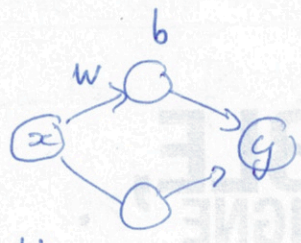
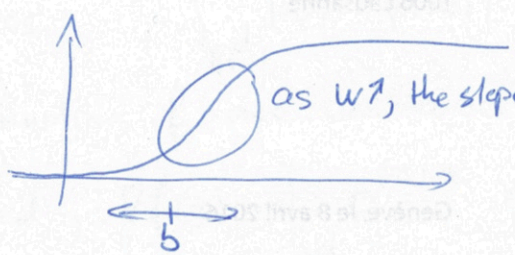


Prblm: A simple proof that a neural net can approximate any function  
 (↔ universality thm) (Michael Nielsen)

Take  $\phi(z) = \frac{1}{1+e^{-z}}$  eg.



• Top hidden neuron outputs  $\phi(wx+b)$ :



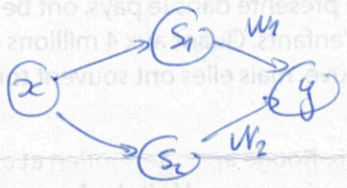
as  $w \uparrow$ , the slope increases here  $\Rightarrow$  step fn, where you want

$1 \text{ } x > s, s \approx -\frac{b}{w}$

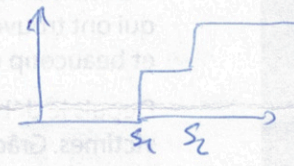
$\Rightarrow$  neuron idealization outputs  $1 \text{ } x > s$ ; single parameter  $s$

CNB: works for "every" activation fn  $\phi(x) \xrightarrow{x \rightarrow -\infty} 0, \phi(x) \xrightarrow{x \rightarrow +\infty} 1$

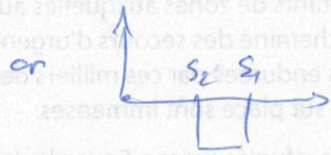
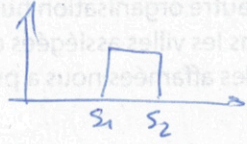
• Nav:



output  $y = w_1 1 \text{ } x > s_1 + w_2 1 \text{ } x > s_2$   
 (forget the last  $\phi \leftarrow$  bias  $b$ )

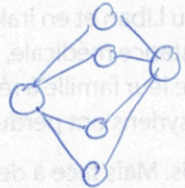


$w_2 = -w_1 \Rightarrow$

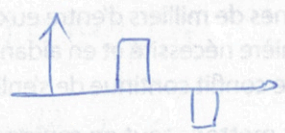


• adjustable height

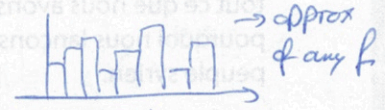
• Two pairs of neurons:



$\Rightarrow$  two bumps:



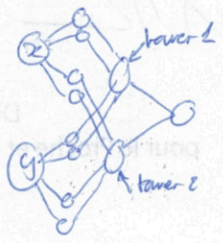
$\Rightarrow$  staircase with 2n neurons  $\leftrightarrow$  n intervals



(cf. Riemann sums!)

Generalization to multivariate  $x$ :

! two hidden layers:



or

! need  $h$  large & bias  $b$  ... (read again)

(?)

$\hookrightarrow$  one layer: circular towers!



enke: sigmoid functions can be used for approx: fine  
(as sines, polynomials)

na. Preliminary on Fourier transform

What kind of function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  will we be able to approximate?

Consider first  $\Gamma$ , the set of functions  $f$  s.t. its FT defined as

$$f(x) = \int_{\mathbb{R}^d} e^{i\omega \cdot x} \tilde{f}(\omega) d\omega \quad (*)$$

satisfies  $\int_{\mathbb{R}^d} |\omega| \cdot |\tilde{f}(\omega)| d\omega < \infty$

NB: integrability of  $\tilde{f}$   
 $\leftrightarrow$  regularity of  $f$

Such functions  $f$  are continuously differentiable  $\rightarrow$

(indeed,  $\nabla f(x) = \int_{\mathbb{R}^d} e^{i\omega \cdot x} i\omega \tilde{f}(\omega) d\omega$  is well defined)  $\& \|\nabla f(x)\| \leq C$

Then consider  $\Gamma_C$  for  $C > 0$ , the set of functions  $f$

such that  $C_f = \int_{\mathbb{R}^d} |\omega| \cdot |\tilde{f}(\omega)| d\omega \leq C$  [examples later?] sect. IX

NB: can be generalized to  $\tilde{f}(\omega) d\omega = \tilde{F}(d\omega)$  a measure, but we will skip this (bel)

Note: even if  $C_f$  is finite, (\*) might be an ill defined integral

skip

remedy to this:  $f(x) - f(0) = \int_{\mathbb{R}^d} (e^{i\omega \cdot x} - 1) \tilde{f}(\omega) d\omega$

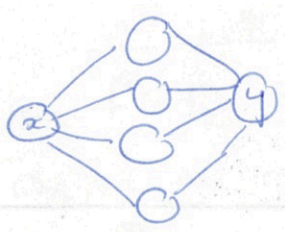
is well defined, as  $|e^{i\omega \cdot x} - 1| \leq |\omega \cdot x| \leq |\omega| \cdot |x|$

With what functions  $f_n$  will we approximate  $f$ ?

Let  $\phi$  be a fixed banded measurable fn on  $\mathbb{R}$  with  $\phi(z) \xrightarrow{z \rightarrow 0} 1$   
 $\phi(z) \xrightarrow{z \rightarrow \infty} 0$ .  
typically continuous!

Typically a cdf, even more typically a sigmoid:  $\phi(z) = \frac{1}{1+e^{-z}}$

One layer network:  $f_n(x) = \sum_{j=1}^n c_j \phi(a_j \cdot x + b_j) + c_0 \quad (**)$



n nodes

In total:  $(d+2)n+1$  parameters ( $O(n)$ )



1717  
1718

14.4.16

2

# Theorem 1 (NB: $B_r \rightarrow$ general $B \ni 0$ possible bounded)

Let us consider:  $\phi$  ~~fixed~~ "sigmoid"  
 $B_r = \{x \in \mathbb{R}^d : |x| \leq r\}, r > 0$  } fixed  
 (Think uniform or empirical meas.)  $\mu$  a probability distribution on  $B_r$   
 Then  $\forall f \in \Gamma_C$  and  $\forall n \geq 1, \exists f_n$  of the form (\*\*)  
 such that  $\int_{B_r} (f(x) - f_n(x))^2 \mu(dx) \leq \frac{(2EC)^2}{n}$  (1)

Why is this theorem important/interesting?  $\& \sum_{j=1}^n |q_j| \leq 2Cr$   
 here, only  $O(nd)$

- classical approx methods (polynomial, Fourier, ...)
- ~~exponential~~  $\rightarrow$  "exponential" number of parameters in  $nd$  for same error  $O(\frac{1}{n})$
- fixed basis of  $O(n)$  functions  $\rightarrow$  error  $O(\frac{1}{n^{2/d}})$  at least  $\rightarrow$  curse of dimensionality! [proof later?] Sect. X
- ( $\Delta C$  also depends on  $d$ , ~~but not on  $n$~~ )
- (practical algorithm to find the parameters  $\rightarrow$  next time!)  
 low-complexity

2 proofs: not iterative / iterative  
 generalization:  $B_r \rightarrow B$

⊕ p.932, col.2:

= only approximation error here  
~~also training error and else~~

$$M = \frac{1}{N} \sum_{j=1}^N \sum_{x_j} : \frac{1}{N} \sum_{j=1}^N (f(x_j) - \hat{f}_n(x_j))^2 \leq \frac{C_f^2}{n}$$

"training" error (?)  $\hookrightarrow$  chosen to minimize the expression / or iterative

⊕ computational complexity

(1) Note:  $f \in \Gamma_C$  is defined on  $\mathbb{R}^d$  (the whole space) differentiable  
 one can as well consider  $f: B_r \rightarrow \mathbb{R}$  st.  $\exists$  an  $\forall$  extension  $\hat{f}$   
 of  $f$  which coincides on  $B_r$  and is in  $\Gamma_C$ .



# Proof

Pisier: "remarks on an unpublished result by Naoray" 1980 (!)

## Lemma 1 (NB: ref of Naoray & Pisier!)

Let  $H$  be a Hilbert space (scalar prod  $\langle \cdot, \cdot \rangle$ , norm  $\|\cdot\|$ )  
 $G \subset H$  be a set with  $\|g\| \leq b \quad \forall g \in G$   
 $\overline{\text{co}(G)}$  denotes the closure of the convex hull of  $G$  (clear?)

Assume  $\bar{f} \in \overline{\text{co}(G)}$ . Then  $\forall n \geq 1$  and  $\forall c' > b^2 - \|\bar{f}\|^2$ ,

$\exists g_1, \dots, g_n \in G$  and  $\alpha_1, \dots, \alpha_n \geq 0$  st  $\sum_{j=1}^n \alpha_j = 1$  with

$$\|\bar{f} - \sum_{j=1}^n \alpha_j g_j\|^2 \leq \frac{c'}{n}$$

NB: The surprise here is that an error  $O(\frac{1}{n})$  can be obtained with only  $n$  points (exactly what we want)

## Proof (probabilistic 😊)

•  $\bar{f} \in \overline{\text{co}(G)}$ , so  $\forall \delta > 0$  &  $\forall n \geq 1$ ,  $\exists f^* \in \text{co}(G)$  st.  $\|\bar{f} - f^*\| \leq \frac{\delta}{n}$

•  $f^* \in \text{co}(G)$ , i.e.  $\exists m$  suff. large st.  $f^* = \sum_{k=1}^m \gamma_k g_k^*$

for some  $\gamma_1, \dots, \gamma_m \geq 0$ ,  $\sum_{k=1}^m \gamma_k = 1$ ,  $g_1^* \dots g_m^* \in G$

• take now  $g$  drawn from  $g_1^* \dots g_m^*$  with probabilities  $\mathbb{P}(g = g_k^*) = \gamma_k$   
and pick  $g_1 \dots g_n$  <sup>randomly</sup> iid  $\sim g$ ,  $f_n = \frac{1}{n} \sum_{j=1}^n g_j$

•  $\mathbb{E}(f_n) = \mathbb{E}(g) = f^*$

•  $\mathbb{E}(\|f_n - f^*\|^2) = \mathbb{E}(\|\frac{1}{n} \sum_{j=1}^n (g_j - f^*)\|^2) = \frac{1}{n^2} \mathbb{E}(\sum_{j=1}^n \langle g_j - f^*, g_j - f^* \rangle)$

(Var)  $= \frac{1}{n^2} \mathbb{E}(\sum_j \|g_j - f^*\|^2) + \frac{1}{n^2} \sum_{j \neq l} \langle \mathbb{E}(g_j - f^*), \mathbb{E}(g_l - f^*) \rangle$

$= \frac{1}{n} \mathbb{E}(\|g - f^*\|^2) = \frac{1}{n} \mathbb{E}(\|g\|^2 - 2\langle g, f^* \rangle + \|f^*\|^2)$

$= \frac{1}{n} (\mathbb{E}(\|g\|^2) - 2\langle f^*, f^* \rangle + \|f^*\|^2) = \frac{1}{n} (\mathbb{E}(\|g\|^2) - \|f^*\|^2)$

$\leq \frac{1}{n} (b^2 - \|f^*\|^2)$   
hyp.

= key!



If the preceding inequality holds in expectation, this means  $\exists$  actual ~~the~~  $g_1, \dots, g_n \in G$ ,  $f_n = \frac{1}{n} \sum_{j=1}^n g_j$  s.t.

$$\|f_n - f^*\| \leq \frac{1}{n} (b^2 - \|f^*\|^2)$$

NB: uniform weights  $\alpha$

(skip details)

Now,  $\|f_n - \bar{f}\| \leq \|f_n - f^*\| + \|f^* - \bar{f}\| \leq \frac{b^2 - \|f^*\|^2 + \delta}{n}$

&  $\|\bar{f} - f^*\| \geq \|\bar{f}\| - \|f^*\| \Rightarrow \|f^*\| \geq \|\bar{f}\| - \frac{\delta}{n} \Rightarrow \|f^*\|^2 \geq \|\bar{f}\|^2 - \frac{2\delta\|f^*\|}{n}$

so  $\|f_n - \bar{f}\| \leq \frac{b^2 - \|\bar{f}\|^2 + \delta(1 + 2\|f^*\|)}{n}$   $\leftarrow$  can be taken arbitrarily close to  $b^2 - \|\bar{f}\|^2$  by taking  $\delta$  small #

Application

$G_\phi = \{ \gamma \phi(a \cdot x + b), |\gamma| \leq 2cr, a \in \mathbb{R}^d, b \in \mathbb{R} \}$

$\Gamma_{c, B_r} = \{ f: B_r \rightarrow \mathbb{R} \text{ s.t. } \exists \text{ regular extension } \tilde{f}: \mathbb{R}^d \rightarrow \mathbb{R} \text{ with } \int_{\mathbb{R}^d} |\omega| |\tilde{f}(\omega)| d\omega \leq c \}$

Theorem 2

$\forall f \in \Gamma_{c, B_r}$ , the function  $f(x) - f(0)$  is in the closure of the convex hull of  $G_\phi$ , where the closure is taken with respect to the  $L^2(B_r, \mu)$  norm:  $\|f - g\|^2 = \int_{B_r} (f(x) - g(x))^2 d\mu(x)$

(therefore the claim) [up to the constant: to come]

Proof steps: 1)  $\Gamma_{c, B_r} \subset \overline{\text{co}} G_{\cos}$

$$G_{\cos} = \left\{ \frac{\gamma}{|\omega|} (\cos(\omega \cdot x + b) - \cos(b)) : \omega \neq 0, |\gamma| \leq c, b \in \mathbb{R} \right\}$$

rem: A convex  $\Rightarrow$  A convex  $\Rightarrow$  ok!

2)  $G_{\cos} \subset \overline{\text{co}} G_{\text{step}}$  (i.e.  $\phi(z) = 1_{\{z \geq 0\}}$ )

$$G_{\text{step}} = \{ \gamma 1_{\{\alpha x \geq t\}} : |\gamma| \leq 2cr, |\alpha| \leq 1, |t| \leq r \}$$

(detail omitted)

3)  $G_{\text{step}} \subset \overline{\text{co}} G_\phi$



# Proof of 1)

~~quadrature~~

Let  $f$  be the extension to  $\mathbb{R}^d$  of  $f^u$ :  $\rightarrow FT \hat{f}(\omega)$

$x \in B$ :  $f(x) - f(0) = \int_{\mathbb{R}^d} (e^{i\omega \cdot x} - 1) \hat{f}(\omega) d\omega$ ,  $C_f = \int_{\mathbb{R}^d} |\omega| |\hat{f}(\omega)| d\omega \leq C$   
by assumption

write:  $\hat{f}(\omega) = e^{i\theta(\omega)} R(\omega)$

Because  $f$  is real-valued,

$$\bar{f}(x) = f(x) - f(0) = \text{Re} \left( \int_{\mathbb{R}^d} (e^{i\omega \cdot x} - 1) e^{i\theta(\omega)} R(\omega) d\omega \right)$$

$$= \int_{\mathbb{R}^d} (\cos(\omega \cdot x + \theta(\omega)) - \cos(\theta(\omega))) R(\omega) d\omega$$

$$= \int_{\mathbb{R}^d} \underbrace{\frac{C_f}{|\omega|} (\cos(\omega \cdot x + \theta(\omega)) - \cos(\theta(\omega)))}_{= g(x, \omega)} \underbrace{\frac{|\omega| \cdot R(\omega)}{C_f}}_{= p(\omega) \text{ prob. density}} d\omega$$

$|g(x, \omega)| \leq \frac{C_f}{|\omega|} \cdot |\omega \cdot x| \leq C_f r \leq C_r$  well def. (\*) (\*\*)

So  $\bar{f} = \text{convex combination of functions in } G_{\text{cos}}$

$\Rightarrow \bar{f} \in \overline{\text{co } G_{\text{cos}}}$

draw iid  $\omega_1, \dots, \omega_n \sim p(\omega)$  and define  $f_n(x) = \frac{1}{n} \sum_{j=1}^n g(x, \omega_j)$

$$\mathbb{E}(\|\bar{f} - f_n\|^2) = \mathbb{E} \left( \int_{B_r} (\bar{f}(x) - f_n(x))^2 \mu(dx) \right)$$

$$= \int_{B_r} \mathbb{E} \left( \left( \bar{f}(x) - \frac{1}{n} \sum_{j=1}^n g(x, \omega_j) \right)^2 \right) \mu(dx)$$

$$= \mathbb{E} \left( \left( \frac{1}{n} \sum_{j=1}^n (\bar{f}(x) - g(x, \omega_j)) \right)^2 \right) \stackrel{\text{indep}}{=} \frac{1}{n} \mathbb{E} \left( (\bar{f}(x) - g(x, \omega))^2 \right) \leq \frac{1}{n} \text{Var}(g(x, \omega)) \leq \frac{(C_r)^2}{n}$$

$\Rightarrow \exists$  an actual sequence  $f_n = \frac{1}{n} \sum_{j=1}^n g(x, \omega_j)$  st.  $\|\bar{f} - f_n\|^2 \leq \frac{(C_r)^2}{n}$ . #

(\*) Any  $g \in G_{\text{cos}}$  satisfies  $\|g\|^2 = \int_{B_r} (g(x))^2 \mu(dx) \leq \sup_{x \in B_r} |g(x)|^2 \leq \frac{(C_r)^2}{b^2}$   
BTW;

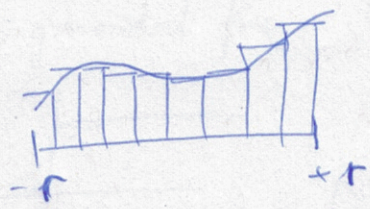
(\*\*)  $\nabla g^u(x, \omega) = \frac{C_f}{|\omega|} \sin(\omega \cdot x + \theta(\omega)) \cdot \omega$ :  $|\nabla g| \leq C_f \in C$



Proof of 2) [  $G_{\cos} \subset C \subset C \subset G_{\text{step}}$  ]

$$f \in G_{\cos} = f \circ g \circ h(x) \begin{cases} g(z) = \frac{\gamma}{|\omega|} (\cos(|\omega|z + b) - \cos(b)) \\ z = h(x) = \frac{\omega \cdot x}{|\omega|} \in [-r, r] \end{cases}$$

$g$  unif. continuous on  $[-r, r]$ ,  $|g'(z)| \leq |\gamma| \leq C \Rightarrow$  total variation of  $g$  on  $[-r, r] \leq 2Cr$   
 $\Rightarrow$  can be uniformly well approximated by step fns



more precisely,  $g(0) = 0$  so for  $z > 0$ ,  $g(z) \approx \sum_{i=1}^{k-1} (g(t_i) - g(t_{i-1})) 1_{z \geq t_i}$

add the two!  $z < 0$   $g(z) \approx \sum_{i=1}^{k-1} (g(-t_i) - g(-t_{i-1})) 1_{z \leq -t_i}$

$$\sum_{i=1}^{k-1} |g(t_i) - g(t_{i-1})| \leq \int_0^r |g'(z)| dz \leq |\gamma| r \leq Cr, \quad \sum_{i=1}^{k-1} |g(-t_i) - g(-t_{i-1})| \leq Cr$$

$g =$  linear combination of functions  $1_{z \geq t_i}, 1_{z \leq -t_i}$  with coeffs  $|c_i| \leq 2Cr$   
 (with  $k \rightarrow \infty$ )

$$\in \overline{CO \{ g(z) = \gamma 1_{z \geq t}, |\gamma| \leq 2Cr, |t| \leq r \}}$$

$$\Rightarrow f \in \overline{CO \{ f_0(z) = \gamma 1_{x \geq t}, |\gamma| \leq 1, |\gamma| \leq 2Cr, |t| \leq r \}} \quad \#$$

NB: all this with respect to  $\|\cdot\|_{\infty}$  norm  
 $\rightarrow$  also with respect to  $L^2(B_r, \mu)$  norm!

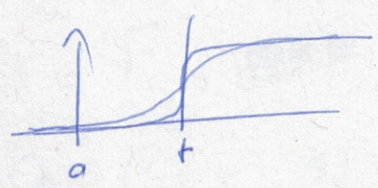


roof of 3)

detail:  $t \in [-r, r] \rightarrow t \in$  continuity points of ~~the distribution~~  
 the distribution of  $Z = d \cdot x$  induced by  $\mu$  on ~~the set~~  $B_r$   
 = dense set!  $\rightarrow G_{step}^M$ !  $G_{cos} \subset \overline{Co G_{step}^M}$

$G_{step}^{(M)} \subset \overline{Co G_\phi}$

take sigmoidal fns:  $\phi(|\alpha|(ax-t))$  with  $|\alpha| \rightarrow \infty$



$\xrightarrow{|\alpha| \rightarrow \infty} 1_{\alpha \cdot x \geq t}$   
 (except for  $\alpha$ 's st  $\alpha \cdot x = t$   
 but  $\mu(\{x : \alpha \cdot x = t\}) = 0$  by restriction on  $t$ ) #

Constants (end of proof)

Claim:  $f \in \Gamma_{c, B_r} \Rightarrow \|f - \bar{f}_n\|_{L^2(B_r, \mu)}^2 \leq \frac{(2rc)^2}{n}$  ~~XXXXXXXXXX~~

$\bar{f}(x) = f(x) - f(0)$

• if  $\|\bar{f}\| = 0$ , then ~~XXXXXXXXXX~~,  $\bar{f} \equiv 0$ : clear ( $c=0$  ok)

• assume  $\|\bar{f}\| > 0$ .

$G_\phi = \{ \gamma \circ \phi(a \cdot x - b), |\gamma| \leq 2rc, a \in \mathbb{R}^d, b \in \mathbb{R} \}$

$f \in G_\phi \Rightarrow \|f\|_{L^2(B_r, \mu)} \leq \|f\|_\infty \leq 2rc$  ( $|\phi| \leq 1$ )  $\rightarrow$  add hyp.

$\Gamma_{c, B_r} \subset \overline{Co G_\phi} \Rightarrow$  Lemma 1:  $\exists f \in G_\phi$  st  $\|\bar{f} - \bar{f}_n\|_{L^2(B_r, \mu)}^2 \leq \frac{(2rc)^2 - \|\bar{f}\|^2}{n}$   
 $\Rightarrow$  one can choose  $c' = (2rc)^2$  #

•  $|\phi| \leq 1$ : ok