



• In other words, we remapped all vectors  $\mathbf{x}$  to  $\mathbf{y}$  such that the classes are linearly separable in the new vector space



- This is a <u>two-layer perceptron</u> or <u>two-layer</u> <u>feedforward neural network</u>
- Each neuron outputs 1 if its weighted sum exceeds its threshold, 0 otherwise

5

7



- Define the <u>p-dimensional unit hypercube</u> as  $H_p = \left\{ \begin{bmatrix} y_1, \dots, y_p \end{bmatrix}^T \in \Re^p, y_i \in [0, 1] \, \forall i \right\}$
- A hidden layer with p neurons maps an  $\ell$ -dim vector  $\mathbf{x}$  to a p-dim vector  $\mathbf{y}$  whose elements are corners of  $H_p$ , i.e.  $y_i \in \{0, 1\} \forall i$
- Each of the p neurons corresponds to an  $\ell$ -dim hyperplane
- The intersection\* of the (pos. or neg.) half-spaces from these p hyperplanes maps to a vertex of  ${\cal H}_p$
- If the classes of  $H_p$ 's vertices are linearly separable, then a perfect two-layer network exists
- I.e. a 2-layer network can separate classes consisting of unions of <u>adjacent</u> polyhedra

\*Also known as polyhedra.

#### What Else Can We Do with Two Layers?



### Three-Layer Networks

- With two-layer networks, there exist unions of polyhedra not linearly separable on  ${\cal H}_p$
- I.e. there exist assignments of classes to points on *H*<sub>p</sub> that are not linearly separable
- Solution: Add a second hidden layer of q neurons to partition  $H_p$  into regions based on class
- Output layer combines appropriate regions
- E.g. including 110 from Slide 6 in  $\omega_1$  is possible using procedure similar to XOR solution
- In general, can always use simple procedure of isolating each  $\omega_1$  node in  $H_p$  with its own second-layer hyperplane and taking disjunction
- Thus, can use 3-layer network to perfectly classify any union polyhedral regions

# The Backpropagation Algorithm

- A popular way to train a neural network
- Assume the architecture is <u>fixed</u> and complete
  - $\cdot \ k_r =$  number of nodes in layer r (could have  $k_L > 1)$
  - $\cdot \ w_{ji}^r =$  weight from neuron i in layer r-1 to neuron j in layer r
  - $v_{j}^{r} = \sum_{k=1}^{k_{r-1}} w_{jk}^{r} y_{k}^{r-1} + w_{j0}^{r}$
  - $\cdot \ y_{j}^{r} = f\left(v_{j}^{r}\right) = \text{output of neuron } j \text{ in layer } r$
- During training we'll attempt to minimize a cost function, so use differentiable activation func. *f*, e.g.:

$$f(v) = \frac{1}{1 + e^{-av}} \in [0, 1]$$
$$\underbrace{OR}{f(v) = c \tanh(av)} \in [-c, c]$$

.



### The Backpropagation Algorithm



# The Backpropagation Algorithm Intuition

- Recall derivation of Perceptron update rule:
  - Cost function:

$$U(\mathbf{w}) = \sum_{i=1}^{\ell} (w_i(t+1) - w_i(t))^2 + \eta \left( y(t) - \sum_{i=1}^{\ell} w_i(t+1)x_i(t) \right)^2$$

- Take gradient w.r.t. w(t + 1), set to 0, approximate, and solve:

$$w_i(t+1) = w_i(t) +$$
$$\eta \left( y(t) - \sum_{i=1}^{\ell} w_i(t) x_i(t) \right) x_i(t)$$

11

9

#### The Backpropagation Algorithm Intuition: Output Layer

- Now use similar idea with *j*th node of output layer (layer *L*):
  - Cost function:

$$U\left(\mathbf{w}_{j}^{L}\right) = \sum_{k=1}^{k_{L-1}} \left(w_{jk}^{L}(t+1) - w_{jk}^{L}(t)\right)^{2} + \eta \left[\underbrace{v_{jk}^{correct}}_{y_{j}(t)} - f\left(\underbrace{\sum_{k=1}^{k_{L-1}} w_{jk}^{L}(t+1)y_{k}^{L-1}(t)}_{k=1}\right)^{2}\right]$$

- Take gradient w.r.t.  $\mathbf{w}_j^L(t+1)$  and set to 0:

$$0 = 2 \left( w_{jk}^{L}(t+1) - w_{jk}^{L}(t) \right)$$
$$- 2\eta \left[ y_{j}(t) - f \left( \sum_{k=1}^{k_{L-1}} w_{jk}^{L}(t+1) y_{k}^{L-1}(t) \right) \right]$$
$$\cdot f' \left( \sum_{k=1}^{k_{L-1}} w_{jk}^{L}(t+1) y_{k}^{L-1}(t) \right) y_{k}^{L-1}(t)$$

13

## The Backpropagation Algorithm Intuition: The Other Layers

- How can we compute the "error term" for the hidden layers r < L when there is no "target vector" y for these layers?
- Instead, propagate back error values from output layer toward input layers, scaling with the weights
- Scaling with the weights characterizes how much of the error term each hidden unit is "responsible for":

$$w_{jk}^{r}(t+1) = w_{jk}^{r}(t) + \eta y_{k}^{r-1}(t) \,\delta_{j}^{r}(t)$$

where

$$\delta_j^r(t) = f'\left(v_j^r(t)\right) \sum_{k=1}^{k_{r+1}} \delta_k^{r+1}(t) \, w_{kj}^{r+1}(t)$$

• Derivation comes from computing gradient of cost function w.r.t.  $\mathbf{w}_{i}^{r}(t+1)$  via chain rule

15

The Backpropagation Algorithm Intuition: Output Layer (cont'd)

• Again, approximate and solve for  $w_{ik}^L(t+1)$ :

$$w_{jk}^{L}(t+1) = w_{jk}^{L}(t) + \eta y_{k}^{L-1}(t) \cdot \left[ y_{j}(t) - f\left(\sum_{k=1}^{k_{L-1}} w_{jk}^{L}(t)y_{k}^{L-1}(t)\right) \right] \cdot f'\left(\sum_{k=1}^{k_{L-1}} w_{jk}^{L}(t)y_{k}^{L-1}(t)\right)$$

• So:

$$w_{jk}^{L}(t+1) = w_{jk}^{L}(t) + \eta y_{k}^{L-1}(t) \underbrace{\left(y_{j}(t) - f\left(v_{j}^{L}(t)\right)\right) f'\left(v_{j}^{L}(t)\right)}_{\delta_{j}^{L}(t) = "\underline{\text{error term}}"}$$

• For 
$$f(v) = 1/(1 + \exp(-av))$$
:  
 $\delta_j^L(t) = a \cdot y_j^L(t) \cdot (y_j(t) - y_j^L(t)) (1 - y_j^L(t))$   
where  $y_j(t) = \text{target}$  and  $y_j^L(t) = \text{output}$ 

14

## The Backpropagation Algorithm Example



eta	0.3		
	trial 1	trial 2	
w ca	0.1	0.1008513	0.1008513
w_cb	0.1	0.1	0.0987985
w_c0	0.1	0.1008513	0.0996498
а	1	0	
b	0	1	
const	1	1	
sum_c	0.2	0.2008513	
y_c	0.5498340	0.5500447	
w_dc	0.1	0.1189104	0.0964548
w_d0	0.1	0.1343929	0.0935679
sum_d	0.1549834	0.1997990	
y_d	0.5386685	0.5497842	
target	1	0	
delta_d	0.1146431	-0.136083	
delta_c	0.0028376	-0.004005	
delta_d(t)	$= y_d(t) * (y(t) + y_d(t))$	- y_d(t)) * (1 -	y_d(t))
delta_c(t)	= y_c(t) * (1 - y_	_c(t)) * delta_	d(t) * w_dc(
w_dc(t+1)	$= w_dc(t) + eta$	* y_c(t) * del	ta_d(t)
w_ca(t+1)	= w ca(t) + eta	* a * delta_c	(t)



# Variations

(cont'd)

- Can implement a "backprop" scheme with EG
- Other nonlinear optimization schemes:
  - Conjugate gradient
  - Newton's method
  - Genetic algorithms
  - Simulated annealing
- Other cost functions, e.g. cross-entropy:

$$-\sum_{k=1}^{k_L} \left( \underbrace{\frac{label}{y_k(t)} \ln \left( y_k^L(t) \right)}_{p_k(t)} + (1 - y_k(t)) \ln \left( 1 - y_k^L(t) \right) \right)$$

"blows up" if  $y_k(t) \approx 1$  and  $y_k^L(t) \approx 0$  or vice-versa (Section 4.8)

21

#### Sizing the Network

Pruning Techniques [Also see Bishop, Sec. 9.5]

- <u>Approach 1</u>: Train with backprop, periodically computing effect of varying  $w_i$  on cost func:
  - From Taylor series expansion (p. 109),

$$\overbrace{\delta J}^{cost\ change} \approx \frac{1}{2} \sum_{i} h_{ii} \, \delta w_i^2 \quad \text{where} \quad h_{ii} = \frac{\partial^2 J}{\partial^2 w_i}$$

- If  $h_{ii} w_i^2/2$  (saliency factor) small, then  $w_i$  doesn't have much impact and is removed
- Now continue training with backprop
- Example (Sec 4.10): 480 wts pruned to 25



## Sizing the Network

- Before training, need to choose appropriate number of layers and size of each layer
  - Too small: Cannot learn what features make same classes similar and separate classes different
  - Too large: Adapts to details of the particular training set and cannot generalize well (called <u>overfitting</u>)
  - Also, increasing size increases complexity
- Approaches:
  - <u>Analytical methods</u>: Use knowledge of data to est. number of needed layers and neurons
  - Pruning techniques: Start with a large network and periodically remove weights and neurons that don't affect output much
  - <u>Constructive techniques</u>: Start with small netw. and periodically add neurons and wts

22

Sizing the Network Pruning Techniques (cont'd) [Also see Bishop, Sec. 9.5]

• <u>Approach 2</u>: Train with backprop, but add to the cost function *J* a term that penalizes large weights:

$$J' = J + penalty$$

- If  $w_i$ 's contribution to network output is small, then its share of J is small
- So penalty term dominates  $w_i$ 's share of J', driving it down
- Periodically prune weights that get too low

#### Sizing the Network Constructive Techniques

Cascade Correlation [Also Bishop, Sec. 9.5]

- Start with no hidden units and train
- After training, if residual error too high, then add a hidden unit:



- Then continue training; if residual error still to high, add another hidden unit:
  - Same as HU1, but connect input units and HU1's output to inputs of HU2
- Limit the number added to avoid overfitting

25



- For arbitrary set of N points, there are  $2^N$  ways to classify them into  $\omega_1$  and  $\omega_2$  (i.e.  $2^N$  dichotomies)
- If classification done by a single hyperplane, then the number of <u>linear dichotomies</u> is

$$O(N, \ell) = 2 \sum_{i=0}^{\ell} {N-1 \choose i}$$
$$= 2^{N} \text{ if } N \le \ell + 1, \text{ else } < 2^{N}$$

14 linear dichotomies

8 linear dichotomies





27

## Generalized Linear Classifiers Section 4.12

- In XOR problem, used linear threshold funcs. in hidden layer to map non-lin. sep. classes to new space where they were lin. sep.
- Output layer gave sep. hyperplane in new space
- Replace hidden-layer lin. thresh. funcs. with family of <u>nonlinear</u> functions  $f_i : \Re^{\ell} \to \Re$ , i = 1, ..., k
- Hidden layer maps  $\mathbf{x} \in \Re^{\ell}$  to  $\mathbf{y} = [f_1(\mathbf{x}), \dots, f_k(\mathbf{x})]^T$ and output layer finds separating hyperplane:



• I.e. approximating separating surface as linear combination of interpolation functions:

$$g(\mathbf{x}) = w_0 + \sum_{i=1}^k w_i f_i(\mathbf{x})$$

26

Generalized Linear Classifiers Cover's Theorem (cont'd)

- Thus if dimensionality  $\ell \ge N-1$  then a <u>perfect</u> separating hyperplane is <u>guaranteed</u> to exist
- Otherwise (N > ℓ+1) the fraction of dichotomies that are linear dichotomies is



 For fixed N, mapping to higher dimensional space increases likelihood of ∃ of sep. hyperplane!





Support Vector Machines [See refs. on slides page]

- Introduced in 1992
- State-of-the-art technique for classification and regression
- Techniques can also be applied to e.g. clustering and principal components analysis
- Similar to polynomial classifiers and RBF networks in that it remaps inputs and then finds a hyperplane
  - Main difference is how it works
- Features of SVMs:
  - Maximization of margin
  - Duality
  - Use of <u>kernels</u>
  - Use of problem convexity to find classifier

Generalized Linear Classifiers Radial Basis Function Networks Choosing the Centers

- Randomly select from the training set
  - Might work well if training set representative of probability distribution over data
- Learn the  $c_i$ 's and  $\sigma_i^2$ 's via gradient descent
  - Frequently computationally complex
- First <u>cluster</u> the data (Chapters 11–16) and use results to find centers
- Use methods similar to constructive and pruning techniques when sizing neural network
  - Add new center when perceived as needed, delete unnecessary centers
  - E.g. if new input vector  ${\bf x}$  far from all current centers and error high, then new center necessary, so add  ${\bf x}$  as new center

34



- A hyperplane's margin  $\gamma$  is the shortest distance from it to any training vector
- Intuition: larger margin ⇒ higher confidence in classifier's ability to generalize
  - Guaranteed generalization error bound in terms of  $1/\gamma^2$
- Definition assumes linear separability (more general definitions exist that do not)

#### Support Vector Machines

Maximum-Margin Perceptron Algorithm

- $\mathbf{w}(0) \leftarrow 0, b(0) \leftarrow 0, k \leftarrow 0, R \leftarrow \max_{1 \le i \le N} \|\mathbf{x}_i\|_2$ (R = radius of ball centered at origin containing training vectors),  $y_i \in \{-1, +1\} \forall i$
- Update <u>slope</u> same as before, update <u>offset</u> differently
- While mistakes are made on training set
  - For i = 1 to N (= # training vectors)
    - \* If  $y_i (\mathbf{w}_k \cdot \mathbf{x}_i + b_k) \leq 0$ 
      - $\cdot \mathbf{w}_{k+1} \leftarrow \mathbf{w}_k + \eta y_i \mathbf{x}_i$
      - $\cdot b_{k+1} \leftarrow b_k + \eta y_i R^2$
      - $\cdot k \leftarrow k + 1$
- Final predictor:  $h(\mathbf{x}) = \operatorname{sgn}(\mathbf{w}_k \cdot \mathbf{x} + b_k)$
- 37

#### Support Vector Machines Duality

• Another way of representing predictor:

$$h(\mathbf{x}) = \operatorname{sgn} (\mathbf{w} \cdot \mathbf{x} + b) = \operatorname{sgn} \left( \sum_{i=1}^{N} (\alpha_i y_i \mathbf{x}_i) \cdot \mathbf{x} + b \right)$$
$$= \operatorname{sgn} \left( \sum_{i=1}^{N} \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{x}) + b \right)$$
$$(\alpha_i = \# \text{ mistakes on } \mathbf{x}_i, \ \eta > 0 \text{ ignored})$$

- So perceptron alg has equivalent dual form:
- $\alpha \leftarrow \mathbf{0}$ ,  $b \leftarrow \mathbf{0}$ ,  $R \leftarrow \max_{1 \le i \le N} \|\mathbf{x}_i\|_2$
- While mistakes are made in For loop

- For 
$$i = 1$$
 to  $N$  (= # training vectors)  
\* If  $y_i \left( \sum_{j=1}^N \alpha_j y_j \left( \mathbf{x}_j \cdot \mathbf{x}_i \right) + b \right) \le 0$   
 $\cdot \alpha_i \leftarrow \alpha_i + 1$   
 $\cdot b \leftarrow b + y_i R^2$ 

• Now data only in dot products

38

#### Kernels

- Duality lets us remap to many more features!
- Let  $\phi: \Re^{\ell} \to F$  be nonlinear map of f.v.s, so

$$h(\mathbf{x}) = \operatorname{sgn}\left(\sum_{i=1}^{N} \alpha_{i} y_{i} \left(\phi\left(\mathbf{x}_{i}\right) \cdot \phi\left(\mathbf{x}\right)\right) + b\right)$$

- Can we compute φ(x<sub>i</sub>) · φ(x) without evaluating φ(x<sub>i</sub>) and φ(x)? <u>YES!</u>
- $\mathbf{x} = [x_1, x_2], \ \mathbf{z} = [z_1, z_2]:$   $(\mathbf{x} \cdot \mathbf{z})^2 = (x_1 z_1 + x_2 z_2)^2$   $= x_1^2 z_1^2 + x_2^2 z_2^2 + 2 x_1 x_2 z_1 z_2$  $= \underbrace{[x_1^2, x_2^2, \sqrt{2} x_1 x_2]}_{\phi(\mathbf{x})} \cdot [z_1^2, z_2^2, \sqrt{2} z_1 z_2]$
- LHS requires 2 mults + 1 squaring to compute, RHS takes 3 mults
- In general,  $(\mathbf{x} \cdot \mathbf{z})^d$  takes  $\ell$  mults + 1 expon., vs.  $\binom{\ell+d-1}{d} \ge \left(\frac{\ell+d-1}{d}\right)^d$  mults if compute  $\phi$  first

39

## Kernels (cont'd)

- In general, a <u>kernel</u> is a function K such that  $\forall \mathbf{x}, \mathbf{z}, K(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{z})$
- Typically start with kernel and take the feature mapping that it yields
- E.g. Let  $\ell = 1, x = x, z = z, K(x, z) = \sin(x-z)$
- By Fourier expansion,

$$\sin(x-z) = a_0 + \sum_{n=1}^{\infty} a_n \sin(nx) \sin(nz) + \sum_{n=1}^{\infty} a_n \cos(nx) \cos(nz)$$

for Fourier coefficients  $a_0, a_1, \ldots$ 

 This is the dot product of two infinite sequences of nonlinear functions:

 $\{\phi_i(x)\}_{i=0}^{\infty} = [1, \sin(x), \cos(x), \sin(2x), \cos(2x), \ldots]$ 

• I.e. there are an infinite number of features in this remapped space!

## Support Vector Machines Finding a Hyperplane

- Can show [Cristianini & Shawe-Taylor] that if data linearly separable in remapped space, then get maximum margin classifier by minimizing  $\mathbf{w} \cdot \mathbf{w}$  subject to  $y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1$
- Can reformulate this into a <u>convex quadratic</u> <u>program</u>, which can be solved optimally, i.e. won't encounter local optima
- Can always find a kernel that will make training set linearly separable, but <u>beware of choosing a</u> <u>kernel that is too powerful</u> (overfitting)
- If kernel doesn't separate, can optimize subject to  $y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1 \xi_i$ , where  $\xi_i$  are <u>slack variables</u> that <u>soften</u> the margin (can still solve optimally)
- If number of training vectors is very large, may opt to <u>approximately</u> solve these problems to save time and space
- Use e.g. gradient ascent and sequential minimal optimization (SMO) [Cristianini & Shawe-Taylor]

41

## Decision Trees [Also Mitchell, ch. 3]



- Start at root and work down tree until leaf reached; output that classification
- E.g.  $\mathbf{x} = [1/2, 1/4]^T$  classified as  $\omega_3$



#### **Decision Trees**

Learning Good Trees [Also Mitchell, ch. 3]

• Feature at root is one that yields highest information gain, equivalent to max. reduction of <u>entropy</u> (class impurity) in training data:

S = set of N feature vectors  $N_i = \text{number in } \omega_i$ 

$$p_i = N_i / N$$
  $Ent(S) = \sum_{i=1}^{M} -p_i \log_2(p_i)$ 

• First partition along dimensions into set *A* of features and places where classes change, e.g.

 $A = \{(x_1, 0), (x_1, 1/4), (x_1, 1/2), (x_1, 3/4), (x_2, 0), (x_2, 1/2), (x_2, 3/4)\}$ 

• For  $a = (x_i, b) \in A$ , define

 $S_a = \{\mathbf{x} \in S : x_i > b\} \qquad S'_a = \{\mathbf{x} \in S : x_i \le b\}$  $Gain(S, a) = Ent(S) - \left(\frac{|S_a|}{|S|}Ent(S_a) + \frac{|S'_a|}{|S|}Ent(S'_a)\right)$  $= 0 \text{ for } (x_1, 1/4)$ 

- Choose a from A that maximizes Gain, place it at root, then recursively call on  $S_a$  and  $S'_a$
- Forms basis of algorithms ID3 and C4.5
- Can avoid overfitting by pruning