



THEORY OF MACHINE LEARNING

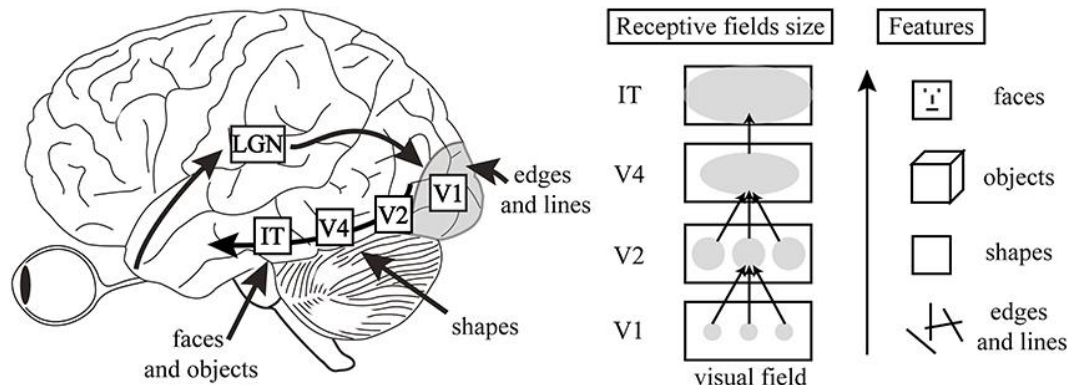
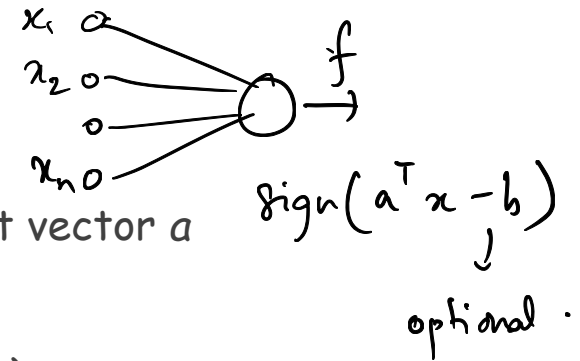
LECTURE 18

NEURAL NETWORKS – REPRESENTATION BASICS



RECAP

- Perceptron, or linear threshold
 - Hypotheses of form $\text{sign}(\langle a, x \rangle)$ for an appropriate weight vector a
 - Generally, $\sigma(a^T x)$ for some "activation function" σ
- Biologically inspired, arithmetic circuit (with threshold gate)
- Idea behind neural nets:
 - Perceptrons detect "basic" or "primitive" features; 'composing' them allows for complex decision-making (can get high level concepts using basic features)
 - Supported by human visual system (V1, V2, ...)



$\sim 6^7$ layers.

RECAP: ARTIFICIAL/DEEP NEURAL NETWORK (DNN)

- **Definition.** A layered "circuit" that takes a vector of input features x , produces output $y = F_r \circ F_{r-1} \circ \dots \circ F_1(x)$, where each F_i is a function of the form $F_i(z) = \sigma(Az + b)$, for some activation function $\sigma()$ (that acts coordinate-wise)

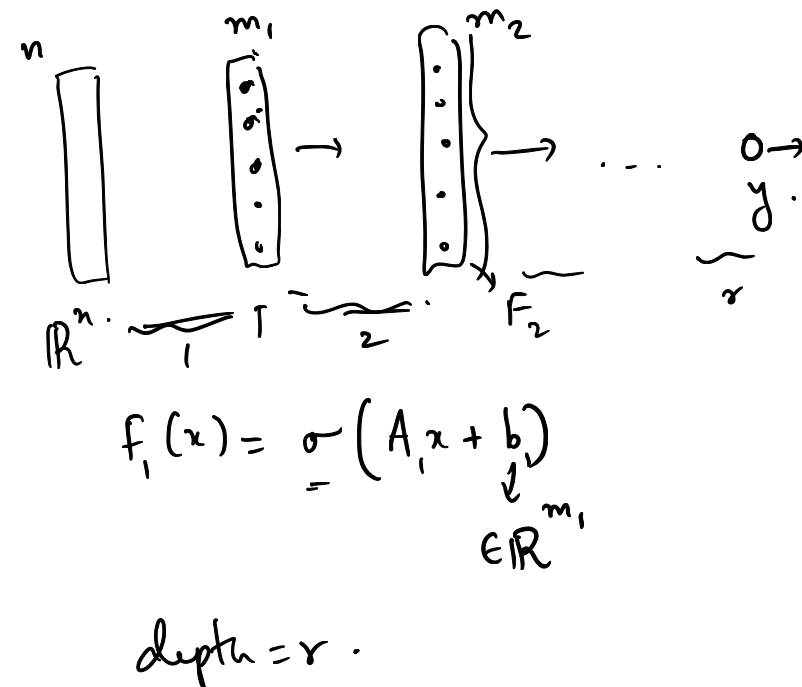
- Common activation functions:

- Threshold

- Sigmoid: (continuous approx.) $\frac{1}{1+e^{-x}}$

- ReLU, Tanh

- ...



BASICS

linear classifiers + + .

- Neural networks are basically a (fairly complex) hypothesis class - takes input x , produces y

- Question (vanilla supervised learning):** given data $(\underbrace{x_1}_{\substack{\text{feature vector} \\ \in \mathbb{R}^n}}, \underbrace{y_1}_{\text{label} \in \mathbb{R}}), (x_2, y_2), \dots$ from some distribution D , find h in this class that minimizes the risk

$$\text{risk}(h) = \mathbb{E}_{x \sim D} \mathbb{1}[h(x) \neq y] \quad \rightsquigarrow \quad \mathbb{E}_{(x,y) \sim D} l(h(x), y) .$$

\downarrow label(x) \downarrow real valued y

- ERM problem usually called neural network "training" - given data, find best hypothesis ($f(x_i) = y_i$) for all i

ERM: Minimize training loss
(using depth, width of network bounded appropriately).

THEORY OF DEEP LEARNING

- Expressibility (how rich is this hypothesis class?)
 - What kinds of functions can be obtained using a DNN? \neq we .
- Training complexity & training dynamics for GD and variants
 - Can the ERM problem be solved efficiently? What guarantees are possible?
- n \neq we .
- Generalization
 - What kind of generalization bounds can we prove? (VC dimension?) $\rightarrow \approx \# \text{parameters}$
 \neq n \neq we .

Key: "easy" answers for all questions, but unsatisfactory for realistic settings

EXPRESSIBILITY BASICS

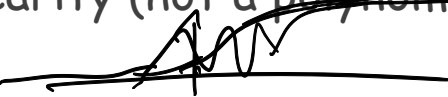
- **Barron's theorem [93]**. Any continuous function f that satisfies an appropriate "niceness" condition (parametrized by C) can be approximated to error ϵ (in L^2 !) by a 2-layer NN with $\sim \frac{C^2}{\epsilon}$ internal nodes

$$\int x |\hat{f}(x)| dx \leq C.$$

- (Nice functions can be approximated by small NNs) (2 layer)

- **Universal approximation [Cybenko, Hornik '87, '91]**. Any continuous function (over a compact domain) can be approximated by a 2-layer NN with any non-linearity (not a polynomial)

"Sigmoidal" non-lin:



(pointwise)



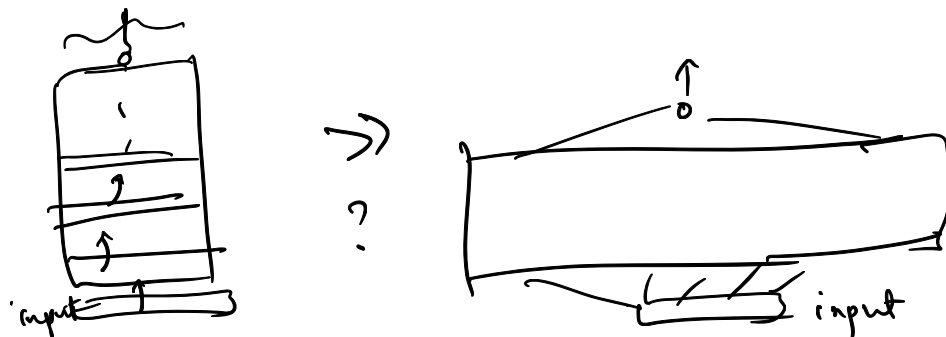
But wait.. who uses infinitely wide 2 layer nets?

→ input dim = $n \Rightarrow$ width needed for "basic" fns

is $\sim 2^n \cdot (\exp(n))$



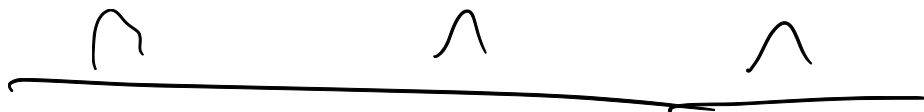
DEPTH VERSUS WIDTH



■ Practical intuition:

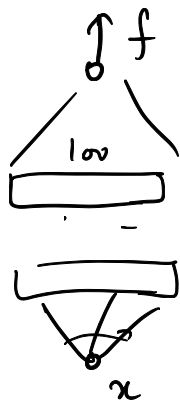
- Depth allows "meaningful features" while width is for "brute force memorization"
- Universality results degrade rapidly with dimensions
 - Curse of dimensionality (width ~~needs to~~ ^{needs to} be exp in dimension). ^{Universality: Any fn is expressible using depth-2 nets.}
 - Modern nets work with high dimensional data
- Does higher depth lead to higher expressibility (with much fewer neurons)?
- Bunch of works ... [Eldan and Shamir (depth 2 vs depth 3)], [Telgarsky], 2015-16

620

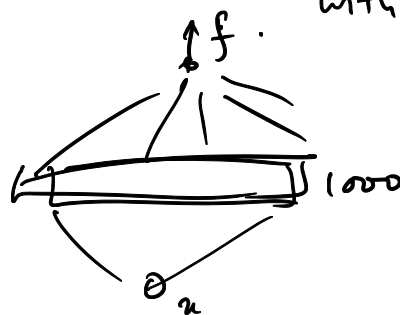


(fix input dim = 1)

$F_1 : \left\{ f : f \text{ is the output of a depth 3 NN with width } \leq 100 \right\}$



$F_2 : \left\{ f : f \text{ is the output of a depth 2 NN with width } \leq 10^8 \right\}$



Qn: Can any fn in F_1 be approximated to error $\leq \epsilon$ using functions in F_2 ?



→ CIRCUIT LOWER BOUNDS: Look at specific f of interest, and

(boolean)
ask if \exists a circuit of size $\leq S$ that computes f ?

[Razborov '85]
[Minsky Pappert '69]

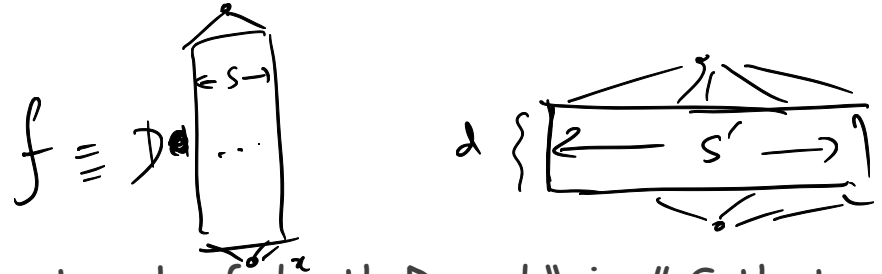
parity function requires "exponential width"

[Håstad '86]

4 inputs $\begin{matrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 \end{matrix} \rightarrow 1$
 $\begin{matrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 \end{matrix} \rightarrow 0$

fn: $\{0,1\}^n \rightarrow \{0,1\}$

POWER OF DEPTH



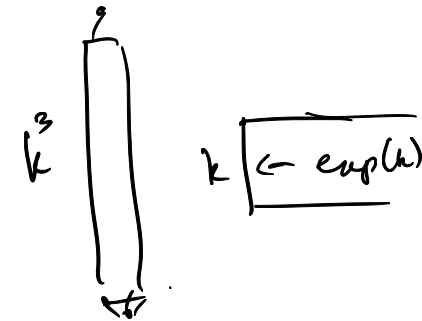
Theorem template. There exists a network of depth D and "size" S that computes some function f that cannot be approximated by the output of any network with depth d and size S' (typically if $d \ll D$, S' will be $\gg S$)

- "Depth versus width" results
- Reminiscent of circuit complexity (original work of Minsky, Pappert)

[Telgarsky 16]. For any $k > 0$, theorem holds with:

$D = S \sim k^3$, $d = k$, and $S' = 2^k$ (and ReLU activations)

$D = k^3$; width $\sim \Theta(1)$ s.t. approximating with depth k requires $\exp(k)$ width.



PROOF OUTLINE

- Consider just one-dimensional inputs and ReLU activations

- **Key insight:**

- depth D lets us achieve $\exp(D)$ many "oscillations" in f
- getting so many oscillations with depth d requires huge width!

↪ width w and depth $= d$, $\# \text{oscillations} \leq$ $\Theta\left(\frac{w^{o(d)}}{d}\right)$

n, d
 $n^{\log d}$
 $2^{\log n}$

$$D = k^2$$

$$; d = k.$$

$$2^{k^2}$$

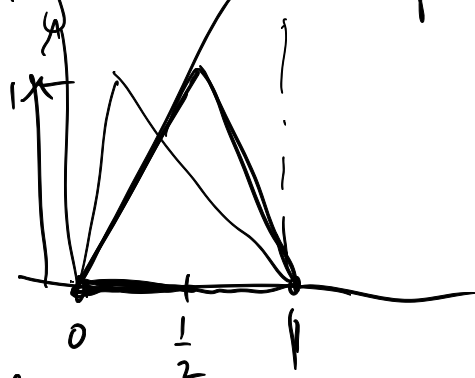
$$w^{10^4}$$

$$\Rightarrow w \geq \frac{2^{k^2}}{10}$$

$$e^{k^2} \Rightarrow (e^k)^k \Rightarrow (1000)^k \Rightarrow e^{(\log 1000)k}$$

Part 1:

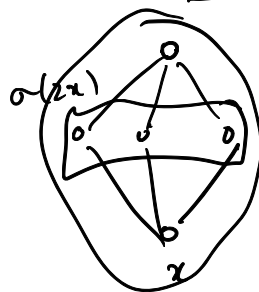
→ Single variable input x : domain = $[0, 1]$



$f(x)$

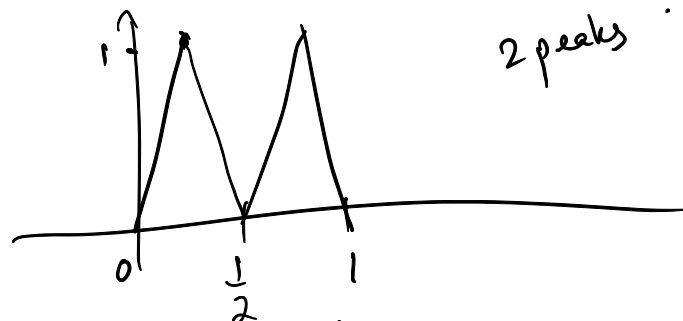
σ : ReLU: ~~max~~ $\sigma(x) = \max(0, x)$

$$\sigma(2x) - \sigma(4(x - \frac{1}{2})) + \sigma(4(x - 1))$$

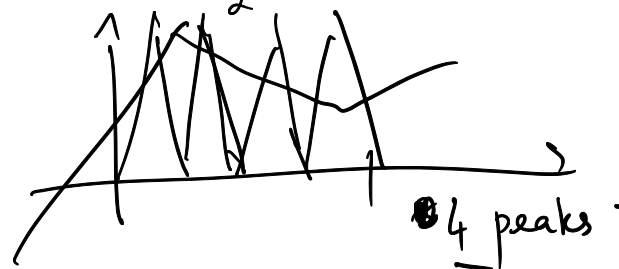


What is Φ ?

→ $f(f(x))$?



→ $f(f(f(x)))$ →



$f^{o(k+1)}(x) \rightarrow 2^k$ peaks.

depth $\sim 2k$ NW with width ≤ 3 , you can implement this..

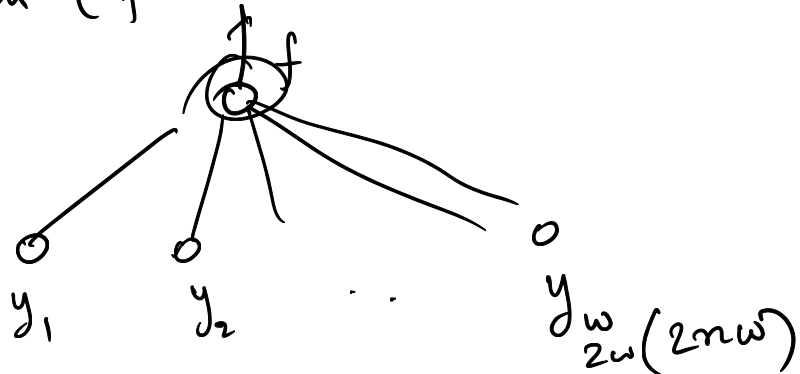
Part 2: Want to show that if $\text{depth} \leq k$, getting this fn requires huge width.

Obsn:

y_1 is a neuron that is a piecewise lin. function (of variable x) with $\leq n_1$ pieces.



Idea:



f is piecewise linear & # pieces in f is $\leq (n_1 + n_2 + \dots + n_w)^2$

→ Part 3: If f_1 and f_2 are p.w.l. with n_1 and n_2 pieces and $n_1 < \frac{n_2}{2}$; then $\|f_1 - f_2\|_1 \geq \frac{1}{8}$.