

## 12 - Neural net as a universal approx.

- BR: neural networks can represent almost anything\*

- NN black box I/O device; content of boxes

- NN originally conceived of emulations/caricatures of 'neurons'

- ANN: perceptron as basic unit - threshold unit

- rewrite a linear combination of inputs with threshold  $\rightarrow$  affine combination

$$\text{from } z = \sum_i w_i x_i - T \rightarrow y = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{else} \end{cases}$$

- Affine (think of as containing intercept)

- replace threshold unit with sigmoid activation 'squashing function'

- W continue assuming activation <sup>function</sup> is  
a threshold unit for this lecture.

- bounded on both sides
- monotonic

- activation function - function that affine combination is put through

- original perceptron was formulated using threshold units

- MLP - network of perceptions (units)

- directed I/O; layered (hidden layers')

- layer - a way of referring to sub-networks of neurons / units that  
do not communicate

- "depth" - what does it mean for a graph to be 'deep'

- sense of universal approximation in Hornik?

BR: - length of longest path in a directed graph  
from source  $\rightarrow$  sink (graph theoretic specification)

- depth  $> 2 \Rightarrow$  'deep'

or  $\geq 3$

- one initial input layer, output layer  $\rightarrow$  not deep

- MLP: inputs: real or Boolean

Output: ...

- MLP can approximate any Boolean function; any real-valued function  
with imitations

4 topics → see slides (MLP)

- classify/organise lecture this way

1. MLPs as universal Boolean functions

LH: AND RH: NOT

universal AND gate

values in circles - thresholds

Bottom: OR

(W): (A) = universal AND (logic)

values edges - weights

- Perceptron also as universal OR (W): (A) - logic

- generalised majority gate

- These 'nearly fooled' Rosenblatt

- (con) Single perceptron cannot compute XOR; but MLPs e.g. 3 perceptrons

- Emulation of XOR with 3 perceptrons (2 hidden perceptions)

- Perceptrons → Boolean gates → Boolean functions.

(Q): MLPs are universal Boolean functions'

(Q): How many layers needed? BR: only 1 hidden layer

- Boolean function is a truth table ultimately

- truth table - 6 conditions for which output is 1

- disjunctive normal formula for Boolean function → 7 disjunctive terms/clauses

- + ~~perceptron~~ <sup>Heuristic</sup> illustration of DNF - 1 hidden layer (see slides)

- Any truth table

(Q): A 1-hidden-layer MLP is a universal Boolean function.

Karnaugh Map

- 7 blocks - 7 clauses

{ reduced DNF }

- get reduced DNF  $O = \bar{Y}\bar{Z} + \bar{W}\bar{X}\bar{Y} + \bar{X}Y\bar{Z}$

{ 3 clauses }

∴ MLP - 3 layers

(Q): (A) - clarify logic of Karnaugh maps

- largest irreducible DNF - see Karnaugh patterns

- largest no. of neurons for which we have irreducible DNF - 8 red squares: 8 neurons

- intermediate cube

- generalisation: largest no. perceptions required for single layer MLP for an N-input variable function:-

(61)

$$2^{N-1} \text{ (exponentially large)}$$

- now about multiple hidden layers?

- 1 XOR - 3 perceptions (in fact 2)

using this logic

m	x	y	z
0	0	0	0
1	0	0	1
2	0	1	0
3	0	1	1
4	1	0	0
5	1	0	1
6	1	1	0
7	1	1	1

$$0 = W \oplus X \oplus Y \oplus Z$$

3XORS ; 3 perception / XOR  $\rightarrow$  9 neurons / perceptions.

- checkerboard pattern is  
XOR  
 $0 = W \oplus X \oplus Y \oplus Z$  (check square)  
- check cube  
(see slide)

(ii): more generally, XOR of N variables requires  $3(N-1)$  perceptions (in deep)

(iii): for 1 hidden layer:  $2^{N-1} + 1$  perceptions (including output) : exponential in N

For deep-network:  $3(N-1)$  perceptions - linear in N!

arranged in  $2 \log_2(N)$  layers ?

(iv) (v): Relation between depth-hidden layers

(vi): depth-perceptions logic

(vii): optimal depth is function of no. inputs

(viii): No. parameters  $\rightarrow$  no. of connections 32 connections

(A6)

- this matters

- exponential no. neurons  $\rightarrow$  superexponential no. of parameters/weights

(A7)

- read slides before lectures

- there are formal results here  $\rightarrow$  XOR as parity problem (A8)  
first et al.

shannon et al.

- network size summary  $\rightarrow$  see course readings ✓

(vi) points here (heuristic observations)

## caveats 2

- threshold gates subst. for Boolean circuits

2. MLPs as universal classifiers over real inputs

- decision boundary classifies  $n$  input patterns in 784 dimensional space

-  $(n-1)$  dimensional hyperplane? (~)

- OR, NOT, XOR

- composition of Boolean perceptrons → complex decision boundaries

- Q: Booleans over reals → pentagon

- And then further compositions (composition of pentagons)

- via fragmentation; compose using subsets

Q: MLP can compose arbitrarily complex decision boundaries

BR: with 1 hidden layer

Q: do you understand the heuristic arguments

- some kind of limit behavior

Q: I don't understand what he is trying to illustrate here

going from pentagon → cylinder

Q: illustration of arbitrarily approximating the two pentagons  
using limiting behavior (universal approx.)

Q: some shapes can be modelled with 1 hidden layer MLP exactly;  
others can use infinitely deep hidden layers (double pentagon approx.) as approx.

Q: need to clarify the claims precisely

optimal depth → see papers

(tradeoff of depth (hidden layers) vs no. of neurons)

- checkerboard → heuristic understanding → dig(1)

Q: the difference in size between the optimal XOR net and  
shallow net increases (with no. neurons) increases with <sup>depth of</sup> pattern  
<sup>complexity</sup>

Q:

How quickly can  
you adapt your methods  
to course

①: deep networks are more expressive

### 3. MLPs as universal approximators

- use 'cylinders' with volume (some kind of calculus/analysis arguments)
- can exploit ranges of activation functions

②: network as a universal map from entire domain of input values to entire range of output activation

- #1  
optimal depth vs width

- sufficiency of architecture

- there are some requisite depth/breadth constraints (capacity)
- not all architectures can represent any function

- after information has propagated through network for 8 (rather than 16) units in output layer; there is some loss in the sense that we can no longer make distinctions within the diamonds

③: there are significant constraints on architecture:-

1. network must be sufficiently wide or deep at each layer

- ④ (AB) : ③ - information propagation

- get his logic (\*) - a lot missed; rework until you understand
- provided you have appropriate activations to capture info missed by lower layers?

width vs activations vs depth

- sufficiency of architecture

- capacity/expressive power has many definitions

- info/storage capacity

- VC dimension (complexity of patterns): bounded by square no. of weights

VC dimension (capacity): See papers  $\emptyset \rightarrow [E]$

## Highly abstract results

The function we are trying to model:- remember:-

1. width of individual layers
2. depth of network

3. Activation functions used -

(\*) For any given activation function, they limit extent of tradeoff.

Next class: now do we set parameters, actually training an MLP