

Admin
Assignment 7

2

1

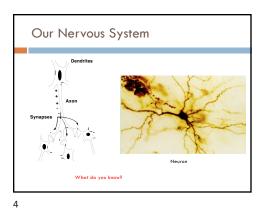
3

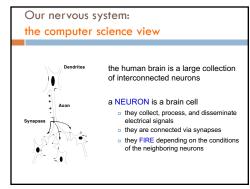
Perceptron learning algorithm

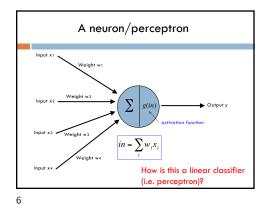
repeat until convergence (or for some # of iterations):
 for each training example $(f_1, f_2, ..., f_{n_t} | \text{label})$:
 $prediction = b + \sum_{i=1}^n w_i f_i$ if $prediction * | \text{label} \le 0$: // they don't agree
 for each w:
 w_i = w_i + f_**| \text{label}

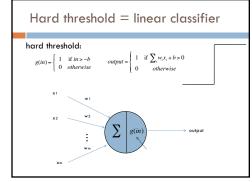
 b = b + label

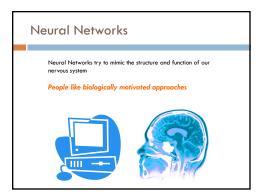
Why is it called the "perceptron" learning algorithm if what it learns is a line? Why not "line learning" algorithm?

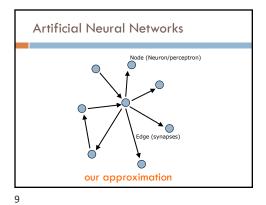










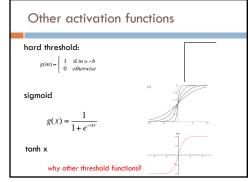


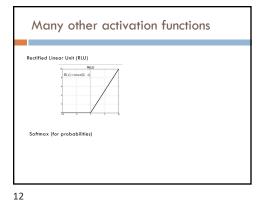
Node A Weight w Node B (perceptron) $output = \begin{cases} 1 & \text{if } \sum_i w_i x_i + b > 0 \\ 0 & \text{otherwise} \end{cases}$ W is the strength of signal sent between A and B.

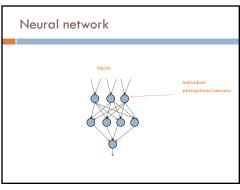
If A fires and w is **positive**, then A **stimulates** B.

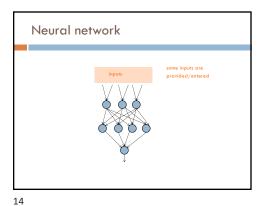
If A fires and w is **negative**, then A **inhibits** B.

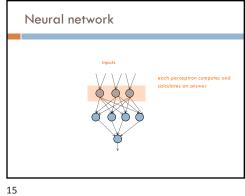
10

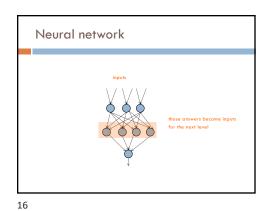


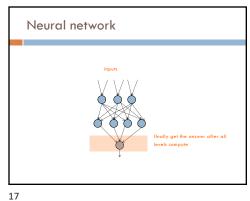




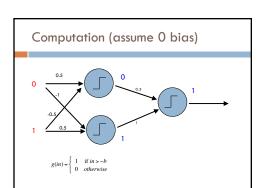


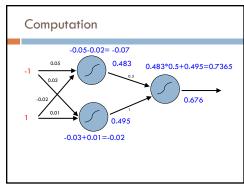


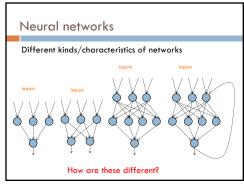


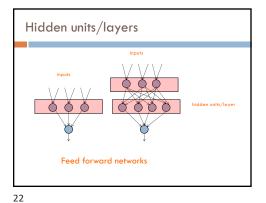


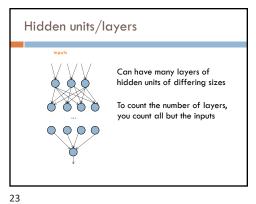


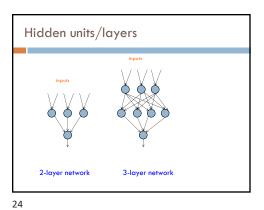


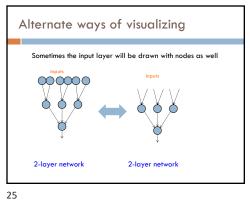


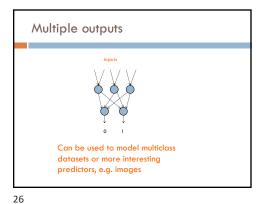


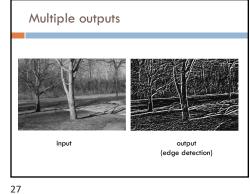


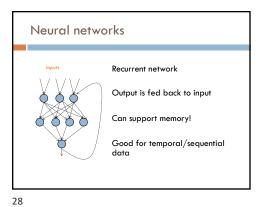


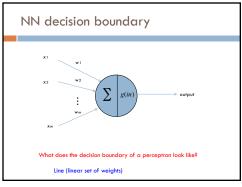


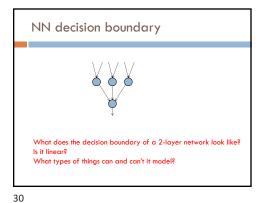


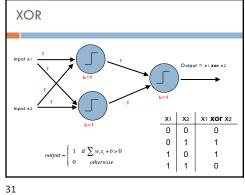


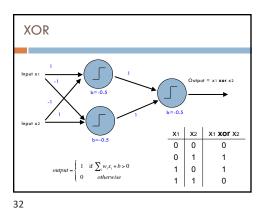


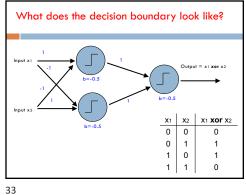


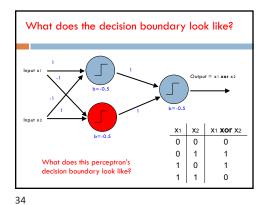


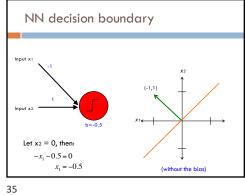


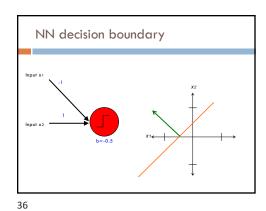


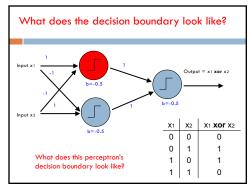


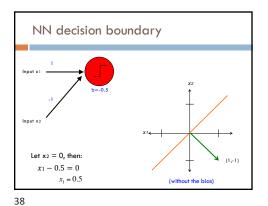


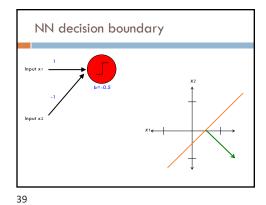


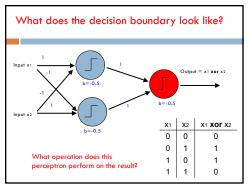


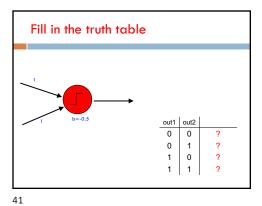


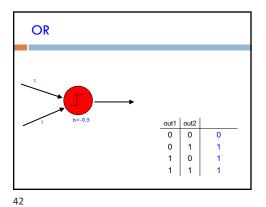






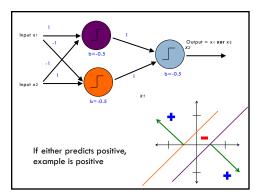


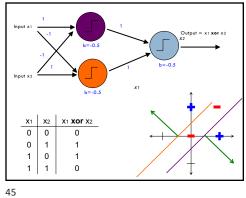


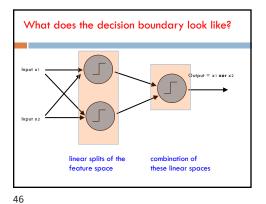


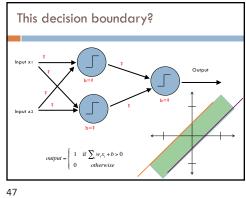
What does the decision boundary look like?

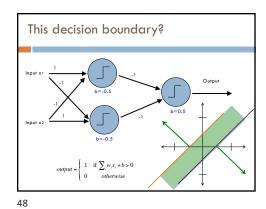
Input x1 b=.0.5 x_1 x_2 x_1

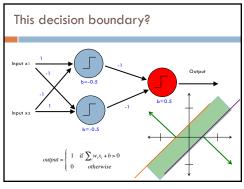


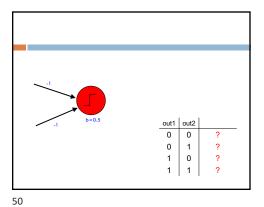


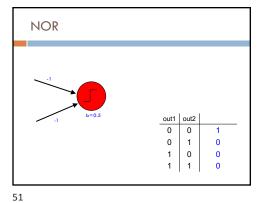


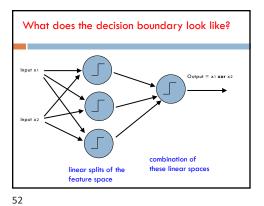


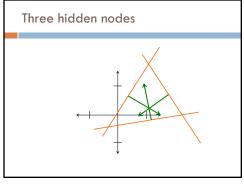












NN decision boundaries

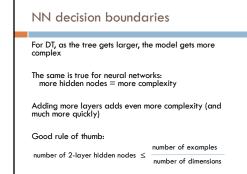
Theorem 9 (Two-Layer Networks are Universal Function Approximators). Let F be a continuous function on a bounded subset of D-dimensional space. Then there exists a two-layer neural network \hat{F} with a finite number of hidden units that approximate F arbitrarily well. Namely, for all x in the domain of F, $|F(x) - \hat{F}(x)| < \varepsilon$.

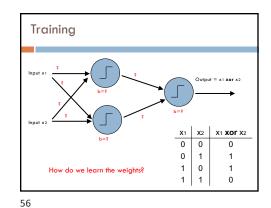
Put simply: two-layer networks can approximate any function

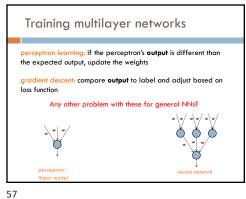
53

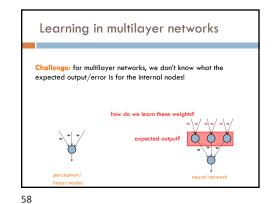
55

54

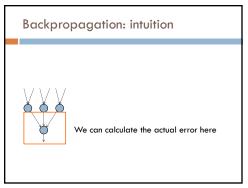




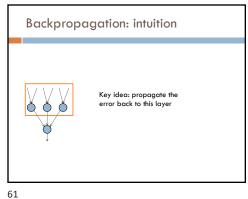


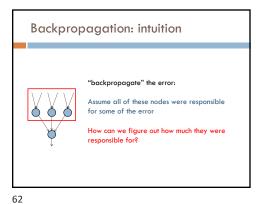


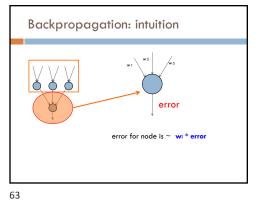
Backpropagation: intuition Gradient descent method for learning weights by optimizing a loss function 1. calculate output of all nodes 2. calculate the weights for the output layer based on 3. "backpropagate" errors through hidden layers

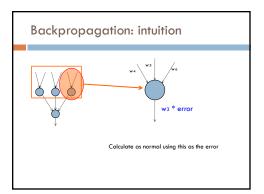


59 60









Backpropagation: the details

Gradient descent method for learning weights by optimizing a loss function

- 1. calculate output of all nodes
- 2. calculate the updates directly for the output layer
- 3. "backpropagate" errors through hidden layers

What loss function?

Backpropagation: the details

Gradient descent method for learning weights by optimizing a loss function

1. calculate output of all nodes

66

- 2. calculate the updates directly for the output layer
- 3. "backpropagate" errors through hidden layers

$$loss = \sum_{x} \frac{1}{2} (y - \hat{y})^2 \quad \text{squared error}$$

65