

Generalization for neural networks

Framework

Binary classification

input space $X = \mathbb{R}^n$

output space $Y = \{0, 1\}$

probability distribution P on $Z = X \times Y$

sample $(z_1 = (x_1, y_1), \dots, z_m = (x_m, y_m)) \in Z^m$ sample size m

Hypothesis class H : neural network

units = neurons

$n = 3$ input units

network

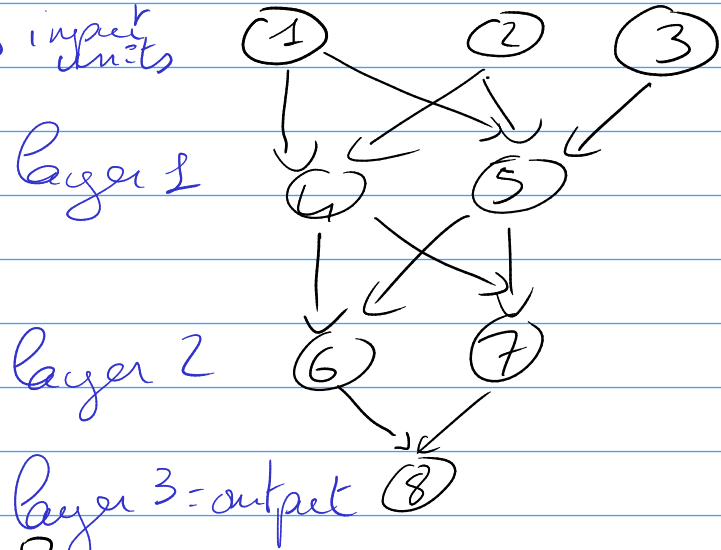
- feed-forward

activation function

f_r for unit r

typically: $f_r(u) = \frac{1}{1 + e^{-u}}$

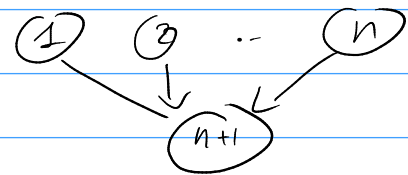
or $f_r(u) = \text{sgn}(u) = \begin{cases} 1 & \text{if } u \geq 0 \\ 0 & \text{otherwise} \end{cases}$



Perceptron: networks with 1 layer

Weights: for each unit r of layer 1, 2, ...

mixed: the outputs of the previous layer are a linear combination before passing thru the activation function f_r



if unit n is fed with entries z_1, \dots, z_d from the previous layer,
its output is $f_n(\sum_{j=1}^d \omega_j z_j - \theta)$

where $(\omega_1, \dots, \omega_d)$ are the corresponding connections weights
and θ is the threshold of unit n .

-> for every $x \in \mathbb{R}^n$, the network produces a prediction $\hat{y}(x)$
which is the output of the last unit.
The goal is to find the parameters (ω_n, θ_n) for all units n of
the network such that $\mathbb{P}(\hat{y}(x) = y)$ is as large as possible.

Linear threshold networks = feed-forward networks with
activation function $f_n(u) = \text{sgn}(u)$ for every unit n of the network

For every hypothesis $h \in H$, we will study its generalization error

we will rely on the training error of h :
$$e_p(h) = \mathbb{P}(q(x, y) \in Z, h(x) \neq y)$$

$$e_t(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{q(h(x_i)) \neq y_i} \quad z = (z_1, \dots, z_n)$$

Best achievable performance: $\text{opt}_p(H) = \inf_{h \in H} e_p(h) = \text{target}$

Formal definition of learning: **Probably Approximately Correct Learning**

A learning algorithm L for H is a function

$$L : \prod_{n=1}^{\infty} Z^n \rightarrow H$$

such that for all tolerance $\epsilon > 0$ and all risk $\delta > 0$, there exists a integer $m_L(\epsilon, \delta)$ such that if $m > m_L(\epsilon, \delta)$,

$$\forall P, \quad P^m \left(e_P(L(Z^m)) < \text{opt}_P(H) + \epsilon \right) \geq 1 - \delta$$

We say that H is learnable if there exists a learning algorithm for H .

The smallest sample size $m_L(\epsilon, \delta)$ reachable by an algorithm L is $m_H(\epsilon, \delta) = \text{sample complexity of } H$.

The smallest $m_L(\epsilon, \delta)$ reachable by algorithm L is called the sample complexity of L .

$$\hookrightarrow e_P(L(Z^m)) = \underbrace{\text{opt}_P(H)}_{\text{approximation error}} + \underbrace{\left(e_P(L(Z^m)) - \text{opt}_P(H) \right)}_{\text{estimation error}}$$

Empirical Risk Minimization:

$$L(h) = \arg \min_{h \in H} \hat{e}_3(h)$$

since $\mathbb{E}[\hat{e}_3(h)] = e_{\mathcal{P}}(h)$: for every $i \in \{1, \dots, n\}$,

$$\mathbb{E}[\mathbb{1}_{\mathcal{D}}(h(x_i) \neq y_i)] = \mathbb{P}(h(x_i) \neq y_i) = e_{\mathcal{P}}(h)$$

Prop 1: for every $h \in H$, $\mathbb{P}(\hat{e}_3(h) > e_{\mathcal{P}}(h) + \varepsilon) \leq \exp(-2m\varepsilon^2)$

Proof: $\hat{e}_3(h) = \sum_{i=1}^m \mathbb{1}_{\mathcal{D}}(h(x_i) \neq y_i)$

Bound: $\mathbb{P}(\hat{e}_3(h) < e_{\mathcal{P}}(h) - \varepsilon) \leq \exp(-2m\varepsilon^2)$ $\stackrel{\text{by } \mathcal{B}(e_{\mathcal{P}}(h))}{\text{by } \mathcal{B}(e_{\mathcal{P}}(h))}$

Prop 2: if $|H| < \infty$, then H is countable.

indeed, $\mathbb{P}^n(\max_{h \in H} |\hat{e}_3(h) - e_{\mathcal{P}}(h)| \geq \varepsilon) = \mathbb{P}^n(\bigcup_{h \in H} |\hat{e}_3(h) - e_{\mathcal{P}}(h)| \geq \varepsilon)$

$$\leq |H| \max_{h \in H} \mathbb{P}^n(|\hat{e}_3(h) - e_{\mathcal{P}}(h)| \geq \varepsilon) \leq 2|H| e^{-2m\varepsilon^2}$$

$$\Rightarrow \mathbb{P}^n(e_{\mathcal{P}}(L(\mathcal{Z})) \geq \hat{e}_3(L(\mathcal{Z})) + \varepsilon) \leq 2|H| e^{-2m\varepsilon^2}$$

$$\rightarrow \text{if } 2|H| \exp(-2m\varepsilon^2) \leq \delta$$

$$\text{that is if } m \geq \sqrt{\frac{\ln(2|H|/\delta)}{2\varepsilon^2}}$$

$$a_P(L(\mathcal{Z})) \leq \hat{e}_3(L(\mathcal{Z})) + \varepsilon$$

$$= \min_{h \in H} \hat{e}_3(h) + \varepsilon$$

$$\leq \hat{e}_3(h^*) + \varepsilon \quad \text{where } h^* \in \arg \min_{h \in H} e_P(h)$$

$$\leq (e_P(h^*) + \varepsilon) + \varepsilon$$

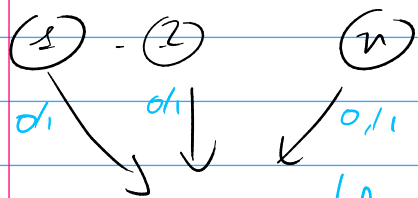
$$= \text{opt}_P(H) + 2\varepsilon$$

\hookrightarrow with prob $\geq 1 - \delta$, L returns $L(\mathcal{Z})$ s.t.

$$e_P(L(\mathcal{Z})) \leq \text{opt}_P(H) + \sqrt{\frac{2 \ln(2|H|)}{m} \frac{1}{\delta}}$$

\rightarrow the ERM has sample complexity at most $\frac{2}{\varepsilon^2} \ln\left(\frac{2|H|}{\delta}\right)$

→ for perceptrons: binary-weight perceptrons: coefficients $\in \{0, 1\}$
 for n entry units



0 ← threshold: $\neq 0, -n$

$$\rightarrow |H| = 2^n (n+1)$$

$$m_L(\epsilon, \delta) \leq \frac{2}{\epsilon^2} \left(n \ln(2) + \ln(n+1) + \ln \frac{2}{\delta} \right)$$

• R-bit perceptron: each coefficient has R bits: $0.b_1 \dots b_R$
 and the threshold where $b_j \in \{0, 1\}$

$$\rightarrow |H| \approx (2^R)^{n+1} = 2^{R(n+1)}$$

$$\hookrightarrow m_L(\epsilon, \delta) = \frac{2}{\epsilon^2} \left(R(n+1) \ln(2) + \ln \frac{2}{\delta} \right)$$

→ have to deal with continuity?

Homeworks . restricted model

• $y_i = t(x_i)$ for some $t \in H$

$\hookrightarrow \text{opt}_P(H) = 0$ reached by hypothesis $t \in H$

Thm. there is a learning algo L for h in the restricted model

with sample complexity $m_L(\epsilon, \delta) \leq \frac{1}{\epsilon} \ln\left(\frac{|H|}{\delta}\right)$

Generalization for infinite hypothesis classes.

Def: the growth function of H is

$$T_H: \mathbb{N} \rightarrow \mathbb{N}$$

$$m \mapsto \max \{ |H|_S | : S \subset X \text{ and } |S| = m \}$$

clearly $T_H(m) \leq 2^m = |\mathcal{A}^m|$

Def: a dichotomy (by h) of a finite subset S is one of the functions $f: S \rightarrow \{0, 1\}$ s.t. $f = h|_S$ for some $h \in H$

Def: H shatters S if $T_H(S) = 2^{|S|}$

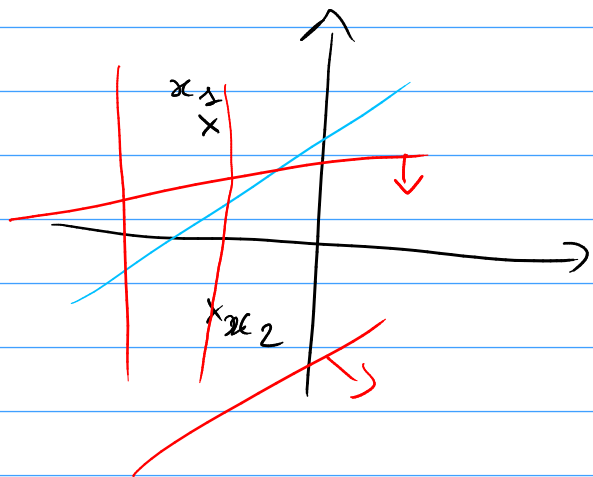
the VC-dimension of H is the largest size of a subset $S \subset X$ shattered by H .

[V.C = Vapnik-Chervonenkis]

Sauer's Lemma: $T_{H,d}(m) \leq \sum_{i=0}^d \binom{m}{i} \leq \left(\frac{em}{d}\right)^d$ polynomial in m

where $\binom{m}{i} = \frac{m(m-1)\dots(m-i+1)}{i!}$ $\binom{m}{0} = 1$ $\binom{m}{i} = 0$ if $i > m$

Ex: $n=2$ $H = \text{perceptrons} = \{x \rightarrow \text{sgn}(\omega_1 x_1 + \omega_2 x_2 - \theta)\} : (\omega_1, \omega_2, \theta) \in \mathbb{R}^3$



$|S|=1$

$f(x_1) = 0$
 $f(x_1) = 1$
 $T_H(1) = 2 = 2^1$

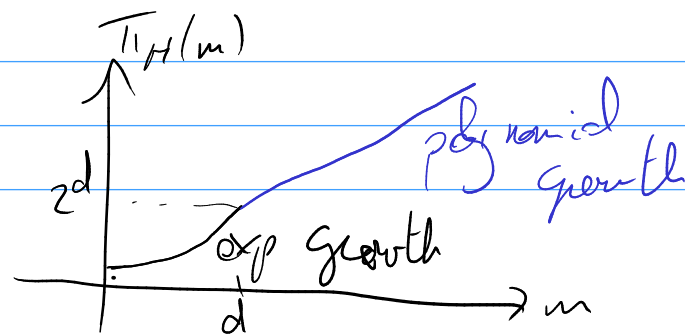
$|S|=2$

$x_1 \rightarrow \begin{pmatrix} 1 & 1 & 0 & 0 \\ x_2 \rightarrow -1 & 0 & 1 & 0 \end{pmatrix}$
 $S = \{x_1, x_2\} \rightarrow \{0, 1\}$

$T_H(2) = 4 = 2^2$

$\rightarrow VC(H) \geq 2$

Home work: find $VC(H)$.



Uniform Convergence and Learnability:

Thm: Suppose that $H : X \rightarrow \{0,1\}$ has finite VC-dimension $d \geq 1$.

Let L be an ϵ and δ learning algorithm.

Then L is a learning algorithm for H with sample complexity

$$m_L(\epsilon, \delta) \leq \frac{64}{\epsilon^2} \left(2d \ln \frac{12}{\epsilon} + \ln \frac{4}{\delta} \right) = \frac{64}{\epsilon^2} \ln \left(\left(\frac{12}{\epsilon} \right)^{2d} \frac{4}{\delta} \right)$$

(Recall: for finite H .

$$|H| = \left(\frac{12}{\epsilon} \right)^{2d}$$

$$m_L(\epsilon, \delta) \leq \frac{2}{\epsilon^2} \ln \left(\frac{2|H|}{\delta} \right)$$

Same lemma:
 $\frac{1}{|H|} \leq \left(\frac{e^{-m}}{d} \right)^d$

Compare with the bounds obtained above for binary-weight and k -bits perceptrons:

Thm: Uniform Convergence Theorem.

$$\text{finite } H \leq 2|H| \exp(-2m\varepsilon^2)$$

$$\frac{m}{N} |\exp(h) - \hat{e}_p(h)| \geq \varepsilon \text{ for some } h \in H \leq 4 \frac{\pi_H(2m)}{\delta} \exp(-\frac{m\varepsilon^2}{8})$$

from this it is easy to prove the previous theorem as in the finite case:

if $\forall h \in H, |\exp(h) - \hat{e}_p(h)| \leq \varepsilon$, then

$$\exp(L(z)) \leq \hat{e}_p(L(z)) + \varepsilon$$

$$= \min_{h \in H} \hat{e}_p(h) + \varepsilon$$

$$\leq \hat{e}_p(h^*) + \varepsilon \quad \text{for any } h^* \text{ st } \exp(h^*) \leq \text{opt}_p(H) + \alpha \quad \alpha > 0$$

$$\leq \exp(h^*) + \varepsilon + \varepsilon \leq \text{opt}_p(H) + 2\varepsilon + \alpha$$

thus for every $\alpha > 0 \rightarrow \exp(L(z)) \leq \text{opt}_p(H) + 2\varepsilon$

if $\varepsilon^2 \geq \frac{8}{m} \ln\left(\frac{\pi_H(2m)}{\delta}\right)$, then $4 \frac{\pi_H(2m)}{\delta} \exp(-\frac{m\varepsilon^2}{8}) \leq \delta$

and with prob $> 1 - \delta$ $\exp(L(z)) < \text{opt}_p(H) + \frac{32}{m} \sqrt{d \ln\left(\frac{2em}{d} + \ln \frac{4}{\delta}\right)}$

it suffices that $m \geq \frac{32}{\varepsilon^2} \left(d \ln m + d \ln \frac{2e}{d} + \ln \frac{4}{\delta} \right)$

→ Sauer's Lemma

since $Q_n \leq \alpha^2 - \ln \alpha - 1$ for all $\alpha, \alpha > 0$,

$$\begin{aligned} \frac{32}{\varepsilon^2} Q_n m &\leq \frac{32d}{\varepsilon^2} \left(\frac{\varepsilon^2}{64d} m + Q_n \frac{64d}{\varepsilon^2} - 1 \right) \\ &= \frac{m}{2} + \frac{32d}{\varepsilon^2} Q_n \frac{64d}{\varepsilon^2} \end{aligned}$$

\rightarrow it suffices that $m \geq \frac{m}{2} + \frac{32}{\varepsilon^2} \left(d \ln \frac{128}{\varepsilon^2} + Q_n \frac{4}{5} \right)$

so $m \geq \frac{64}{\varepsilon^2} \left(2d \ln \frac{128}{\varepsilon^2} + Q_n \frac{4}{5} \right)$ suffices.

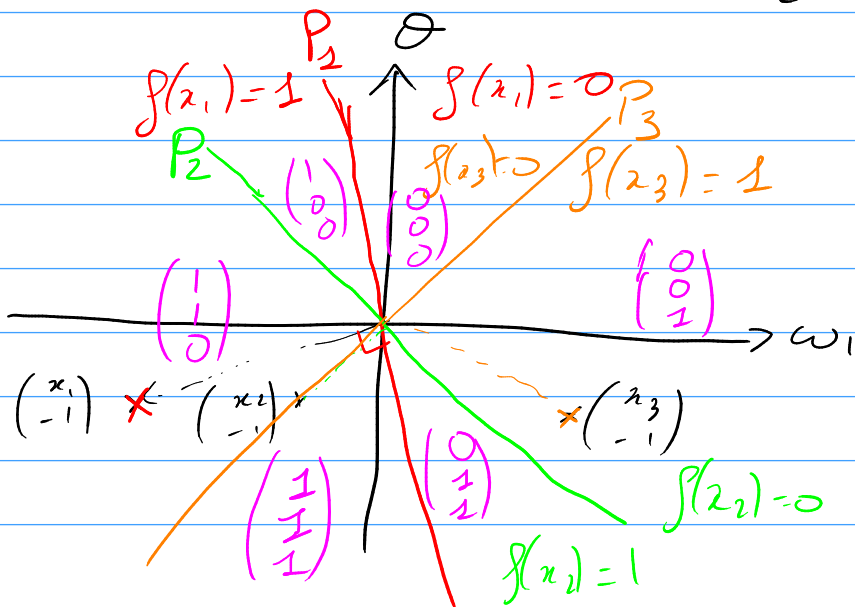
Growth Function and VC-dimension for neural networks

I. Perceptrons:

Thm: Let N be the real-weight perceptron with $n \in \mathbb{N}$ ^{real} inputs and H be the set of functions it computes.

$$H = \left\{ \phi_{\mathbf{w}, \theta} : \mathbb{R}^n \rightarrow \{-1, 1\} \mid \phi_{\mathbf{w}, \theta}(\mathbf{x}) = \text{sgn}(\omega_1 x_1 + \dots + \omega_n x_n - \theta) : (\omega_1, \dots, \omega_n, \theta) \in \mathbb{R}^{n+1} \right\}$$

then $\Pi_H(m) = 2 \sum_{k=0}^m \binom{m}{k}$



$n=1$

$$P_j = \{ \omega, \theta \mid \omega x_j - \theta = 0 \}$$

$$\text{sgn}(\omega_1 x_1 - \theta)$$

$$\Pi_H(3) = 6 = 2 \sum_{k=0}^3 \binom{3}{k}$$

$\{x_1, x_2, x_3\}$ is not shattered by H .

Lemma: for a set $S = \{a_1, \dots, a_m\} \subseteq \mathbb{R}^n$, let $P_1, \dots, P_m \subseteq \mathbb{R}^n$

the hyperplanes $P_i = \{(\omega, 0) \in \mathbb{R}^{n+1} : \omega^T a_i - 0 = 0\}$


Then $|H_S| = \mathcal{CC}(\mathbb{R}^{n+1} \setminus \bigcup_{i=1}^m P_i)$

↑
no. of connected components -

Proof: Homework 3.

Remark: $\alpha_i = \begin{pmatrix} a_i \\ -1 \end{pmatrix} \rightarrow P_i = \{v \in \mathbb{R}^{n+1} : \alpha_i^T v = 0\}$

Def: a set of points in \mathbb{R}^n is in general position if no subset of $k+1$ points lies on a $(k-1)$ -plane for all $k = 1, \dots, n$.

ex:  \mathbb{R}^n - not in general position

Prop: The $(a_i)_{i=1}^m$ are in general position iff every subset of up to $n+1$ points in $\{r_1, \dots, r_m\}$ is linearly independent.

Lemma: for $m, d \in \mathbb{N}$, suppose $T = \{r_1, \dots, r_m\} \subset \mathbb{R}^d$ has every subset of no more than d points linearly independent.

Let $P_i = \{v \in \mathbb{R}^d : v^T a_i = 0\}$ for $i = 0, \dots, m$ and

define $C(T) = CC(\mathbb{R}^d, \bigcup_{i=1}^m P_i)$

Then $C(T)$ depends only on m and d ; $C(T) = C(m, d)$

and for all $m, d \geq 1$: $C(m, d) = 2 \sum_{k=0}^{d-1} \binom{m-1}{k}$

[then we apply the Bra with $d = n+1$].

Proof

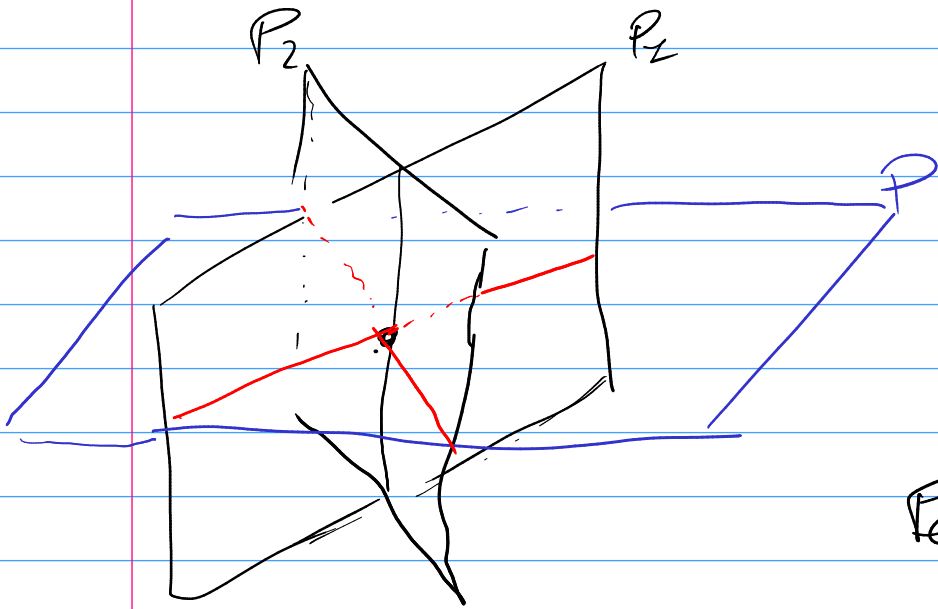
$$c(m, 1) = 2 \quad \text{for } m \geq 1$$

$$c(1, d) = 2 \quad \text{for } d \geq 1$$

Proof by induction: assume claim is true for all $T \subset \mathbb{R}^d$ $H_i \leq m$
 $j \leq d$

Suppose P_1, \dots, P_m satisfy the independence condition

and let P be another hyperplane also satisfying the independence conditions -



each intersection of a P_i with P is a $(d-2)$ plane in the $(d-1)$ hyperplane P and all these planes satisfy the independence condition

Each of them cuts a cc of $\mathbb{R}^d \setminus \bigcup_{i=1}^m P_i$ in two -

$$\begin{aligned}
\text{Hence, } C(m+1, d) &= C(m, d) + C(m, d-1) \\
&= 2 \left(\sum_{k=0}^{d-1} \binom{m-1}{k} + \sum_{k=0}^{d-2} \binom{m-1}{k} \right) \\
&= 2 \left(\binom{m-1}{0} + \sum_{k=1}^{d-1} \left(\binom{m-1}{k} + \binom{m-1}{k-1} \right) \right) \\
&= 2 \left(\binom{m-1}{0} + \sum_{k=0}^{d-1} \binom{m}{k} \right) = 2 \sum_{k=0}^{d-1} \binom{m}{k}
\end{aligned}$$

→ if S is in general position, then
 $|H_S| = C(m, n+1) = 2 \sum_{k=0}^n \binom{m-1}{k}$

if S is not in general position, it can be seen that
 $|H_S| \leq C(m, n+1)$

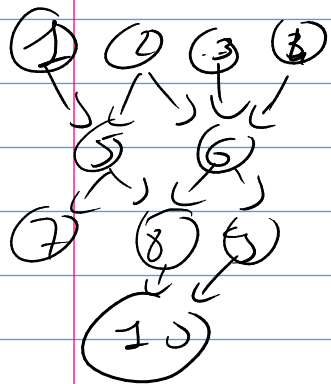
$$\rightarrow |H_S(m)| = 2 \sum_{k=0}^n \binom{m-1}{k} = \begin{cases} 2^m & \text{if } n \geq m-1 \\ 2^m - 2 \sum_{k=n+1}^{m-1} \binom{m-1}{k} < 2^m & \text{otherwise} \end{cases} \rightarrow \boxed{\text{VCDim}(N) = n+1}$$

(no. of parameters $(\omega_1, \omega_2, \dots, \omega_n, 0)$)

Growth function of Perceptron threshold networks

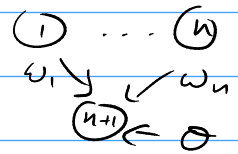
activation function $\rho(u) = \text{sgn}(u)$

k units $1, 2, \dots, k$ such that the output of i is an input of j only if $i < j$



I denote by W the total number of weights and thresholds

(In a perceptron $k = n + 1 = W$)



$H =$ the set of functions that are LTN with k units and W parameters

Thm: For $m \geq W$ the growth function of H is

$$\overline{N}_H(m) \leq \left(\frac{e m k}{W} \right)^W$$

hence $\text{vc dim}(H) \leq 2W \ln \left(\frac{2k}{\ln(2)} \right)$

Proof

• first consider unit (1) (below the entry layer):
it is a perceptron, hence the number of dichotomies it computes is at most

$$D_1(s) = \binom{e_m}{d_1}^{d_1}$$

where d_1 is the number of parameters attached to unit 1.

• consider now unit l $2 \leq l \leq h$.
its output depends on the output of unit $j \leq l$
and on its d_l parameters.

\hookrightarrow for each of the $D_{l-1}(s)$ states of the network on units $1, \dots, l-1$, it can read at most

$$D_l(s) \leq D_{l-1}(s) \binom{e_m}{d_l}^{d_l}$$

different states
 \Rightarrow by induction $D_h(s) \leq \prod_{l=1}^h \binom{e_m}{d_l}^{d_l}$

$$\text{and } \ln D_h(s) \leq \sum_{l=1}^h d_l \ln \binom{e_m}{d_l}$$

→ $T_H(m) = \max_R D_R(S)$ over all possible strings S :

$$\ln T_H(m) \leq \sum_{e=1}^R d_e \ln \frac{em}{d_e}$$

$$\frac{1}{W} \ln T_H(m) + \ln \left(\frac{W}{em} \right) \leq \sum_{e=1}^R \frac{d_e}{W} \ln \frac{em}{d_e} + \ln \frac{W}{em}$$

$$= \sum_{e=1}^R \frac{d_e}{W} \ln \frac{em}{d_e} + \sum_{e=1}^R \frac{d_e}{W} \ln \frac{W}{em}$$

$$= \sum_{e=1}^R \frac{d_e}{W} \ln \frac{W}{d_e} = \sum_{e=1}^R p_e \ln \frac{1}{p_e} = H(f)$$

where $p_e = \frac{d_e}{W}$ $\sum_{e=1}^R p_e = \frac{\sum d_e}{W} = 1$

where $H(p) = \sum_{e=1}^R p_e \ln \frac{1}{p_e} \leq \ln R$ (reached for $\frac{d_e}{W} = \frac{1}{R}$)

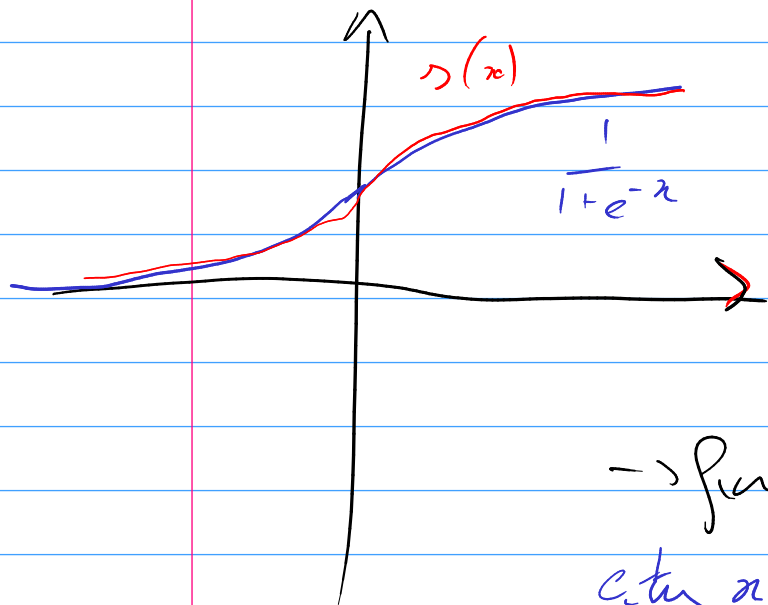
$$\rightarrow T_H(m) \leq \left(\frac{emR}{W} \right)^W$$

NB: one can prove that there exist neural architectures N such that $\forall c \dim(H) \geq \frac{3W}{5}$

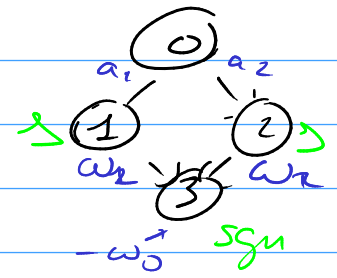
what for other activation functions? \rightarrow conditions are needed!

then: define $s(x) = \frac{1}{1+e^{-x}} + c x^3 e^{-x^2} \sin(x)$ for $c > 0$

one can show that if c is small enough, $s(x) \rightarrow \frac{1}{1+e^{-x}}$
 $s''(x) = \begin{cases} < 0 & \text{if } x > 0 \\ > 0 & \text{if } x < 0 \end{cases}$ $s(x) \rightarrow 0$ as $x \rightarrow -\infty$



We consider the neural net N
 one real entry
 two 1st-layer computation units
 one output unit



\rightarrow functions of N are of the form: $x \rightarrow \text{sgn}(w_0 + w_1 s(a_1 x) + w_2 s(a_2 x))$
 entry x and parameters $a_1, a_2, w_1, w_2, w_0 \in \mathbb{R}$

then VC-dim $(M) = +\infty$

Lemma: the function $F = \sum_{n=1}^{\infty} \frac{\sin(ca)}{n}$, $c \in \mathbb{R}$ of functions defined on \mathbb{N} has VC-dim(F) = ∞ .

This proves the previous theorem since:

$$h_c(x) = \sum_{n=1}^{\infty} \frac{\sin(ca)}{n} + \sum_{n=1}^{\infty} \frac{\sin(-ca)}{n} - 1$$

$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow \quad \uparrow$
 $\omega_1=1 \quad a_1=a \quad \omega_2=1 \quad a_2=-a \quad \omega_0=-1$

$$= \frac{1}{1+e^{-ca}} + \frac{1}{1+e^{ca}} - 1$$

$\underbrace{\hspace{10em}}_1$

$$+ 2c(ca)^3 e^{-ca^2} \sin(ca)$$

$$= 2c(ca)^3 \sin(ca) e^{-ca^2}$$

for $a > 0$ and $c > 0$, has the same sign as $\sin(ca)$ \square

Proof of the Q2:

for $d \in \mathbb{N}$, choose $a_i = 2^{i-1}$ for $i = 1, \dots, d$.

We show that $\{a_1, \dots, a_d\}$ is sheltered by F

since d is arbitrary, this shows that $\text{UCd}_-(F) = \infty$

for $(b_1, \dots, b_d) \in (0, 1)^d$ let $c = \sum_{j=1}^d 2^{-j} b_j + 2^{-(d+1)} = 0.b_1 b_2 \dots b_d 1$

take $a = 2\pi c$

for all i :

$$\begin{aligned} \text{sgn}(\sin(a a_i)) &= \text{sgn} \left(\sin \left(\frac{2\pi}{1} \left(\sum_{j=1}^d 2^{-j} b_j + 2^{-(d+1)} \right) / 2^{i+1} \right) \right) \\ &= \text{sgn} \left(\sin \left(\sum_{j=1}^d 2^{i-j} \pi b_j + \pi b_i + \sum_{j=i+1}^d 2^{i-j} \pi b_j + 2^{i-d-1} \pi \right) \right) \\ &= \text{sgn} \left(\sin \left(\pi \left(b_i + \underbrace{\sum_{j=1}^d 2^{-j} b_{i+j}}_{\in (0, 1)} + 2^{-(d-i+1)} \right) \right) \right) \\ &= 1 - b_i \end{aligned}$$