

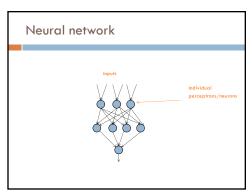
Admin

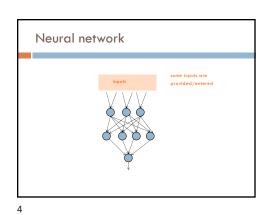
Assignment 7

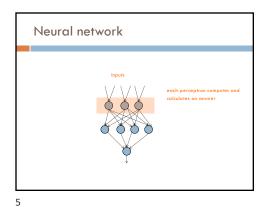
Assignment 8 released on Monday. Start ASAP!

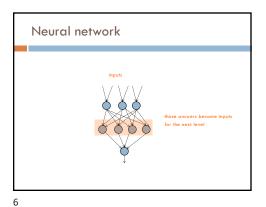
2

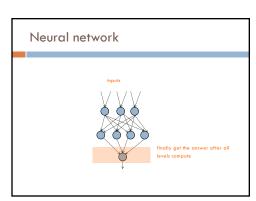
1

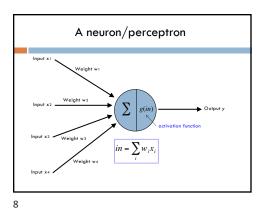


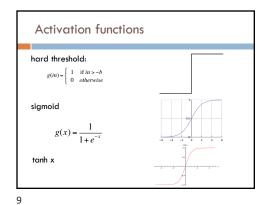


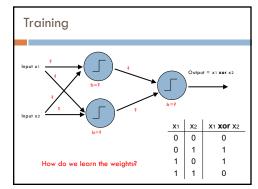












Learning in multilayer networks

Challenge: for multilayer networks, we don't know what the expected output/error is for the internal nodes!



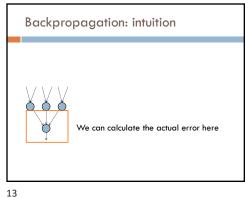
Backpropagation: intuition

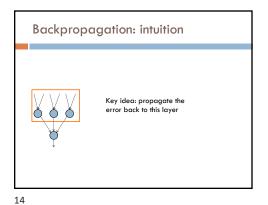
Gradient descent method for learning weights by optimizing a loss function

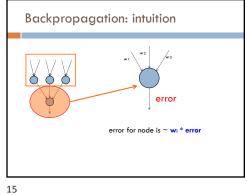
- 1. calculate output of all nodes
- calculate the weights for the output layer based on the error
- 3. "backpropagate" errors through hidden layers

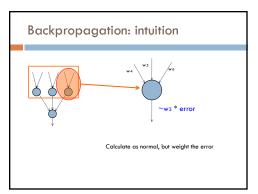
11

12









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

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

Notation:

m: features/inputs
d: hidden nodes
hk: output from
hidden node k

How many weights (ignore bias for now)?

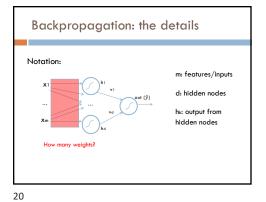
18

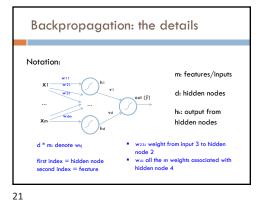
17

19

Notation:

m: features/inputs
d: hidden nodes
hi: output from
hidden nodes
d weights: denote vs.



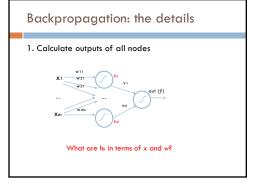


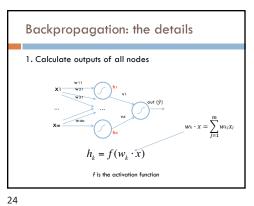
Backpropagation: the details

Gradient descent method for learning weights by optimizing a loss function $\underset{xy, \sum_{x} \frac{1}{2}(y-\hat{y})^2}{\operatorname{argmin}_{w, y} \sum_{x} \frac{1}{2}(y-\hat{y})^2}$ 1. calculate output of all nodes
2. calculate the updates directly for the output layer
3. "backpropagate" errors through hidden layers

22

۷.





Backpropagation: the details

1. Calculate outputs of all nodes

$$h_k = f(w_k \cdot x) = \frac{1}{1 + e^{-w_k}}$$

Backpropagation: the details 1. Calculate outputs of all nodes

What is out in terms of h and v?

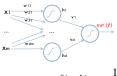
25

27

26

Backpropagation: the details

1. Calculate outputs of all nodes



$$out = f(v \cdot h) = \frac{1}{1 + e^{-v \cdot h}}$$

Backpropagation: the details 2. Calculate new weights for output layer Want to take a small step towards decreasing loss. How?

Recall: derivative chain rule

$$\frac{d}{dx}(f(g(x)) = ?$$

Recall: derivative chain rule

$$\frac{d}{dx}(f(g(x)) = f'(g(x))\frac{d}{dx}g(x)$$

29

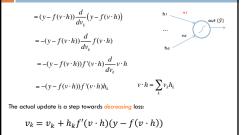
31

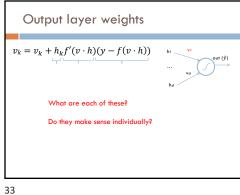
30

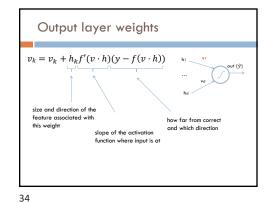
Output layer weights

$$\begin{aligned} & \operatorname{argmin}_{wy} \sum_{x} \frac{1}{2} (y - \hat{y})^{2} \\ & \xrightarrow{h_{1}} \underbrace{\frac{dloss}{dv_{k}} = \frac{d}{dv_{k}} \left(\frac{1}{2} (y - \hat{y})^{2} \right)}_{\text{hd}} \\ & = \frac{d}{dv_{k}} \underbrace{\frac{1}{2} (y - f(v \cdot h))^{2}}_{\text{d}v_{k}} \underbrace{\hat{y} = f(v \cdot h)}_{\text{f}} \\ & = (y - f(v \cdot h)) \frac{d}{dv_{k}} (y - f(v \cdot h)) \end{aligned}$$

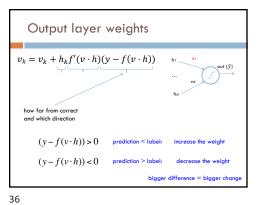
Output layer weights

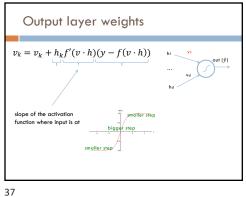


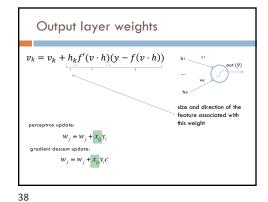




Output layer weights $v_k = v_k + h_k f'(v \cdot h)(y - f(v \cdot h))$ how far from correct and which direction $(y - f(v \cdot h)) > 0$ $(y - f(v \cdot h)) < 0$



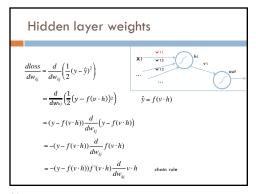


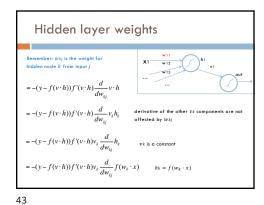


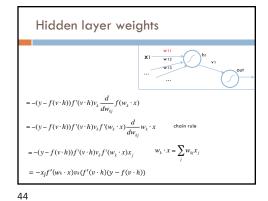
Backpropagation: the details Gradient descent method for learning weights by optimizing a loss function 1. calculate output of all nodes 2. calculate the updates directly for the output layer "backpropagate" errors through hidden layers

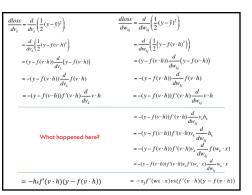
39

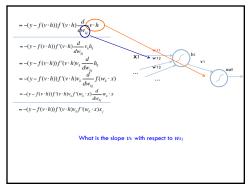
Backpropagation 3. "backpropagate" errors through hidden layers Want to take a small step towards decreasing loss. How? 40

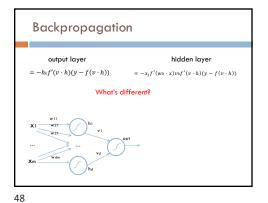


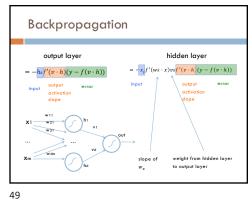


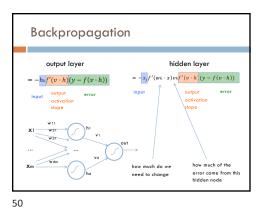


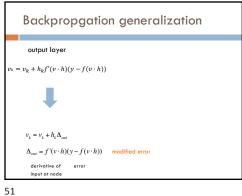


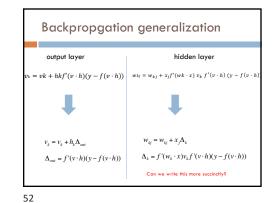


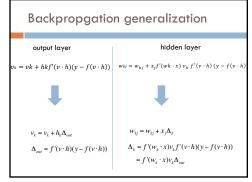


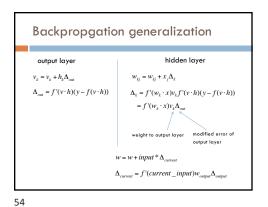


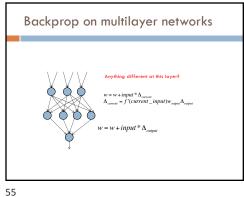


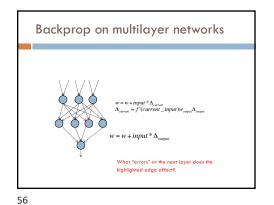


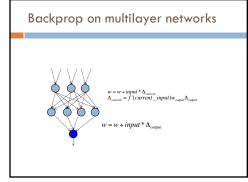


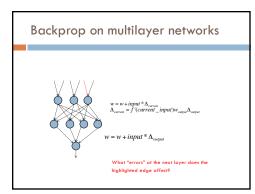


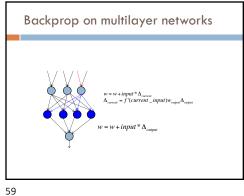


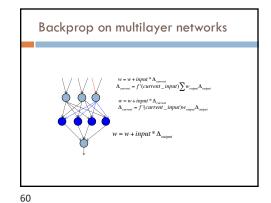


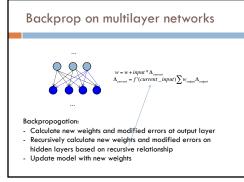


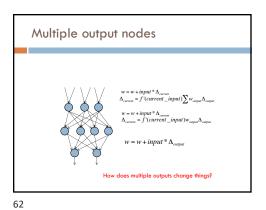


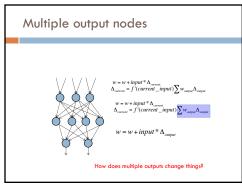












Backpropagation implementation

Output layer update: $v_k = v_k + h_k (y - f(v \cdot h)) f'(v \cdot h)$ Hidden layer update: $w_{ij} = w_{ij} + x_j f'(w_k \cdot x) v_k f'(v \cdot h) (y - f(v \cdot h))$ Any missing information for implementation?

63

65

Backpropagation implementation

Output layer update: $v_k = v_k + h_k(y - \overline{f(v \cdot h)}) f'(v \cdot h)$ Hidden layer update: $w_{ij} = w_{ij} + x_j \overline{f'(w_k \cdot x)} v_k \overline{f'(v \cdot h)} (y - \overline{f(v \cdot h)})$ 1. What activation function are we using

2. What is the derivative of that activation function

Activation function derivatives $s(x) = \frac{1}{1 + e^{-x}}$ s'(x) = s(x)(1 - s(x)) tonh $\frac{d}{dx} \tanh(x) = 1 - \tanh^2 x$

66

Learning rate

Output layer update:

 $v_k = v_k + \frac{\eta}{\eta} h_k (y - f(v \cdot h)) f'(v \cdot h)$

Hidden layer update:

67

69

 $w_{kj} = w_{kj} + \eta x_j f'(w_k \cdot x) v_k f'(v \cdot h) (y - f(v \cdot h))$

- Like gradient descent for linear classifiers, use a learning rate
- Often will start larger and then get smaller

Backpropagation implementation

Just like gradient descent!

for some number of iterations:

randomly shuffle training data

for each example:

- Compute all outputs going forward
- Calculate new weights and modified errors at output layer
- Recursively calculate new weights and modified errors on hidden layers based on recursive relationship
- Update model with new weights

68

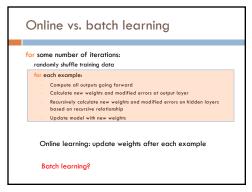
Handling bias

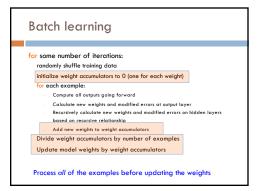
x1 w21 hb1 v1

xm wdm+1 hd 1

1. Add an extra feature hard-wired to 1 to all the

2. For other layers, add an extra parameter whose input is always 1





Many variations

Momentum: include a factor in the weight update to keep moving in the direction of the previous update

Mini-batch:

Compromise between online and batch
Avoids noisiness of updates from online while making more educated weight updates

Simulated annealing:

With some probability make a random weight update
Reduce this probability over time

Picking network configuration

Can be slow to train for large networks and large amounts of data

Loss functions (including squared error) are generally not convex with respect to the parameter space

73 74

History of Neural Networks

McCulloch and Pitts (1943) – introduced model of artificial neurons and suggested they could learn

Hebb (1949) – Simple updating rule for learning

Rosenblatt (1962) - the perceptron model

Minsky and Papert (1969) – wrote Perceptrons

Bryson and Ho (1969, but largely ignored until 1980s-Rosenblatt) – invented backpropagation learning for multilayer networks



http://www.nytimes.com/2012/06/26/technol ogy/in-a-big-network-of-computers-evidence-of-machine-learning.html?_r=0

77