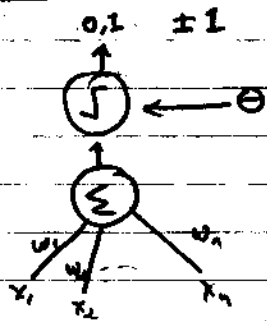


$$\delta_x = \delta_{y_1} + \delta_{y_2}$$

Now, can ~~we~~ back propagate & get  $\delta w_i$  &  $\delta c_i$   
& this is gradient of loss w.r.t weights.

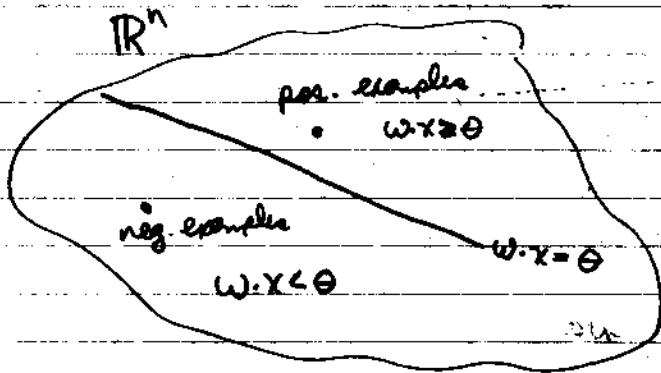


## Perceptrons

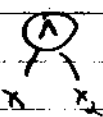


$$\sum w_i x_i \geq \theta \Rightarrow \text{output} = +1$$

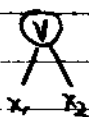
$$\text{else} \Rightarrow \text{output} = -1$$



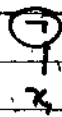
## Boolean Functions



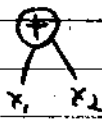
$$x_1 + x_2 \geq 2$$



$$x_1 + x_2 \geq 1$$



$$-x_1 \geq 0$$



XOR  
not linear

Thm: # boolean threshold functions is  $2^{\Theta(n)}$

(note: total # boolean functions is  $2^{2^n}$ )

Linear separator work well:

Consider two types of examples which are generated by multivariate gaussian distribution:



optimal classifier is a hyperplane  
(perpendicular bisector of segment joining the centers)

### Perceptron training algorithm

Assumptions:

- assume no examples on hyperplane,  $w \cdot x \neq 0$  for  $\forall w \in \text{any } x \text{ we see}$

- assume,  $\Theta = 0$  (w.l.o.g.)

Pr: add new coordinate,  $x_i$ , to all examples with  $x_i = 1$

$$-w_i = \Theta$$

- assume all examples are positive examples (w.l.o.g.)

(negate  $x_i$  on all negative examples)

## Perceptron Alg.

- Initially:  $\vec{w} = (0, 0, \dots, 0)$
- Given an infinite sequence  $x^1, x^2, \dots$  of examples:
  - predict using  $\text{sign}(w \cdot x^i)$
  - set correct label for  $x^i$

if  $w \cdot x^i > 0$  do nothing  
else update  $w$  by  $w \leftarrow w + x^i$

(note:  $w' = w + x$ )

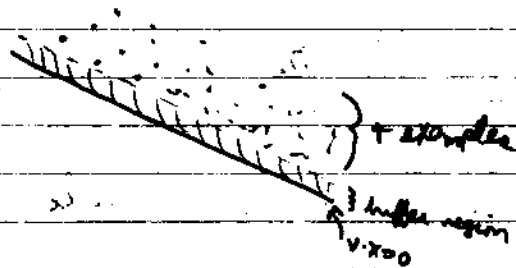
$$w' \cdot x = wx + |x|^2$$

$\hookrightarrow \infty$

(no more likely to correctly classify)

- Assume  $|x^i| = 1$  ( $\forall i$ ) (this can be weakened, but some upper bound is needed)

- Assume  $(\exists \delta)(\exists v \in \mathbb{R}^n) v \cdot x^i > \delta$  ( $\forall i$ )



## Thm: Perceptron Convergence Theorem

Given assumptions, perceptron procedure only makes a finite # mistakes

pt: assume  $(w, \xi)$   $|w| = 1$  (size, scale & change  $\delta$ )

$$f(w) = \frac{v \cdot w}{|w|} \quad (\cong \cos \alpha \text{ where } \begin{array}{c} \nearrow v \\ \searrow w \end{array})$$
$$\leq 1$$

How does  $f(w)$  change when  $w' \leftarrow w + \xi$ ?

numerator:  $v \cdot w' = v \cdot (w + \xi) = v \cdot w + v \cdot \xi \geq v \cdot w + \delta$

(increases by at least  $\delta$  each update)

denominator:  $|w'|^2 = w' \cdot w' = (w + \xi) \cdot (w + \xi)$

$$= |w|^2 + 2w \cdot \xi + |\xi|^2$$

$\uparrow$   $\hookrightarrow = 1$   
 $< 0$

because we made a mistake

$$\leq |w|^2 + 1$$

after  $t$  updates:  $f(w) \geq \frac{t \cdot \delta}{\sqrt{t}} = \delta \sqrt{t}$

note: yield contradiction if  $t > 1/\delta^2$

$\Rightarrow$

At most  $1/\delta^2$  mistakes

## Variations

$$w \leftarrow w - \lambda(w \cdot x)x$$

$$0 < \lambda < 2$$

(note: if  $\lambda=1$ , moves "just enough" to correctly classify  $x$ ;  
procedure converges for any  $0 < \lambda < 2$ )

Note: Boolean threshold functions can have  $1/8 \approx 2^{-3}$

e.g.  $x_1 \vee (x_2 \wedge (x_3 \vee (x_4 \wedge \dots)))$

Note:  $2^{\Theta(n^2)}$  threshold functions (boolean)

Halving alg. yield only  $\Theta(n)$  mistakes (possibly much better than above).

- Can we bridge this gap?

## VC-dimension

Let  $X$  be some domain

Let  $\mathcal{C}$  be some concept class on  $X$

Define The Vapnik-Chervonenkis dimension ( $VC\text{-dim}(\mathcal{C})$ ) is the

largest  $d$  s.t.  $\exists d$  distinct pts  $x^1, x^2, \dots, x^d \in X$

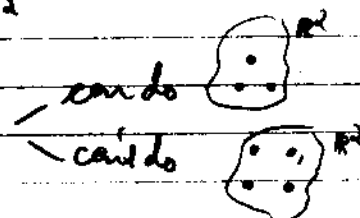
s.t. for any  $d$ -labelling of these pts,  $\exists$  a concept  $c \in \mathcal{C}$   
that realizes that labelling.

• If  $|C| < \infty$ , then  $VC\text{-dim}(C) \leq \lg |C|$

•  $C = \text{all concepts}$ ,  $|X| = \infty \Rightarrow VC\text{-dim} = \infty$

•  $C = \text{half-spaces in } \mathbb{R}^2$

$VC\text{-dim} = 3$



•  $C = \text{half-spaces in } \mathbb{R}^n$

$VC\text{-dim} = n+1$

Thm  $\text{opt}(C) \geq VC\text{-dim}(C)$

pf: take  $d$  points which yield  $VC\text{-dim}(C)$ , & force learning alg. to make  $d$  mistakes.

Suppose  $x^1, x^2, \dots$  are drawn independently according to some distribution  $P$  on  $X$ , and labelled according to some target concept  $c$ . Let  $M_A(t)$  denote probability that learning alg.  $A$  makes a mistake on  $t^{\text{th}}$  trial.

Thm:  $\exists A \text{ s.t. } (\forall t)(\forall P)(\forall c)$

$$M_A(t) \leq \frac{2 \cdot VC\text{-dim}(C)}{t} \quad (\text{Haussler, Littlestone, Warmuth})$$

(tight to within constant factors)