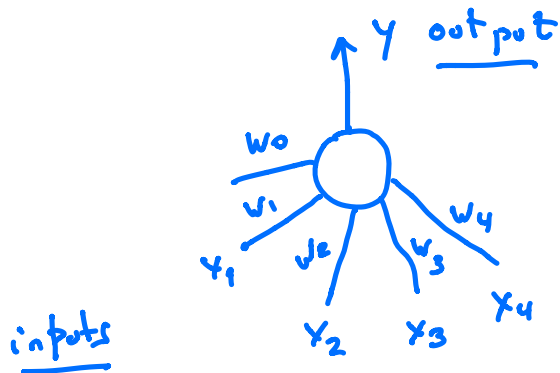


W08 - Neural Networks . Learning as Inference

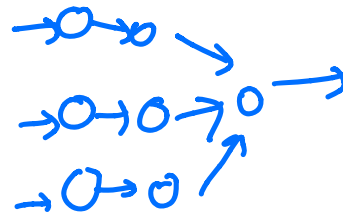
Mackay's lectures 15, 16. Chapters 39, 41, 42.

Feedforward Networks

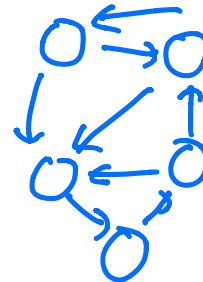
A single neuron



Feedforward network



Feedback network



Elements of the network

• inputs x_1, \dots, x_I

• output $y = \text{activity}$

• weights (parameters)

$\gamma > \text{threshold}$, neuron fires

w_0, w_1, \dots, w_I

This is not a probabilistic process

Given the inputs, how does the neuron produce y ?

i) neuron adds all weights

$$a = w_0 + \sum_{i=1}^I w_i x_i = \text{activation}$$

↑
bias

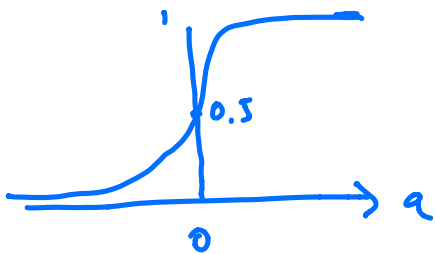
ii) the activity of the neuron y is a function of the activation a , $\gamma(a)$

$\gamma(a)$ is generally a "threshold" function

Several commonly used forms for the activity $\gamma(a)$

are:

The linear logistic function

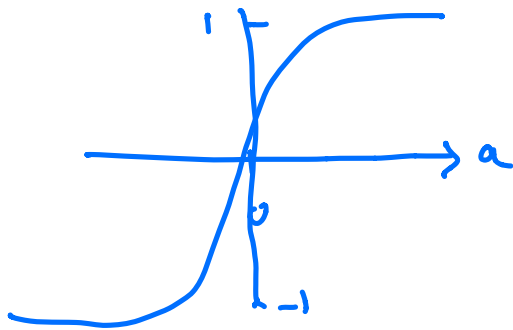


$$\gamma(a) = \frac{1}{1 + e^{-a}} = \frac{1}{1 + e^{-\bar{w}\bar{x}}}$$

$$\bar{w}\bar{x} = \sum_i w_i x_i + w_0$$

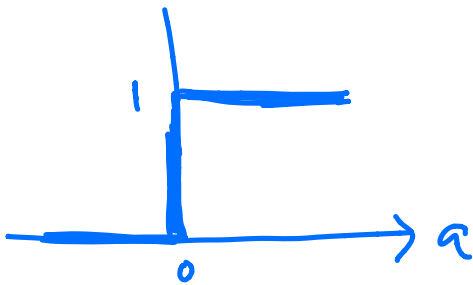
$$= \sum_{i=0}^I w_i x_i \quad x_0 = 1$$

• the sigmoidal function



$$\gamma(a) = \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$

• The step function



$$\gamma(a) = \begin{cases} 1 & a > 0 \\ 0 & a \leq 0 \end{cases}$$

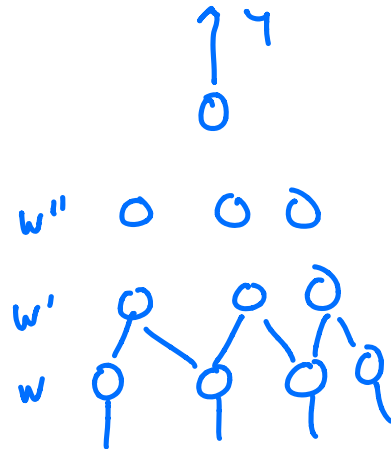
Recap neural network

i) architecture

w_i

ii) activity $a = \sum_{i=0}^I w_i x_i = \bar{w} \bar{x}$

iii) the activity rule $\gamma(a) = \frac{1}{1+e^{-a}} = \frac{1}{1+e^{-\bar{w}\bar{x}}}$



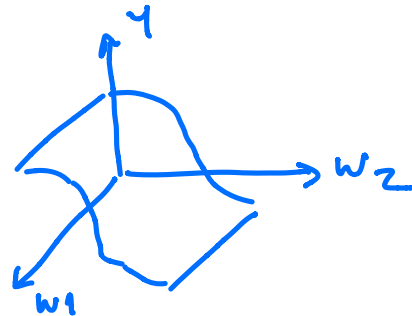
the activity $y(a)$ can be interpreted as:

The probability according to the neuron (weights)
that the inputs (\bar{x}) deserves a response

$y \approx 1$ response

$y \approx 0$ no response

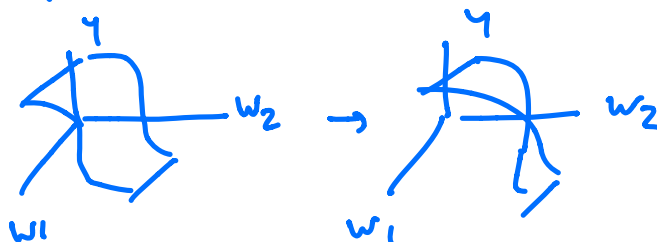
The Space of Weights



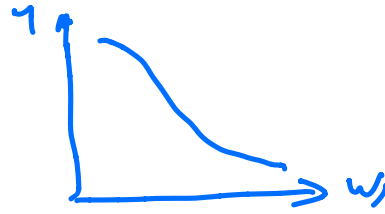
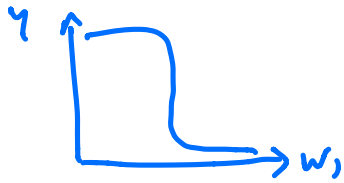
+ change in bias w_0



+ change in w_1/w_2 (a twist)



+ scale of weight



Contour plots



What can a single neuron learn?

to be a classifier!

Idea of supervised learning is:

given a number of examples of
input vectors $\bar{x}^{(1)} \dots \bar{x}^{(n)}$
and their outputs $t^{(1)} \dots t^{(n)}$

make the network learn their relationship.

Find the values of the weights so that $y^{(n)} \approx t^{(n)}$

"learning" ~ "finding parameters"

Classification problem

A $\rightarrow t=1$

0 $\rightarrow t=0$

w?

$$y(A|w) \approx 1$$

$$y(0|w) \approx 0$$

$$\text{error} = \left\{ y_w^{(1)} - t^{(1)}, \dots, y_w^{(n)} - t^{(n)} \right\}$$

learning = adjust w's so that error is small.

The Error Function

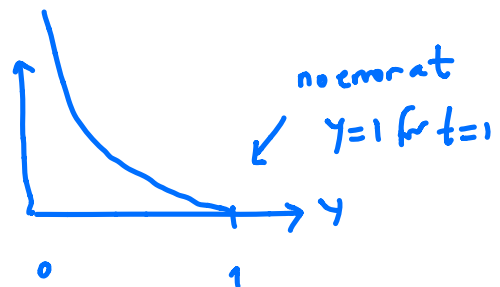
$$\text{inputs: } \left\{ \bar{x}^{(n)}, t^{(n)} \right\}_{n=1}^N \quad \left\{ x_i^{(n)} \right\}_{i=0}^I$$

$$\text{outputs } y^{(n)}(\bar{x}^{(n)}, \bar{w}) = \frac{1}{1 + e^{-\bar{w} \bar{x}^{(n)}}}$$

Introduce the error function

$$\zeta(w) = - \sum_{n=1}^N \left[t^{(n)} \cdot \log y^{(n)} + (1 - t^{(n)}) \log (1 - y^{(n)}) \right]$$

why? $t^{(n)} = 1 : -\log y^{(n)}$



$t^{(n)} = 0 : -\log(1 - y^{(n)})$



$$\underbrace{\zeta(w)} = - \sum_{n=1}^N \left[t^{(n)} \log \gamma^{(n)}(\bar{w}, \bar{x}) + (1 - t^{(n)}) \log (1 - \gamma^{(n)}(\bar{w}, \bar{x})) \right]$$

Back propagation

want to minimize error $\zeta(w)$

$$\zeta(\bar{w}) \geq 0 \quad \parallel \quad \zeta(w) = 0 \Leftrightarrow \gamma^{(n)} = t^{(n)}$$

let's take the derivative of $\zeta(\bar{w})$

$$\begin{aligned} \frac{\delta \zeta(\bar{w})}{\delta w_j} &= - \sum_n \left[\frac{t^{(n)}}{\gamma^{(n)}} - \frac{1 - t^{(n)}}{1 - \gamma^{(n)}} \right] \cdot \frac{\delta \gamma^{(n)}}{\delta w_i} \\ &= - \sum_n \frac{t^{(n)}(1 - \gamma^{(n)}) - \gamma^{(n)}(1 - t^{(n)})}{\gamma^{(n)}(1 - \gamma^{(n)})} \cdot \frac{\delta \gamma^{(n)}}{\delta w_i} \\ &= - \sum_n \frac{t^{(n)} - \gamma^{(n)}}{\gamma^{(n)}(1 - \gamma^{(n)})} \cdot \frac{\delta \gamma^{(n)}}{\delta w_i} \end{aligned}$$

$$y^{(n)}(\bar{w}, \bar{x}) = \frac{1}{1 + e^{-\bar{w} \bar{x}^{(n)}}}$$

$$\frac{\delta y^{(n)}}{\delta w_i} = (-1) (1 + e^{-\bar{w} \bar{x}^{(n)}})^{-2} \cdot \frac{\delta}{\delta w_i} \left[e^{-\bar{w} \bar{x}^{(n)}} \right]$$

$$= \frac{-1}{\left[1 + e^{-\bar{w} \bar{x}^{(n)}} \right]^2} \cdot (-w_i) e^{-\bar{w} \bar{x}^{(n)}}$$

$$= x_i^{(n)} \cdot \frac{e^{-\bar{w} \bar{x}^{(n)}}}{\left(1 + e^{-\bar{w} \bar{x}^{(n)}} \right)^2} = x_i^{(n)} y^{(n)} (1 - y^{(n)})$$

$$\frac{\delta h(w)}{\delta w_j} = - \sum_n (t^{(n)} - y^{(n)}) x_j^{(n)} = - \sum_n e^{(n)} x_j^{(n)}$$

$$\left(\frac{\delta h}{\delta w_0}, \dots, \frac{\delta h}{\delta w_I} \right) = \underbrace{\vec{g} = - \sum_n e^{(n)} \bar{x}^{(n)}}_{\text{gradient vector}}$$

gradient vector

Back propagation

Update weights by a quantity η in the opposite direction to the gradient

$$\begin{aligned}\bar{w}^{(old)} \\ \bar{w}^{(new)} &= \bar{w}^{(old)} + \eta \sum_n e^{(n)} \bar{x}^{(n)} \\ &= \bar{w}^{(old)} + \eta \sum_n (t^{(n)} - \gamma(\bar{w}^{old}, \bar{x})) \bar{x}^{(n)}\end{aligned}$$

η = learning rate " a free parameter

Different ways to update the weights

→ batch gradient descent learning
→ on-line gradient descent learning

batch learning

all weights are updated by looking at all data points at the time

- $\bar{w}^{(0)}$

- $\bar{w}^{(1)} = \bar{w}^{(0)} + \eta \sum_n (t^{(n)} - y_0^{(n)}) \bar{x}^{(n)}$

$$y_0^{(n)} = y^{(n)}(\bar{w}^{(0)}, \bar{x}^{(n)})$$

On-line learning

change all the weights by looking at one data point at the time

- $\bar{w}^{(0)}$

- take $m \in M$

$$\bar{w}^{(1)} = \bar{w}^{(0)} + \eta \left[t^{(m)} - y^{(m)}(\bar{w}^{(0)}, \bar{x}^{(m)}) \right] \bar{x}^{(m)}$$

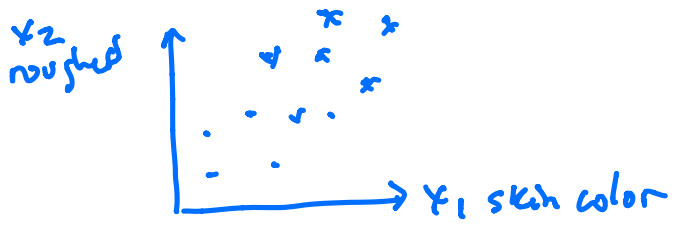
How well does the batch learning algorithm do?

Apples/oranges

2 inputs	x^1 - skin color	w_0
	x^2 - surface roughness	w_1
		w_2

$$y(\bar{w}, \bar{x}) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$

$N=10$ data points



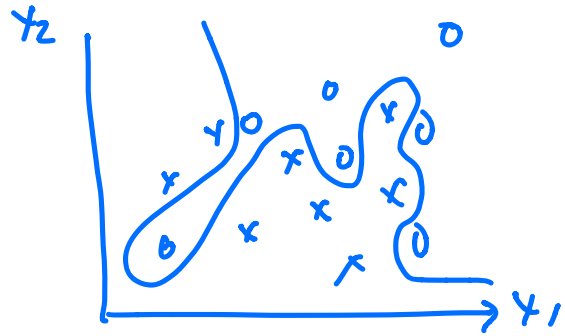
$\eta = 0.01$ $w_0^{(0)} = w_1^{(0)} = w_2^{(0)} = 0$



→ class code

Regularization

To avoid overfitting



$\mathcal{L}(w)$

$$M(\bar{w}) = \mathcal{L}(\bar{w}, \bar{x}) + \alpha R(\bar{w})$$

$$R(\bar{w}) = \frac{1}{2} \sum_i w_i^2$$

$$\frac{\delta M}{\delta \bar{w}_i} = \frac{\delta \mathcal{L}}{\delta w_i} + \frac{\delta R}{\delta w_i} = \frac{\delta \mathcal{L}}{\delta w_i} + \alpha w_i$$

$$\bar{w}^{(new)} = \bar{w}^{(old)} (1 - \eta \alpha) + \eta \sum_n (t^n - y_i^{(n)}) \bar{x}^{(n)}$$