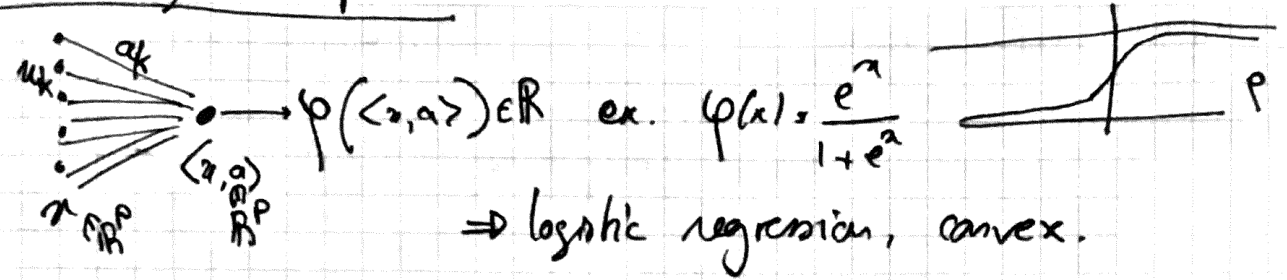


# Multilayer Perceptron with 1 Hidden Layer

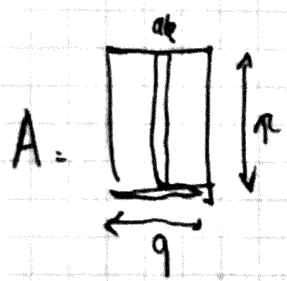
①

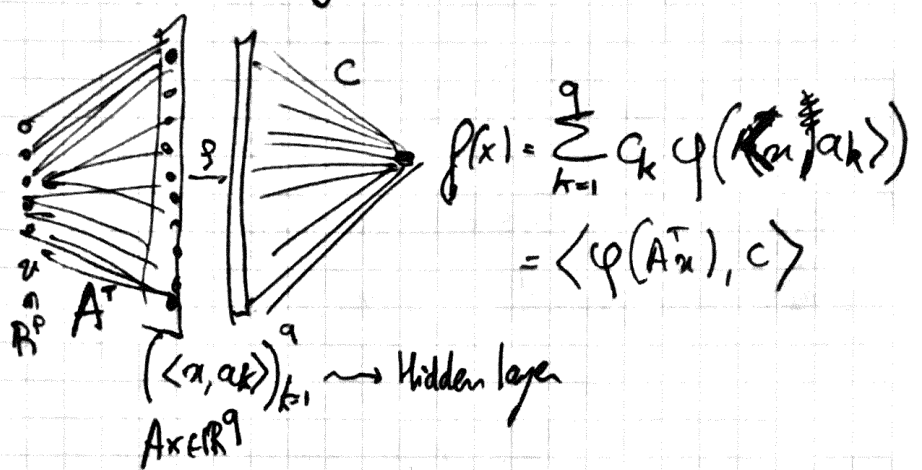
## 0 hidden layer (receptor)



$\Rightarrow$  logistic regression, convex.

## 1 hidden layer



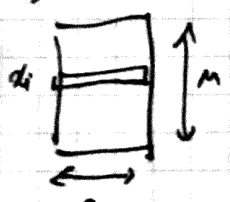


$(\langle x, a_k \rangle)_{k=1}^q \leadsto$  Hidden layer  
 $A^T x \in \mathbb{R}^q$

Req: include a bias

$$\underbrace{\langle x, a_k \rangle}_{\mathbb{R}} + \underbrace{b_k}_{\mathbb{R}} = \underbrace{\langle [x; 1], [a_k, b_k] \rangle}_{\mathbb{R}^{p+1}}$$

Training  
(Regression)



$$\min_{A, c} \tilde{E}(A, c) \triangleq \frac{1}{2m} \sum_{i=1}^m |f(x_i) - y_i|^2 = \frac{1}{m} \sum_{i=1}^m |\langle \phi(A^T x_i), c \rangle - y_i|$$

$$\tilde{E} = \frac{1}{2m} \left\| \underbrace{\phi\left(\underbrace{X \times A}_{\mathbb{R}^{m \times p} \times \mathbb{R}^{p \times q}}\right)}_{\mathbb{R}^q} + \underbrace{c}_{\mathbb{R}^q} - \underbrace{y}_{\mathbb{R}^m} \right\|^2$$

Prob:  $\tilde{E}$  is convex wrt  $c$  but non convex wrt  $A$  ...

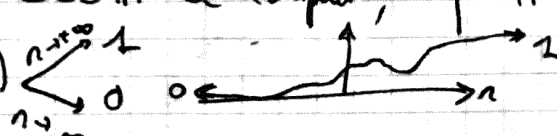
Gradient descent:  $\nabla_c \tilde{E}(A, c) = \underbrace{\phi(XA)^T}_{\mathbb{R}^m \times \mathbb{R}^q} (\underbrace{\phi(XA)c - y}_{\mathbb{R}^m})$

$$\nabla_A \tilde{E}(A, c) = \underbrace{X^T}_{\mathbb{R}^{p \times m}} \left[ \underbrace{\phi'(XA)}_{\mathbb{R}^{m \times q}} \odot \underbrace{(Rc^T)}_{\mathbb{R}^m \times \mathbb{R}^q} \right] \in \mathbb{R}^{p \times q}$$

Approximation theory:  $f(x) = \langle c, \varphi([A; b]^T [x; 1]) \rangle$

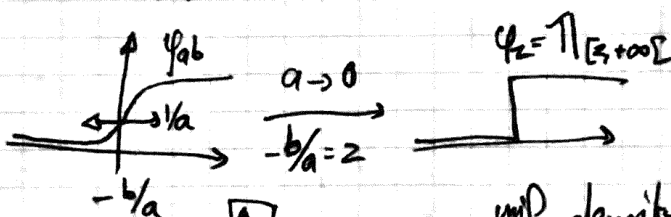
$$= \sum_{k=1}^q c_k \underbrace{\varphi(\langle a_k, x \rangle + b_k)}_{\triangleq \varphi_{a_k, b_k}(x)}$$

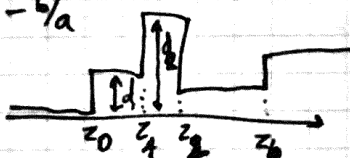
linear approx<sup>o</sup> in an adaptive dictionary  $(\varphi_{a_k, b_k})_{k=1}^q$  of  $q$  atoms

Thm: [Cybenko, 1989] let  $\Omega \subset \mathbb{R}^n$  be compact,  $\varphi: \mathbb{R} \rightarrow \mathbb{R}$  continuous "universal" and bounded,  $\varphi(x)$  


Then finite sums  $\sum_{k=1}^q c_k \varphi_{a_k, b_k}$  are dense in  $L^\infty(\Omega)$ .

Intuition in 1-D,  $p=1$



$$\sum_{k=1}^q d_k (\varphi_{z_k} - \varphi_{z_{k-1}}) = \text{step function}$$


unif. density of piecewise const<sup>o</sup> in  $L^\infty(\Omega)$   
 (unif. continuity on compact)

Req: Thm means that  $\forall f \in C(\Omega), \forall \epsilon > 0, \exists q, \exists A, b, c$  tq. 

$$\max_{x \in \Omega} \left| f(x) - \sum_{k=1}^q c_k \varphi(\langle x, a_k \rangle + b_k) \right| \leq \epsilon$$

~ Pbm: how does  $q$  depend on  $f$  and  $p$  (convergence speed)

~ Require smoothness hypothesis on  $f$

Thm: [Barron 93]: if  $\int_{\mathbb{R}^n} \|f(\omega)\| d\omega \leq C$  then

$$\int_{\|x\|_2 \leq r} \left\| f(x) - \sum_{k=1}^q c_k \varphi(\langle x, a_k \rangle + b_k) \right\|_2^2 dx \leq \|f\| (B(0, r)) \frac{(2\pi C)^2}{q}$$

# Proof of the universality theorem of Cybarko.

Prop 1: If  $\varphi$  is such that

$$\left[ \forall (a,b) \int_{\Omega} \varphi(\langle a, x \rangle + b) d\mu = 0 \right] \Rightarrow \mu = 0 \quad (\text{Discr})$$

$\uparrow \mu \in \mathcal{M}(\Omega) = \mathcal{C}(\Omega)^*$  Radon measures  
banded  $\mu(\Omega) < +\infty$

Then the universality thm is true.

Proof: Let  $\mathcal{J} = \left\{ \sum_k c_k \varphi_{a_k b_k} : \begin{matrix} q \in \mathbb{N} \\ a_k \in \mathbb{R}^p \\ b_k \in \mathbb{R} \end{matrix} \right\} \subset \mathcal{C}(\Omega)$  is a linear space

Let  $\bar{\mathcal{J}}$  be its closure in  $\mathcal{C}(\Omega)$  for  $\|\cdot\|_{\infty}$  (which is a Banach sp)

If  $\bar{\mathcal{J}} \neq \mathcal{C}(\Omega)$ , pick  $g \neq 0, g \in \mathcal{C}(\Omega) \setminus \bar{\mathcal{J}}$

We define a linear ~~operator~~  $L$  on  $\mathcal{J} \oplus \text{Span}(g)$  by  
 $\forall s \in \bar{\mathcal{J}}, L(s + \lambda g) = \lambda$  (so  $L \equiv 0$  on  $\bar{\mathcal{J}}$ )

$L$  is a bounded linear form, so by Hahn-Banach theorem, it can be extended in a bounded linear form  $\bar{L}: \mathcal{C}(\Omega) \rightarrow \mathbb{R}$

Since  $\bar{L} \in \mathcal{C}(\Omega)^* \sim \mathcal{M}(\Omega)$ , it can be represented as

$$\bar{L}(f) = \int_{\Omega} f(x) d\mu(x) \quad \text{with } \mu \neq 0$$

But  $\bar{L} \equiv 0$  on  $\bar{\mathcal{J}}$ , so  $\int_{\Omega} \varphi_{ab} d\mu = 0 \quad \forall (a,b)$ , so by

(Discr),  $\mu = 0$ , contradiction  $\blacksquare$

Prop 2: if  $\varphi \xrightarrow[r \rightarrow 0]{\infty}$  is continuous, then it satisfies (Discr).

Proof:

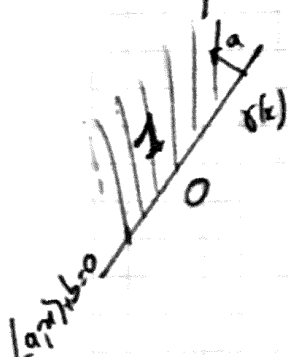
$$\varphi_{\lambda a, \lambda b + t}(x) = \varphi(\lambda(\langle a, x \rangle + b) + t) \xrightarrow{\lambda \rightarrow +\infty} \begin{cases} 1 & \text{if } \langle a, x \rangle + b > 0 \\ 0 & \text{if } \langle a, x \rangle + b < 0 \\ \varphi(t) & \text{if } \langle a, x \rangle + b = 0 \end{cases} \triangleq \gamma(x)$$

By Lebesgue dominated convergence, (banded by 1 on compact set)

$$\int \varphi_{\lambda a, \lambda b + t} d\mu \xrightarrow{\lambda \rightarrow +\infty} \int \gamma d\mu = \varphi(t) \cdot \mu(\Pi_{ab}) + \mu(H_{ab})$$

where  $\Pi_{ab} = \{x: \langle a, x \rangle + b = 0\}$   $H_{ab} = \{x: \langle a, x \rangle + b > 0\}$

If  $\mu$  is such that  $\int \varphi_{a', b'} d\mu = 0 \quad \forall (a', b')$  then  $\forall (a, b, t), \varphi(t) \mu(\Pi_{ab}) + \mu(H_{ab}) = 0$



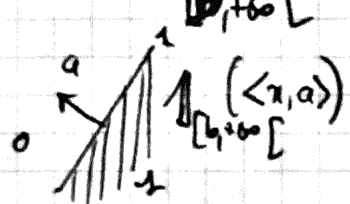
(4)

By selecting  $(t, t')$  such that  $\varphi(t) \neq \varphi(t')$ , one has that  $\forall (a, b) \begin{cases} \mu(H_{a,b}) = 0 \\ \mu(\Pi_{a,b}) = 0 \end{cases}$   
 One now needs to show that  $\mu = 0$ .

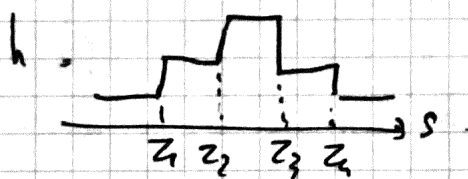
For  $\{h \in L^\infty(\mathbb{R}), \text{ s.t. } h \in \mathbb{R}\}$ , let  $F(h) \triangleq \int_{\mathbb{R}} h(\langle a, x \rangle) d\mu$ .

$F: L^\infty(\mathbb{R}) \rightarrow \mathbb{R}$  is a bounded linear form since  $\|F(h)\| \leq \|F\|_{\infty} \cdot \mu(\Omega)$

let  $h = \mathbb{1}_{[b, +\infty[}$  then  $F(\mathbb{1}_{[b, +\infty[}) = \int_{[b, +\infty[} d\mu = \mu(\{x: \langle a, x \rangle - b \geq 0\})$   
 $= \underbrace{\mu(\Pi_{a,-b})}_0 + \underbrace{\mu(H_{a,-b})}_0 = 0$



Similarly, for in 1D, for  $h(s) = \sum \epsilon_i \left( \mathbb{1}_{[z_i, +\infty[}^{(s)} - \mathbb{1}_{[z_{i+1}, +\infty[}^{(s)} \right)$



One has  $F(h) = 0$  for all piecewise constant  $f \in L^\infty$ .  
 By density in  $L^\infty(\mathbb{R})$ ,  $F(h) = 0 \quad \forall h \in L^\infty(\mathbb{R})$ .

Taking  $h(x) = e^{iax}$ , one has

$$\int \exp(i\langle x, a \rangle) d\mu(x) = \hat{\mu}(a) = 0 \quad \forall a.$$

By injectivity of the Fourier transform,  $\mu = 0$  ■