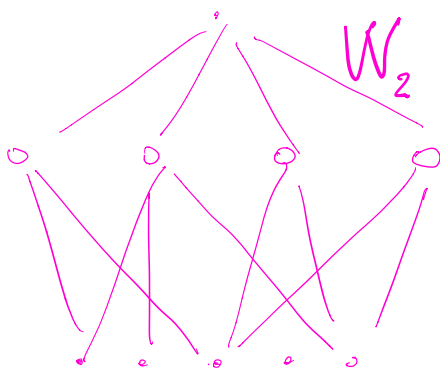


Why deep nets ?

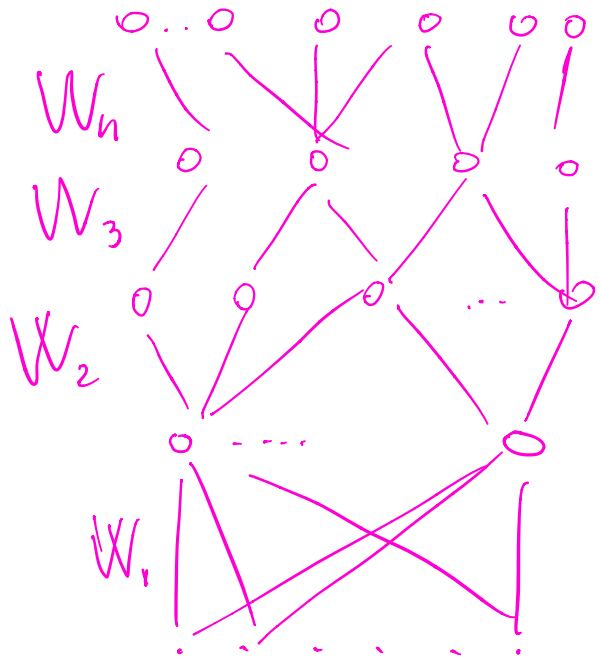
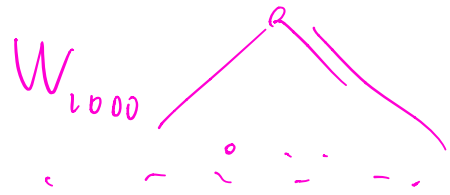
Is deep better than shallow ?

When ?

Why ?



Shallow



Deep

## Some history

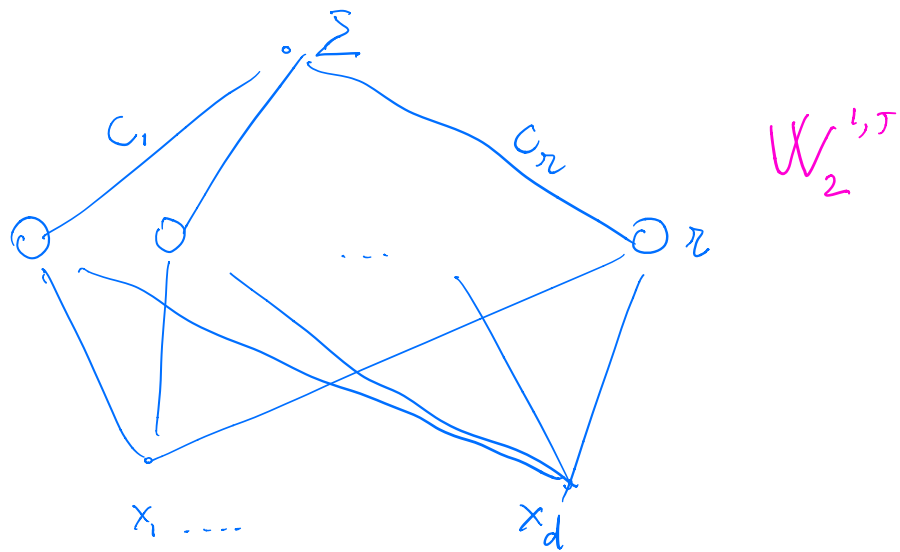
80': S.N.N.

2018:

# Approximation Theory

why is depth better?

# Shallow networks

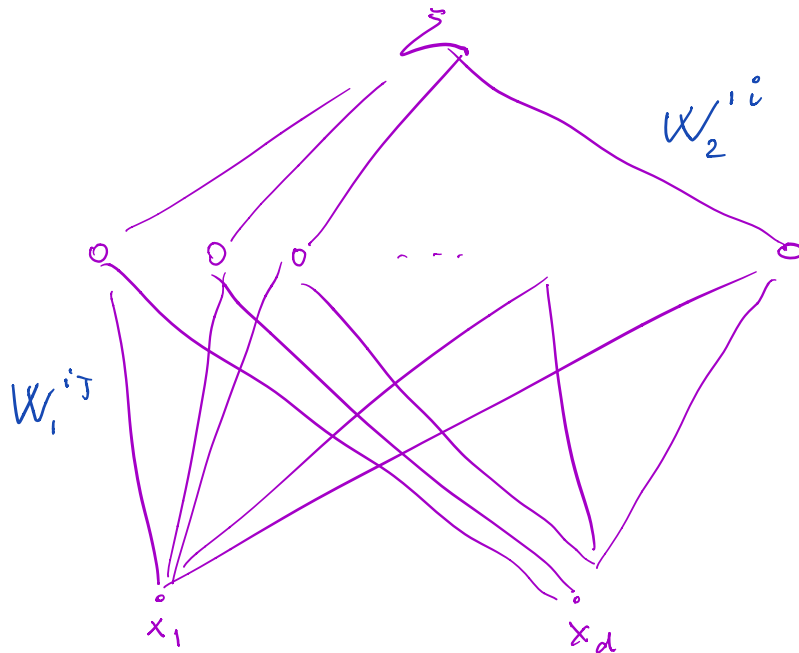


Main example: Kernel machine

$$y = \sum_{i=1}^n c_i K(x, x_i)$$

$$\min_{f \in H_n} \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_n^2$$

Another example

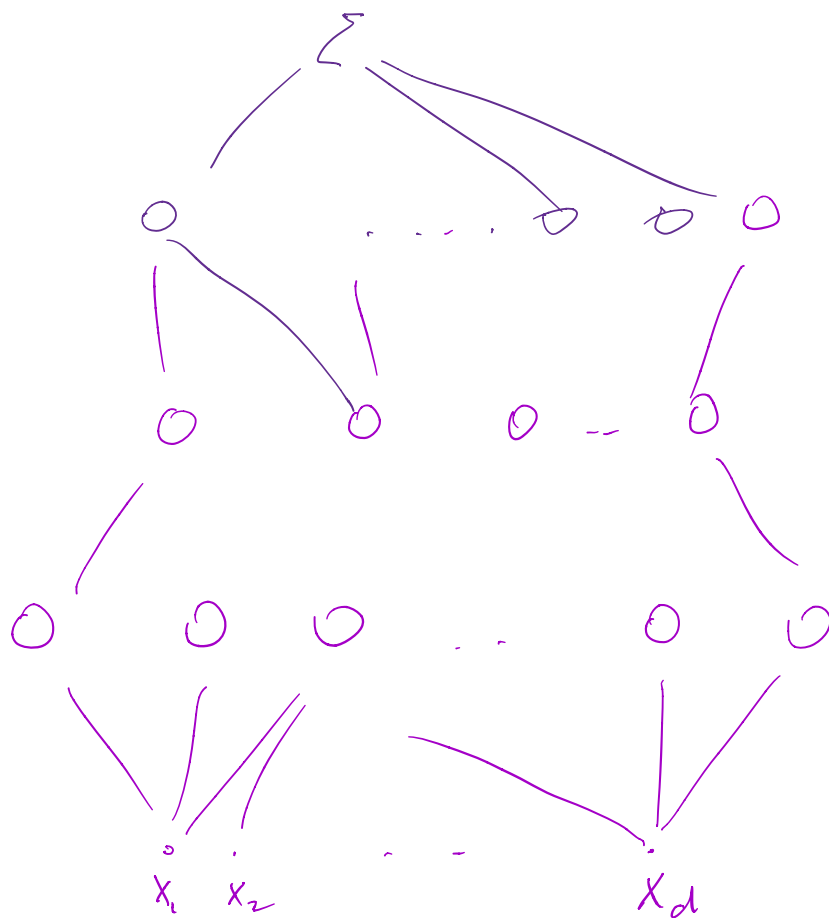


$$\delta(z) = [z]_+ = \max(0, z) = \begin{array}{c} | \\ \hline z \end{array}$$

$$y = \sum_i w_2^{i1} \delta\left(\sum_j w_1^{ij} x_j\right)$$

usually  $\delta(w \cdot x + b)$  but I eliminate  $b$  by  
 assuming one of the components is  $x_d = 1$  so  
 $w^{iJ} = b_i$

# Deep Networks



$$y = W_L^J \sigma \left( W_{L-1}^{JK} \sigma \left( W^{kP} x_P \right) \right)$$

E. summation convention

# Networks to approximate/ represent functions

- Are deep nets better than shallow ones?
- The answer in the 80' was : no!

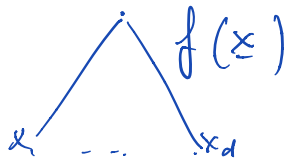
We will see

- ① proof of above
- ② a new answer : deep nets can be much better for certain  $f$ .



# Key ideas in approximation theory

## Functions and approximators



Example  $\forall \epsilon \exists \delta \text{ s.t. } |f(x+\delta) - f(x)| < \epsilon$   
 $f \in C(\mathbb{R}^d)$

$$g \in V_n \rightarrow \text{span}_n \delta(w^i x^T)$$

# Density

$\forall f \in C(\mathbb{R}^d) \forall \text{compact } K,$

$K \in \mathbb{R}^d$

$\forall \varepsilon > 0$

$$\exists g \in V \text{ s.t. } \sup_{x \in K} |f(x) - g(x)| < \varepsilon$$

set of

networks

networks

# Degree of approximation

$$\forall f \in C(\mathbb{R}^d)$$

$$\inf_{g \in V_n} |f - g_\alpha| < \varepsilon \implies$$

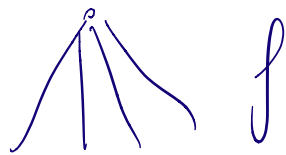
how large is  $\alpha$ ?

# Shallow nets:

## density + degree approximation

Consider target functions

$$f \in W_m^d \quad \|f\| = \sum_{i \in [k] \subseteq m} \|D^k f\| \leq 1$$



no structure assumed

Theorem

$$g(x) = \sum_{i=1}^N c_i \langle \underline{w}^i, \underline{x} \rangle$$

$$\forall f \in W_m^d, \exists g \in V_N \quad \text{s.t.}$$

$$|g(x) - f(x)| < \varepsilon \quad \text{with}$$

$$N = O(\varepsilon^{-d/m})$$

# Curse of dimensionality

Bellman's term:



□ optimization  
cannot be done by  
RS

□ function approximation  
requires  $(\frac{1}{\epsilon})^d$  evaluations  
for  $f$  Lipschitz order  $\epsilon$

□ integration ...

# Blessings of

1) Smoothness

Barron's Gurosi

2) compositionality

Examples :

$$\textcircled{1} P_k^d \text{ has } \binom{d+k}{k} = \frac{d+k!}{d! k!} =$$

$\textcircled{2} K^d$  monomials

$\textcircled{3}$  A function of 10 variables corresponds to a 10D table. If each dimension is discretized in just 10 partitions I have table with  $10^{10}$  entries. If  $d = 100$  pixels

then  $10^{100}$  entries

$$\textcircled{a} \text{ If } f \in W_m^d \rightarrow N = O\left(\varepsilon^{-\frac{d}{m}}\right)$$

$$\text{For } \varepsilon = 10^{-1} \quad d = 100 \quad N = O(10^{100})$$



# Summary Proof

$$\textcircled{1} \sum_i^k c_i \delta(a_i x + b_i) \approx p(x) \in \mathbb{T}_{k-1}$$

$$\textcircled{2} \sum_i^{\mathcal{N}} p_i^k(\langle W, x \rangle) \in P_k \text{ (total degree } k \text{ in } d \text{ variables)}$$

$$\mathcal{N} \sim k^d \Rightarrow k \sim \mathcal{N}^{1/d}$$

$$\textcircled{3} \text{ Sobolev } E(W_d^m, \mathbb{R}_{\mathcal{N}}^d, L_p) \leq C k^{-m}$$

*space of ridge functions  
with  $\mathcal{N}$  elements*

$$\textcircled{2} + \textcircled{3} \Rightarrow E \leq C \mathcal{N}^{-m/d}$$

$$\mathbb{R}_{\mathcal{N}}^d = \left\{ \sum_{\alpha_i \in \mathbb{R}^d} g_i(\langle \alpha_i, x \rangle) : g_i \in C(\mathbb{R}) \right\}$$

Logic of ① + ② + ③

- Networks approximate univariate poly
  - Univariate poly in  $(w, x)$  represent multivariate poly
  - Multivariate poly approximate Sobolev functions
- Thus theorem.

① Any univariate  $p(x)$   
 can be represented as linear  
 combination of smooth ReLU

proof

$$\lim_{h \rightarrow 0} \frac{\delta(a+h)x + b - \delta(ax+b)}{h} =$$

$$= \frac{d}{da} \delta(ax+b) \Big|_{a=0} = x \delta'(b)$$

Theorem

If  $\delta$  is not a polynomial,  
 $\delta \in C^\infty$ , the closure of  $N_r$  (span  
 of  $r \delta$ ) contains the linear  
 space of  $\mathbb{T}^{2-1}$

Second derivative which  
needs 3 terms gives  $x^2$

...

Thus  $N_{\mathbb{R}}$  is dense in  $C(H)$   
because of Weierstrass Theorem.

# From 1D to dD

$$P(x) = \sum_i^r P_i(\langle \underline{W}^i, \underline{x} \rangle) \in \mathbb{R}^d$$

$\ell_i$ : univariate  
 $\in H_{\kappa}^d$  (homogeneous pol in  $d$  variables of degree  $\kappa$ )

$$r = \dim H_{\kappa}^d = \binom{d+\kappa-1}{\kappa} = \frac{d+\kappa-1!}{(d-1! \kappa!)} =$$
$$\sim \kappa^d$$

o thus  $H_{\kappa}^d$  pol. can be represented by a network with  $r \sim \kappa^d$  units

We want to show that  
if  $P(x)$  pol on  $\mathbb{R}^d$  then

$$P(x) = \sum_i^r P_i(\langle W_i, \underline{x} \rangle)$$

for some choice of  $r$   $W_i$  and

$P_i$

No general proof but  
consider following:

assume network with units

s.t.

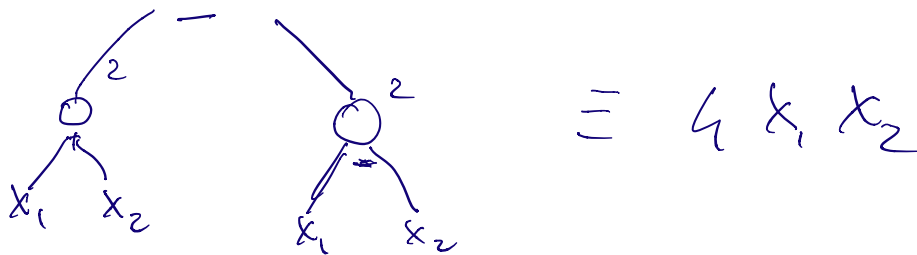


Can't synthesize  $P(x)$  of

in  $\mathbb{R}^d$  variables? I can

get  $x_1^2, x_2^2, \dots$

how do I get  $x_1 x_2$ ? Well



### Remark

$H_n^d$  hom pol degree  $n$  in  $\mathbb{R}^d$   
 $P_n^d$  " " " " "

$$\dim H_n^d = \binom{d-1+n}{n} = \mathcal{N}$$

$$\dim P_n^d = \binom{d+n}{n}$$

③ Define  $E(\mathcal{B}, X, L_p) =$   
 $= \inf_{P \in X} \|P - f\|_{L_p} \quad \forall f \in \mathcal{B}$

Theorem

$$E(W_d^m, N_r, L_p) \leq C r^{-\frac{m}{d-1}}$$

$$\uparrow \approx \sum_{i=1}^r \delta(\underline{W}^i, \underline{x})$$

proof

Classically  $E(W^m, P_k, L_p) \leq C k^{-m}$

Since  $r \sim k^{d-1} \Rightarrow k \sim r^{\frac{1}{d-1}} \Rightarrow$

$$E(W^m, N_r, L_p) \sim E(W^m, P_k, L_p) \sim$$

$$\sim C r^{-\frac{m}{d-1}}$$



# Remarks

1) Even without  $\textcircled{3}$  a shallow net can represent arb. well polynomials in  $\mathbb{F}_{\mathbb{R}}^d$  with  $n$  units  $n \sim H^d$

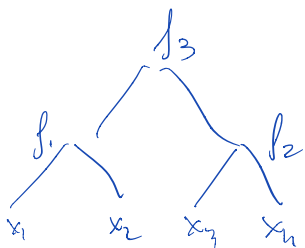
# Depth and curse

For general functions  
shallow and deep nets  
suffer curse of dimensionality

But... for Local (Hierarchical)  
Compositional functions  
deep nets - unlike shallow  
ones - do not have curse

# LHC functions

Simplest example



$$f(x_1, x_2, x_3, x_n) = \\ = f_3(f_1(x_1, x_2), f_2(x_3, x_n))$$

$$f \in W_m^{d,2} \text{ with } f_i \in W_m^2$$

Another example

$$f(x_1, x_2) = (A x_1 x_2 + B x_1 + C x_2)^2$$

shallow net require  $\sim 2^d$  units

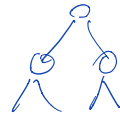
a deep net  $\sim 3 + 3 \cdot 10$

## Intuition



shallow net  $\left(\frac{1}{\epsilon}\right)^4$  units

deep net  $\left(\frac{1}{\epsilon}\right)^2$  for each node



for total  $\sim 3\left(\frac{1}{\epsilon}\right)^2$  units

## Another example

$$y = \sqrt{\sin(x_1 + x_2) \cdot (x_3 - x_4)^2} =$$
$$= h_6 \left( h_5 \left( h_2(x_1, x_2), h_3(x_3, x_4) \right) \right)$$

# Theorem

Deep nets with same graph approximate functions in  $d$  variables in  $W_m^{d,2}$  with  $\neq$  units per node  $\sim O\left(\frac{1}{\varepsilon}\right)^{2/m}$  for

total units  $O\left((d-1)\left(\frac{1}{\varepsilon}\right)^{2/m}\right)$

## Proof



Each  $h$  can be approximated with  $O\left(\varepsilon^{-2/m}\right)$  units. We assume

each  $h$  is Lipschitz continuous

that is  $|h(x) - h(x+\varepsilon)| \leq L\varepsilon$

By hypothesis  $\|h - P\| < \varepsilon$   $|h_1 - P_1| < \varepsilon$

$|h_2 - P_2| < \varepsilon$ . Then

$$|h(h_1, h_2) - P(P_1, P_2)| =$$

$$|h(h_1, h_2) - h(P_1, P_2) + h(P_1, P_2) - P(P_1, P_2)|$$

$$\leq |h - h| + |h - P| \leq \varepsilon + \varepsilon \sim O(\varepsilon)$$

Minkowski

Lipschitz

hypothesis

$$|f(x_1) - f(x_2)| \leq \kappa |x_1 - x_2|$$

= More general theorems for

DAGs



This theorem may explain why deep nets are successful and why all the really good ones are CNNs



Locality is key, not

weight sharing! weight

sharing helps but not exponentially