

Neural Networks (Part 1)

Mark Craven and David Page
Computer Sciences 760
Spring 2018

www.biostat.wisc.edu/~craven/cs760/

Some of the slides in these lectures have been adapted/borrowed from materials developed by Tom Dietterich, Pedro Domingos, Tom Mitchell, David Page, and Jude Shavlik

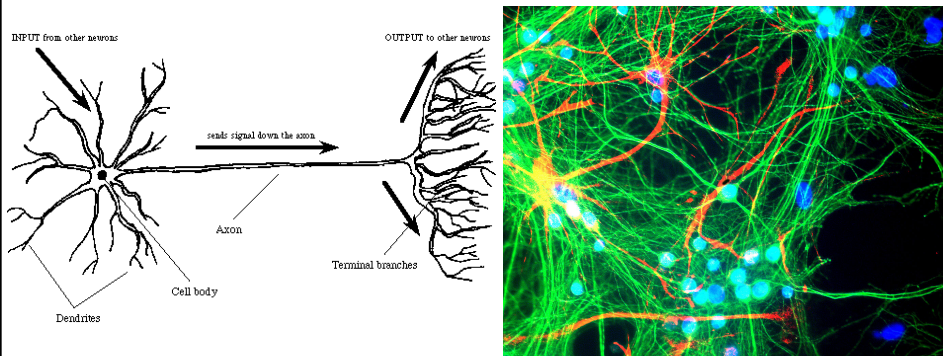
Goals for the lecture

you should understand the following concepts

- perceptrons
- the perceptron training rule
- linear separability
- hidden units
- multilayer neural networks
- gradient descent
- stochastic (online) gradient descent
- activation functions
 - sigmoid, hyperbolic tangent, ReLU
- objective (error, loss) functions
 - squared error, cross entropy

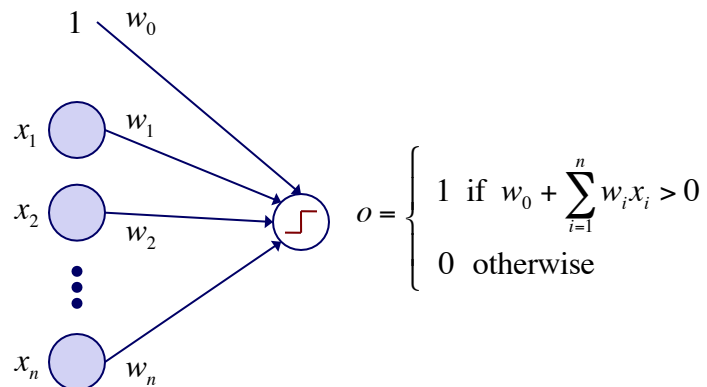
Neural networks

- a.k.a. *artificial neural networks*, *connectionist models*
- inspired by interconnected neurons in biological systems
 - simple processing units
 - each unit receives a number of real-valued inputs
 - each unit produces a single real-valued output



Perceptrons

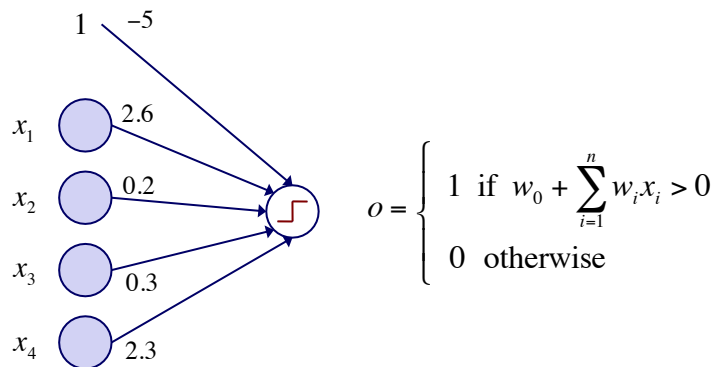
[McCulloch & Pitts, 1943; Rosenblatt, 1959; Widrow & Hoff, 1960]



input units:
represent given x

output unit:
represents binary classification

Perceptron example



features, class labels are represented numerically

$$\begin{aligned} \mathbf{x} &= \langle 1, 0, 0, 1 \rangle & w_0 + \sum_{i=1}^n w_i x_i &= -0.1 & o &= 0 \\ \mathbf{x} &= \langle 1, 0, 1, 1 \rangle & w_0 + \sum_{i=1}^n w_i x_i &= 0.2 & o &= 1 \end{aligned}$$

Learning a perceptron: the perceptron training rule

1. randomly initialize weights
2. iterate through training instances until convergence

2a. calculate the output
for the given instance

$$o = \begin{cases} 1 & \text{if } w_0 + \sum_{i=1}^n w_i x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

2b. update each weight

$$\Delta w_i = \eta (y - o) x_i$$

η is learning rate;
set to value $\ll 1$

$$w_i \leftarrow w_i + \Delta w_i$$

Representational power of perceptrons

perceptrons can represent only *linearly separable* concepts

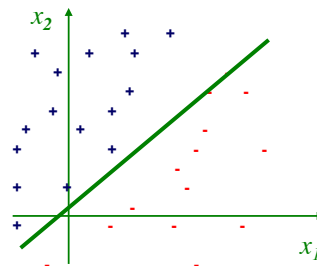
$$o = \begin{cases} 1 & \text{if } w_0 + \sum_{i=1}^n w_i x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

decision boundary given by:

$$1 \text{ if } w_0 + w_1 x_1 + w_2 x_2 > 0$$

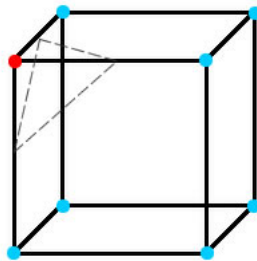
$$w_1 x_1 + w_2 x_2 = -w_0$$

$$x_2 = -\frac{w_1}{w_2} x_1 - \frac{w_0}{w_2}$$



Representational power of perceptrons

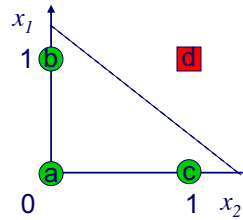
- in previous example, feature space was 2D so decision boundary was a line
- in higher dimensions, decision boundary is a hyperplane



Some linearly separable functions

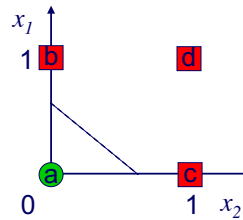
AND

	x_1	x_2	y
a	0	0	0
b	0	1	0
c	1	0	0
d	1	1	1



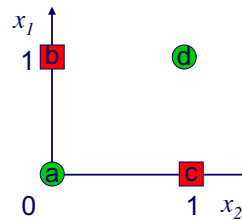
OR

	x_1	x_2	y
a	0	0	0
b	0	1	1
c	1	0	1
d	1	1	1

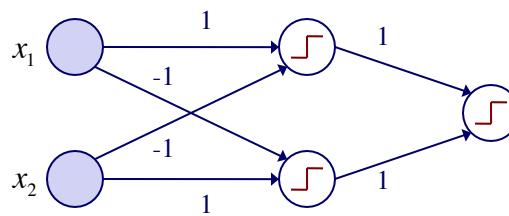


XOR is not linearly separable

	x_1	x_2	y
a	0	0	0
b	0	1	1
c	1	0	1
d	1	1	0



a multilayer perceptron
can represent XOR



assume $w_0 = 0$ for all nodes

Example multilayer neural network

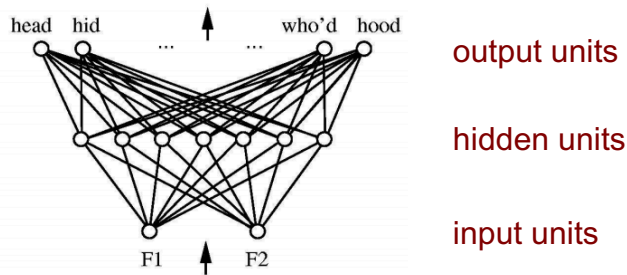


figure from Huang & Lippmann, *NIPS* 1988

input: two features from spectral analysis of a spoken sound

output: vowel sound occurring in the context “h__d”

Decision regions of a multilayer neural network

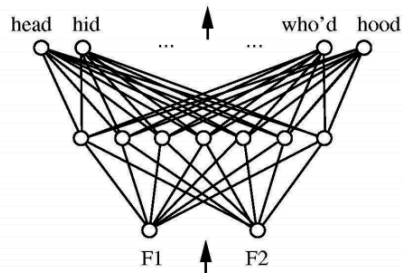
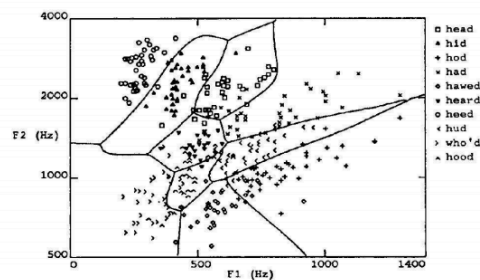


figure from Huang & Lippmann, *NIPS* 1988

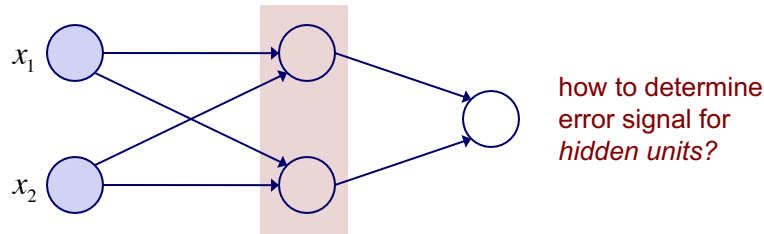


input: two features from spectral analysis of a spoken sound

output: vowel sound occurring in the context “h__d”

Learning in multilayer networks

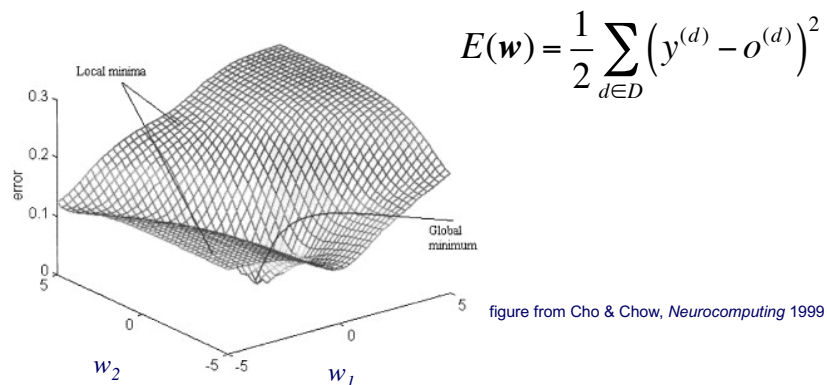
- work on neural nets fizzled in the 1960's
 - single layer networks had representational limitations (linear separability)
 - no effective methods for training multilayer networks



- revived again with the invention of *backpropagation* method [Rumelhart & McClelland, 1986; also Werbos, 1975]
 - key insight: require neural network to be differentiable; use *gradient descent*

Gradient descent in weight space

Given a training set $D = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$ we can specify an error measure that is a function of our weight vector \mathbf{w}



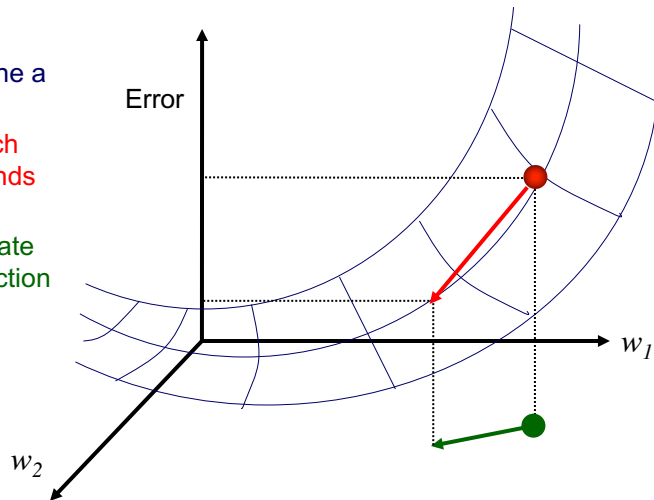
This *objective function* defines a surface over the model (i.e. weight) space

Gradient descent in weight space

gradient descent is an iterative process aimed at finding a minimum in the error surface

on each iteration

- current weights define a point in this space
- find direction in which error surface descends most steeply
- take a step (i.e. update weights) in that direction



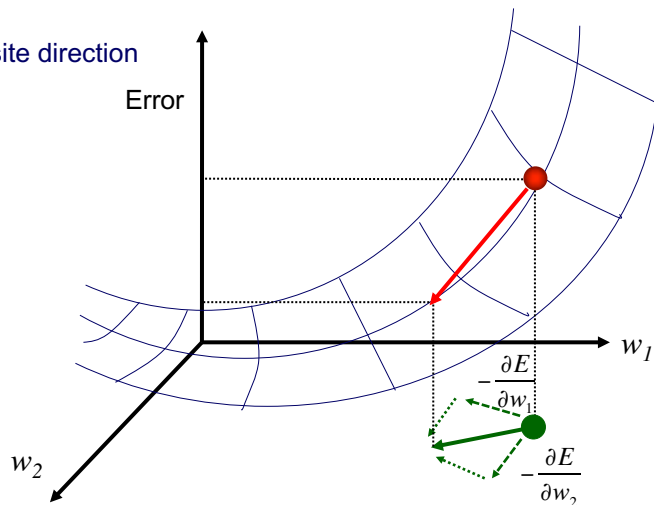
Gradient descent in weight space

calculate the gradient of E : $\nabla E(\mathbf{w}) = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$

take a step in the opposite direction

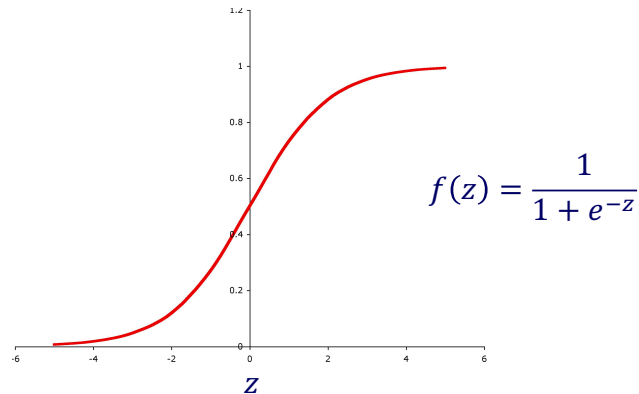
$$\Delta \mathbf{w} = -\eta \nabla E(\mathbf{w})$$

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$



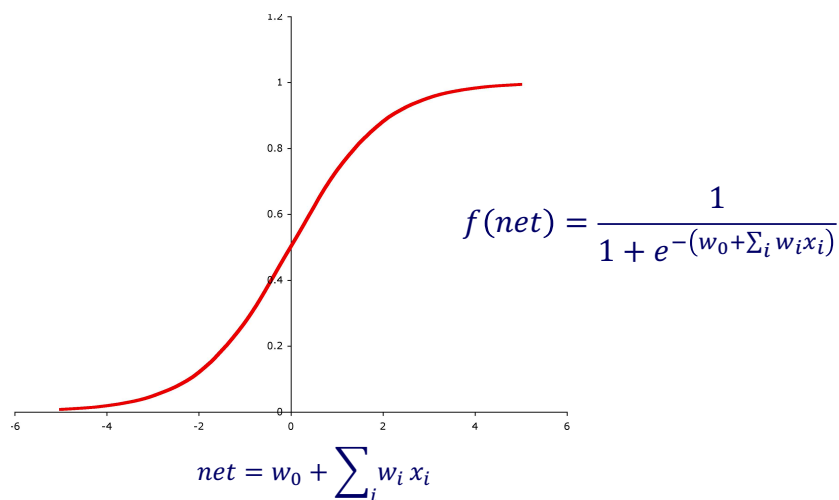
The sigmoid function

- to be able to differentiate E with respect to w_i , our network must represent a continuous function
- to do this, we can use *sigmoid functions* instead of threshold functions in our hidden and output units



The sigmoid function

for the case of a single-layer network



Batch neural network training

given: network structure and a training set $D = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$

initialize all weights in \mathbf{w} to small random numbers

until stopping criteria met do

 initialize the error $E(\mathbf{w}) = 0$

 for each $(\mathbf{x}^{(d)}, y^{(d)})$ in the training set

 input $\mathbf{x}^{(d)}$ to the network and compute output $o^{(d)}$

 increment the error $E(\mathbf{w}) = E(\mathbf{w}) + \frac{1}{2}(y^{(d)} - o^{(d)})^2$

 calculate the gradient

$$\nabla E(\mathbf{w}) = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

 update the weights

$$\Delta \mathbf{w} = -\eta \nabla E(\mathbf{w})$$

Online vs. batch training

- Standard gradient descent (batch training): calculates error gradient for the entire training set, before taking a step in weight space
- *Stochastic gradient descent* (online training): calculates error gradient for a single instance (or a small set of instances, a “mini batch”), then takes a step in weight space
 - much faster convergence
 - less susceptible to local minima

Online neural network training (stochastic gradient descent)

given: network structure and a training set $D = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$

initialize all weights in \mathbf{w} to small random numbers

until stopping criteria met do

for each $(\mathbf{x}^{(d)}, y^{(d)})$ in the training set

input $\mathbf{x}^{(d)}$ to the network and compute output $o^{(d)}$

calculate the error $E(\mathbf{w}) = \frac{1}{2} (y^{(d)} - o^{(d)})^2$

calculate the gradient

$$\nabla E(\mathbf{w}) = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

update the weights

$$\Delta \mathbf{w} = -\eta \nabla E(\mathbf{w})$$

Other activation functions

- the sigmoid is just one choice for an *activation function*
- there are others we can use including

hyperbolic tangent

$$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$$

rectified linear (ReLU)

$$f(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$

Other objective functions

- squared error is just one choice for an *objective function*
- there are others we can use including

cross entropy

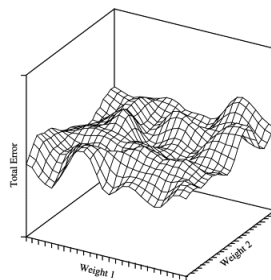
$$E(\mathbf{w}) = \sum_{d \in D} -y^{(d)} \ln(o^{(d)}) - (1 - y^{(d)}) \ln(1 - o^{(d)})$$

multiclass cross entropy

$$E(\mathbf{w}) = - \sum_{d \in D} \sum_{i=1}^{\# \text{ classes}} y_i^{(d)} \ln(o_i^{(d)})$$

Convergence of gradient descent

- gradient descent will converge to a minimum in the error function
- for a multi-layer network, this may be a *local minimum* (i.e. there may be a “better” solution elsewhere in weight space)



- for a single-layer network, this will be a global minimum (i.e. gradient descent will find the “best” solution)
- Recent analysis suggests that local minima are probably rare in high dimensions; saddle points are more of a challenge [Dauphin et al., NIPS 2014]