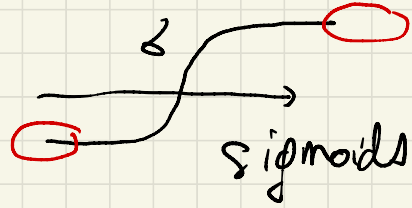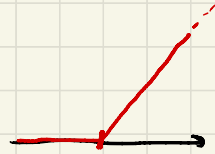# Recap : MLP

$$x_k \in \mathbb{R}^{M_k}$$
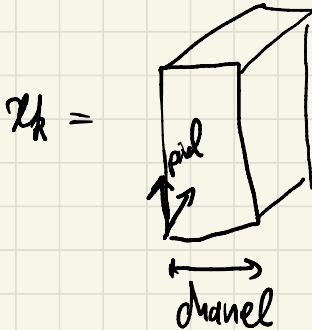
$$x_{k+1} = \sigma(W_k x_k + b_k)$$

$$\Theta_k = (W_k, b_k)$$

sigmoids

$$\sigma(x) = \begin{cases} \text{Arctan}(x) \\ \dfrac{e^x}{1+e^x} \end{cases} \qquad \sigma(x) = \text{Relu}(x)$$

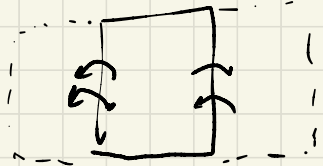Design choices : "Weight sharing" → convolution

audio / image / video

$$x_k =$$

paddel

channel

$$u_k = \left(x_k[\circ, \circ, s]\right)_{s=1}^{s}$$

"Translat° inv. lin operator ⟺ convolution"

$$T_\tau : x[.,.,.] \rightsquigarrow x[.-\tau_1, .-\tau_2, .]$$

$$\underset{\text{pixel}}{\underbrace{\phantom{x[.,.}}} \underset{\text{channel}}{\underbrace{\phantom{.]}}}$$

**Thm:** $W : x \in \mathbb{R}^{M \times M \times S} \longrightarrow \mathbb{R}^{M \times M \times T}$

that commutes with transl$^o$

$$W \circ T_\tau = T_\tau \circ W$$

$$\begin{array}{ccc} x & \xrightarrow{W} & Wx \\ \downarrow{T_\tau} & & \uparrow{T_{-\tau}} \\ & \xrightarrow{W} & \end{array}$$
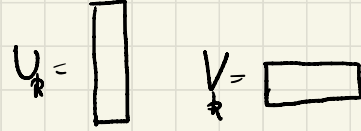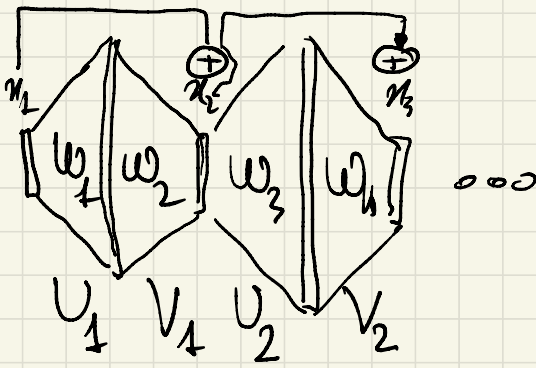
if and only if $\{\psi_{st}\}_{\substack{s=1...S \\ t=1...T}}$ so that

$$Wx = \left( \sum_s \psi_{st} * x[.,.,s] \right)_t$$

<u>Example</u>: $1.1$ conv$^o$

2nd design choice: skip connex$^o$

Res-Net

$x_2$   $x_3$

$\omega_1$ | $\omega_2$   $\omega_3$ | $\omega_4$   $\circ\circ\circ$

$U_1$ $V_1$ $U_2$ $V_2$

$U_k = \boxed{\phantom{x}}$   $V_k = \boxed{\phantom{xx}}$

$$x_{k+1} = x_k + \sigma\left[ V_k \sigma\left[ U_k x_k \right] \right]$$

↑
removed/batch norm.



[ Generative net
[ U-Net → sept. generative model
          diff.
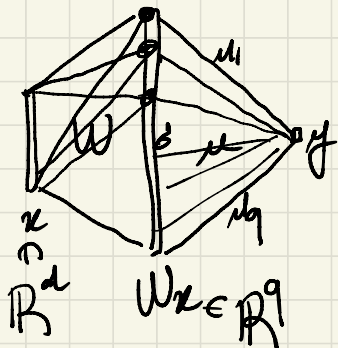
Transformers : convolut^= + nonlinearity → Attent^=

"Theory": Universality / Expressivity

2 layers NN

"Linear" single layer        $q = \#\text{neurons}$



$$y = \langle \mu, \sigma(Wx) \rangle$$
$$\mathbb{R}^q \quad \underbrace{\in \mathbb{R}^q}$$

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

$Wx \in \mathbb{R}^q$

$$W = \begin{bmatrix} w_1 \\ w_2 \\ \cdot \\ \cdot \end{bmatrix}$$

$$y = \sum_k \mu_k \, \sigma(\langle x, w_k \rangle + b_k)$$

$$\underbrace{\qquad\qquad\qquad\qquad}$$

$$\varphi(x) = \sum_k \mu_k \, \varphi_{w_k, b_k}(x)$$

$$\varphi_{w,b}(x) = \sigma(\langle x, w \rangle + b)$$    "Ridge funct°"

1-D



$\frac{1}{|w|}$

$-\frac{b}{w}$

small $w$

large $w$

$$\langle w, x \rangle + b = 0$$
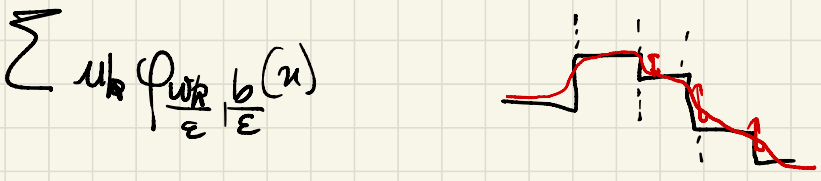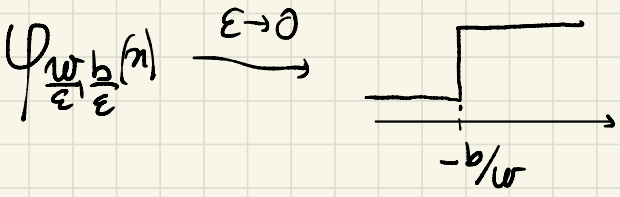
Bounded $\sigma$

Thm: [Cibenko] if $f(x)$ continuous, $\varepsilon > 0$ precision,

$R > 0$, $\exists q = \#$ neurons, $\exists \begin{matrix} w \\ b \\ u \end{matrix}$

such that, $\forall \|x\| \leq R$ $\quad |f(x) - \sum u_k \overset{+c}{\underset{w_k b_k}{\varphi}}(x)| \leq \varepsilon$

$\underbrace{\qquad\qquad}_{\text{N.N. } q \text{ neurons.}}$

$$\varphi_{\frac{w}{\varepsilon}, \frac{b}{\varepsilon}}(x) \xrightarrow{\varepsilon \to 0}$$

$-b/w$

$$\sum u_k \varphi_{\frac{w_k}{\varepsilon}, \frac{b}{\varepsilon}}(x)$$

## Speed of approx° · No free lunch

Barron's function: $\hat{f}(\omega) = \int f(x) \cdot e^{-i\langle\omega,x\rangle} dx$

                                Frequency

Barron's norm: $\|f\|_B \triangleq \int \|\omega\| \cdot |\hat{f}(\omega)| \, d\omega$

Sobolev norm: $\|f\|_{W^k}^2 = \int \|\omega\|^2 \cdot \|\hat{f}(\omega)\|^2 \, d\omega$ ♯

$$\hat{f'}(\omega) = i\omega \cdot \hat{f}(\omega) \quad \Bigg\}$$

$$= \int \|\hat{f'}(\omega)\|^2 \, d\omega$$

$$= \int |f'(x)|^2 dx = \|f'\|_{L^2}^2$$

                            $R > 0$ radius

__Thm:__ If $\|f\|_B < +\infty$ then $\forall q > 0$, $\exists$ a naral net

$$f_q(x) = \sum_{k \leq q} f_{\nu_k, b_k}(x) \quad \text{such that}$$

B(R) with radius R

$$\sqrt{\frac{1}{|B(R)|} \int_{B(x)} |f(x) - f_q(x)|^2 \, dx} \leq \frac{\|f\|_B}{\sqrt{q}} = \varepsilon$$

                                   Loss

With pol.

$$\frac{\|f\|_{\infty}}{q^{1/d}} = \varepsilon \qquad \leftarrow \text{curse of the dim}^c.$$

$\oplus$ No curse of dim, explicit
$\ominus$ 2 layers.
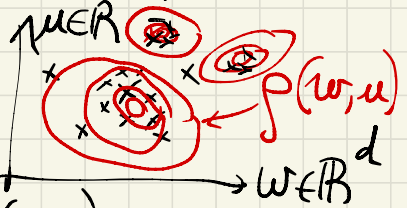$\ominus$ existence, how to find $f_q$

Proof: mean field analysis.

$$f_q(x) = \frac{1}{q} \sum_{k=1}^{q} u_k \, \sigma \left( \langle x, w_k \rangle \right)$$

$$\downarrow q \to +\infty$$

$$\theta_k = (u_k, w_k) \in \mathbb{R}^{d+1}$$



$$u \in \mathbb{R}$$
$$\rho(w, u)$$
$$w \in \mathbb{R}^d$$

$$f(x) = \int u \sigma \left( \langle x, w \rangle \right) d\rho(w, u).$$

Difficult: $\|f\|_{\mathcal{B}} < +\infty \iff \exists \rho \;,\; f = f^{\rho}$

$\underline{Q^\circ}$ : optimiz$^\circ$. chizat / Bach.

if $q \to \infty$, GD on the loss

$\underline{Q^\circ}$ : How many neurons $(q = ??)$



Loss $= 0$

$\underline{Q^\circ}$ : Generaliz$^\circ$ $\longrightarrow$ Implicit bias.

Deep & Wide        ResNet.

$$u_{k+1} \overset{?}{=} u_k + \frac{1}{k} V_k \cdot \sigma(U_k u_k)$$

$k = 1 \dots \boxed{K}$           $\in \mathbb{R}^q$.
        Depth

$\begin{cases} K = \text{depth} \\ q = \text{width} \text{ "} \to \infty \text{"} \quad \text{mean field} \end{cases}$

$\circledast$ $\quad \dfrac{u_{k+1} - u_k}{1/K} = \underset{\theta_k}{v}(u_k) \qquad \theta_k = (U_k, V_k)$

$$v_\theta(x) = U\sigma(Vx)$$



$x(t=1)$

$x(t)$

$x(t=0)$

"$K \to +\infty$"

$$\dot{x}(t) = \underset{\underline{\theta(t)}}{v}(x(t)) \qquad \text{out: } x(1)$$

$$x(t=0) = x_0$$

inf. depth limit: $\quad \psi_\theta: x(0) \longrightarrow x(1)$

$\uparrow$?? learning ??

$$(x_i)_{i=1}^{N} \longrightarrow (y_i)_{i=1}^{N}$$

learn: $\quad \underset{\{\theta(t)\}_{t=0}^{1}}{\min} \quad E(\theta) = \dfrac{1}{N} \sum_{i=1}^{N} \ell\left( \psi_\theta(x_i), y_i \right)$ $\Big|^2$

$\varphi^{\star}$

OPE$_1$  OPE$_2$  OPE$_3$  OPE$_4$

A = Id    B = Id

Thm: Bardoni Neurips 2022

$$C'(\theta) \cdot \|DE(\theta)\|^2 \overset{②}{\leqslant} E(\theta) - E(\theta^*) \overset{①}{\leqslant} C(\|\theta\|) \times \|DE(\theta)\|^2$$

P-L
Pollak-Łojajevitch.



① $\Rightarrow$ No loc. min. but $\omega \to +\infty$

② $\Rightarrow$ if $\theta$ "close" to $\theta^*$ $\omega \not\to +\infty$

$$A \in \mathbb{R}^{d \times D}$$

$$\underline{\underline{D}} \text{ large enough}$$

$$x_{k+1} = \mho_{\theta_k}(x_k)$$

$$X_k = \left(x_k^1, x_k^2 \; -- \; x_k^s\right) \quad \text{\#token}$$



$$X_{k+1} = X_k + \frac{1}{k} A_{\theta_k}(X_k)$$

$$\dot{X}(t) = A_{\theta(t)}(X(t))$$

$$\dot{x}^i(t) = \mho_{\theta(t)}\left(x^i(t); X(t)\right)$$

$$x_i \rightarrow \frac{x_i - m_i}{\| \quad \|}$$

$$v_{Q(k)}\left(x^i_j \, (x^{\tilde{j}})\right) = \sum_j A_{ij} \cdot (V x_j)$$

$$(K, Q, V)$$

$$\sum_j A_{ij} \right) \quad \text{Head}$$
$$\frac{}{2} \, 0$$

$$A_{ij} = e^{\frac{\langle K x_i, Q x_j \rangle}{\sqrt{d}}}$$