

How Many Hidden Layers?

(1)

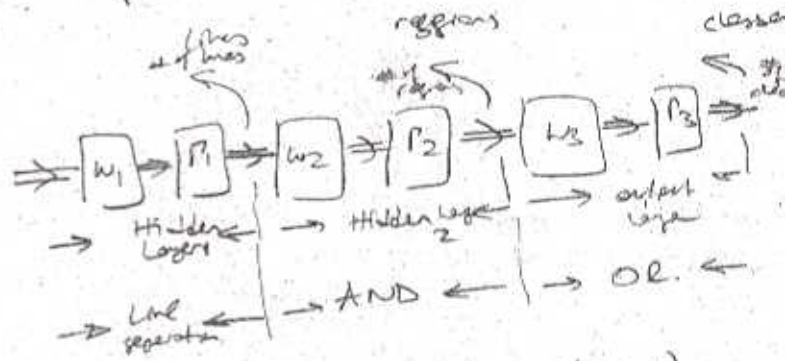
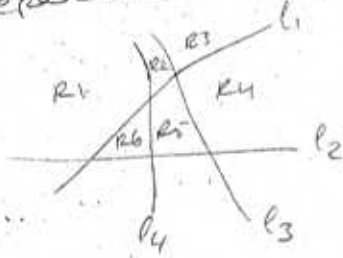
1) Too many layers slow down the convergence.

Too few may result in divergence.

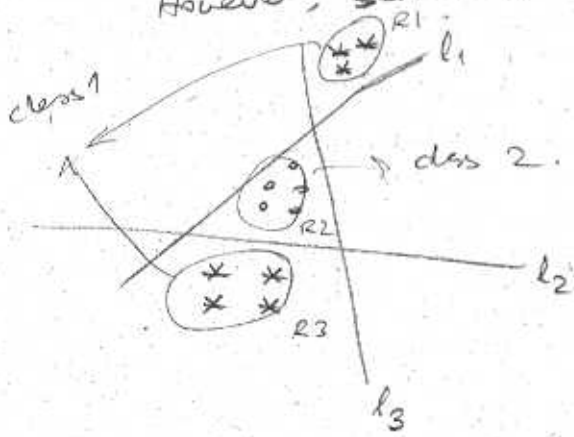
(CLASSIFICATION MILESTONE)

We know that at most 3 layers (2 hidden + one output) is sufficient to classify linearly

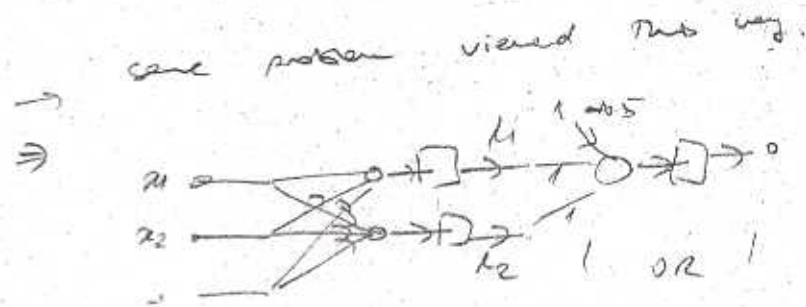
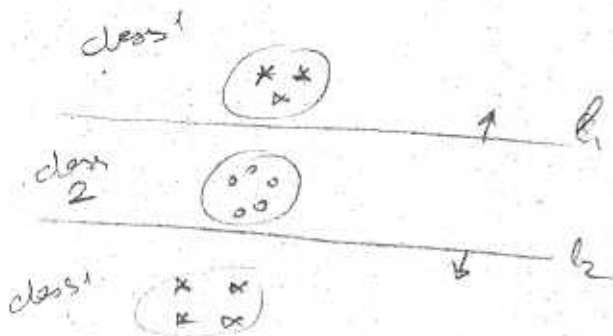
separable classes



So in this approach 2 hidden layers (+1 output layer) is sufficient and safe for many applications. However, sometimes one of the layers may not be necessary.



⇒ In this case region separation (AND) and union (OR) blocks are necessary.



$h_1=1$ OR $h_2=1$ ⇒ class 1

(AND block is not necessary!)

FUNCTION APPROXIMATION VIEWPOINT

One hidden layer + one output layer is sufficient for approximating any continuous function.

Hence 2 hidden layers should be used at next.

For approximation, 1 hidden layer is sufficient.

Advantage of one hidden layer: Simple structure, but may not be sufficient for some problem.

Disadvantage: Too many neurons may be required. That may slow down the convergence.

Advantage of no hidden layer: sufficient for many applications.

total # of weights may be less than single layer one.

Disadvantage: slightly more complex structure.

RECOMMENDATION: May start with single layer. If error cannot be reduced below a prescribed level, add another layer.

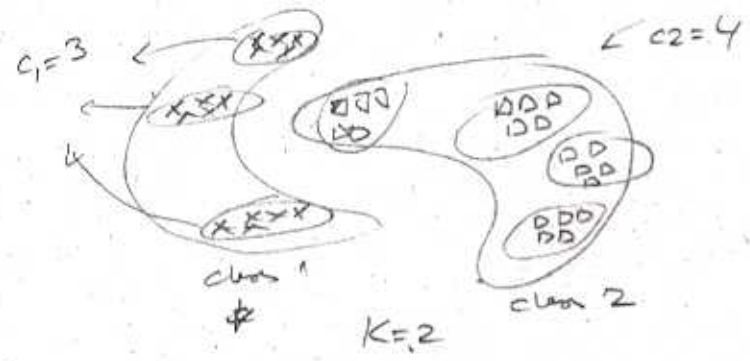
OR: For simple-looking problems (function approximation, classification with few classes, use single layer)

For seemingly complex problems (classification with large number of classes) use two hidden layers.

How many neurons in single hidden layer

K : Number of distinct classes

Each class is assumed to be linearly separable C_i subregions.
 total linearly separable subregions: $C = C_1 + C_2 + \dots + C_K$



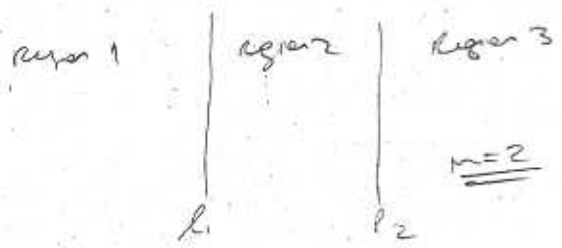
$K=2$ (# of classes)
 $C = C_1 + C_2 = 7$ (# of regions)

Each neuron output separates the input space into 2 regions.

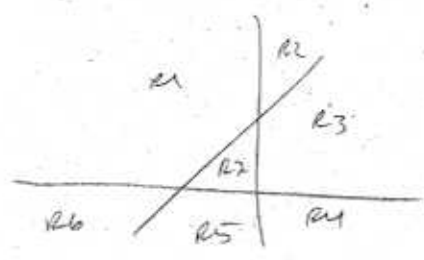
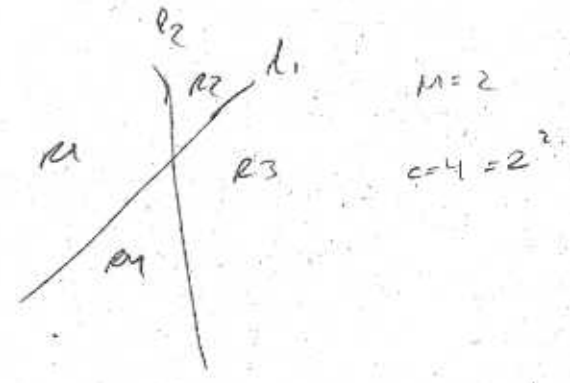
$\rightarrow M$ neurons : at most : 2^M regions
 at least : $M+1$ regions.

At most: Each line cuts the previous spaces into 2.

At least: Lines are parallel.



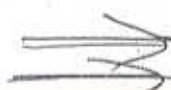
Total regions: $3 = 2 + 1$



$M=3$
 $C=7 < 2^3=8$

\rightarrow No added line does not cut R_2 into 2!

In general: upper bound may not be achieved (It is almost impossible for a newly added line to cut all of the previous regions into 2)



(pessimistic)
(least regions)

(optimistic)
least regions

$$M+1 \leq C \leq 2^M$$

$$\Rightarrow \boxed{M < C \leq 1}$$

$$\boxed{M > \frac{\ln c}{\ln 2}}$$

~~$\frac{1}{2}$~~
Take average!

$$\boxed{M_{ave} = \frac{1}{2} \left(C-1 + \frac{\ln c}{\ln 2} \right)}$$

For large c $\begin{matrix} c \gg \ln c \\ c \gg 1 \end{matrix} \rightarrow \boxed{M_{ave} \approx \frac{1}{2} C}$

(if we knew c !)

In general, we know K , but we hardly know C .
One assumption, on the average, each class K contains 2 subclasses that are linearly separable.

⇒ An estimate for $C \Rightarrow \boxed{C = 2K}$

⇒ Use $\boxed{M = \frac{\ln 2K}{\ln 2}}$ neurons.

For highly (complex) nonlinear cases (?) we may use

$\boxed{C = pK}$ $2 \leq p \leq 4$

⇒ $\boxed{M = \frac{\ln pK}{\ln 2}}$ $2 \leq p \leq 4$

→ A realistic estimate for regions:

The upper bound on the number of regions is

$C \leq 2^M$

This is highly unrealistic (because we assume that a new neuron ⇒ a new separating hyperplane

⇒ this plane cuts all of the previous regions into two.

Assume

$C = r^M$

$1 < r < 2 \Rightarrow$ use $r = 3/2$

→ $C = \left(\frac{3}{2}\right)^M$

→ $\boxed{M = \frac{\ln C}{\ln 3/2}}$

For c , use $c=2K$ as before

$$\rightarrow \boxed{M_{\text{ave}} = \frac{\ln 2K}{\ln 3/2}} \quad ! \quad \approx \underline{\underline{pk}}$$

For two hidden layers, the previous heuristic will not work. Instead use the following:

L : # of neurons in the first hidden layer
 M : # of " " " " second " "

$$L = \sqrt{2K} + p$$

$$M = \sqrt{2K} - p$$

$$\boxed{2 \leq p \leq \frac{\sqrt{2K}}{2}}$$

select p such that $L \approx 2M$. (eg. choose $\boxed{p \approx \frac{\sqrt{2K}}{3}}$)

FOR OUTPUT LAYER.

J : # of neurons in output layer

For small K , the best choice is $\boxed{J=K}$

($K \leq 16$). For $K > 16$, use $\boxed{J = \log_2 K = \frac{\ln K}{\ln 2}}$

(assume binary code at the output)

The first choice represents the maximum selection
The second " " " " binary coding.

Hypothetical Design:

Design a MLP to recognize a set of 100 symbols.
Each pattern is represented with a code length 60.
($x \in \mathbb{R}^{60}$)

$K=100$

Single layer case: $C=2K=200!$

Size of the hidden layer: The permutative case:

$M \leq 2K-1=199$

Optimistic case: $M = \frac{\ln 2K}{\ln 3/2} \approx 12$

take average \Rightarrow $M_{ave} = \frac{12+199}{2} = 105$

output layer: $K=100 > 16 \Rightarrow$ use binary code
 $(2^7=128 > 100, 2^6=64, \text{ not sufficient})$

$J=7$

Two layer case:

$\sqrt{2K} \approx 14, \frac{2K}{3} \approx 4 \Rightarrow$

$L=14+4=18$
 $M=14-4=10$

\Rightarrow # of neurons in the first layer
 \rightarrow # of neurons in the second layer

$J=7$

Total weights:

$$60 \times 12 + 12 \times 7 = 804$$

(the rest apartment)

$$60 \times 199 + 199 \times 7 = 13333$$

(the rest apartment) - 13333

$$60 \times 185 + 185 \times 7 = 7035$$

→ average) ~ 735
= 7035

The layer are:

Total weights:

$$60 \times 18 + 18 \times 10 + 10 \times 7 = 1330$$

(reasonable).