# 3.4 VC dimension of linear classifiers and neural networks

- In this section we want to compute or bound the VC dimension of important hypotheses classes

- We start with the class of linear hypotheses on $\mathcal{X} = \mathbb{R}^d$

$$\mathcal{L}_d = \left\{ h_{\mathbf{w},b}(\mathbf{x}) := \mathrm{sgn}(\mathbf{w} \cdot \mathbf{x} + b) \mid \mathbf{w} \in \mathbb{R}^d, \ b \in \mathbb{R} \right\}$$

and later also consider feedforward neural networks

- As an intermediate step we also consider the subclass of separating hyperplanes $H_{\mathbf{w},0}$ through the origin $\mathbf{0} \in \mathbb{R}^d$

$$\mathcal{L}_d^0 := \left\{ h_{\mathbf{w}}(\mathbf{x}) := \mathrm{sgn}(\mathbf{w} \cdot \mathbf{x}) \mid \mathbf{w} \in \mathbb{R}^d \right\} \subset \mathcal{L}_d$$

# VC dimension of linear hypotheses

**Lemma 3.34:**

A set $M = \{\mathbf{x}_1, \ldots, \mathbf{x}_m\} \subseteq \mathcal{X}$ is shattered by $\mathcal{L}_d^0$ **if and only if** the vectors $\mathbf{x}_j \in \mathbb{R}^d$, $j = 1, \ldots, m$, are linearly independent.

**Theorem 3.35:**

We have $\mathrm{VCD}(\mathcal{L}_d^0) = d$ as well as $\mathrm{VCD}(\mathcal{L}_d) = d + 1$.

**Consequences:**

- The sample complexity $m_{\mathcal{L}_d}$ grows only linearly with number of features $d$

$$m_{\mathcal{L}_d} \in \mathcal{O}(d)$$

- Therefore, linear classifiers are suitable for large number of features $d \gg 1$, e.g., learning classification rules in text analysis

Proof of lemma 3,34:

→ let $M = x_1 .. x_m \leq \mathbb{R}^d$ with $x_j; j = 1, \ldots, m$ being linearly independent. then: $A = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix} \in \mathbb{R}^{m \times d}$ - then the matrix has

a rank of $(m)$ which have to be smaller than $d$ to be linearly independent $(Rank(A) = m \leq d)$ hence for any labeling $b = (b_1, \ldots b_m) \in \{-1, +1\}^m \in \mathbb{R}^m$ means transpose

we can find a solution $w \in \mathbb{R}^d$ of linear system $Aw = b$ becaus matrix has full rank

Thus: the hypothesis of $h_b(x) = sgn(w_b \cdot x)$ yields

$h_b(x_1) = h_1, \ldots h_b, h_b(x_m) = b_m$. Since $b \in \{-1, +1\}^{km}$ was

arbitrary $\mathcal{L}_d^o$ shatters $M$. $(|\mathcal{L}_d^o|_m = 2^{|m|})$.

Rank number of independent numbers of columns of matrix

$\Rightarrow$ now let $\mathcal{L}_d^o$ shatter $M = \{x_1, \ldots x_m\}$ we then argue that the $x_i, j = 1, \ldots m$ can not be linearly independent by contradiction

if the points of $x_i \in m$ are linearly independent then
$\exists \, a = (a_1, \ldots, a_m)^T \neq 0$ such that $a_1 x_1 + \cdots + a_m x_m = 0$
then let $I_t = \{ i \, ; \, a_i > 0 \}$ and

according to linear Algebra

$I_- = \{ i : a_i < 0 \}$ & then we make an extension:

case 1: $I_t = \emptyset$ (null). then let $b \in \{-1, +1\}^m$ with
$b_i = +1$ for $i \in I_-$. this labeling of $m$ cannot be
reproduced by a hypothesis $h_w \in \mathcal{H}$ because

$$0 > \sum_{i \in I_-} a_i \cdot (w \cdot x_i) = w \cdot \left( \sum_{i \in I_-} a_i x_i \right)$$

$\geq 0$ becaus $b_i = 0$

vector here is $0$ which is contradiction

because of contradiction, this contradicts that $M$ is shattered by $\mathcal{L}_d^o$.

- Case 2: $I_t \neq \emptyset$, then let $b = \{-1, +1\}^m$ with $b_i = +1 \ \forall i \in I_t$ & $b_i = -1 \ \forall i \in \overline{I}_-$. again this labeling of $M$ can't be reproduced by a $h_w \in \mathcal{L}_d^o$, because then

$$0 \leq \underbrace{\sum_{i \in I_t}}_{>0} a_i \underbrace{(w \cdot x_i)}_{>0} = w \cdot \Big( \underbrace{\sum_{i \in I_t} a_i x_i}_{= -\sum_{i \in \overline{I}_-} a_i x_i} \Big)$$

$$= -\sum_{i \in I_-} \underbrace{a_i}_{<0} \underbrace{(w x_i)}_{<0} < 0 \quad \text{⨇ a contradiction.}$$

this contradicts that $M$ is shattered by $h_d$

Thus $x_1, \ldots x_m$ can't be linearly independent.

otherwise, we do have $L$ that can't be

reproduced.

Proof of Thm 3.35:

$VCD(\mathcal{H}_d) = d$ follows by lemma 3.34 & $VCD(\mathcal{H}_d) = d+1$ which can be showed by lemma 3.34 as well. to the end note that any $h_{w,b} \in \mathcal{H}_d$,
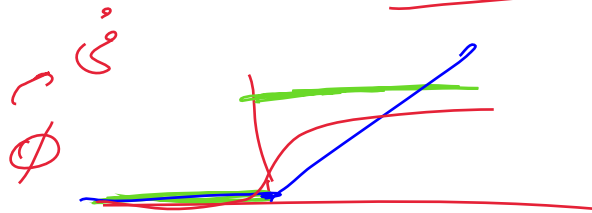
$h_{w,b}(x) = \text{Sgn}(b+wx)$ which corresponds to

$h_{w,b}(x) = h_{(wb)}((x,1)^T)$, $h_{(w,b)} \in \mathcal{H}_{d+1}$

thus $VCD(\mathcal{H}_d) \leq VCD(\mathcal{H}_{d+1}) = d+1$, Since $\mathcal{H}_d$ shatters any set $M = \{x_1, \dots x_m\} \cup \{0\} \subseteq \mathbb{R}^d$

with linearly independable + $x_j \in \mathbb{R}^d_j$ because

$$A = \begin{bmatrix} x_1^T & 1 \\ \vdots & \vdots \\ x_m^T & \vdots \\ 0 & \end{bmatrix}$$ has full rank $(A) = m+1$ then

we set $VC() (d_{cl}) = d + 1$



blue is ReLu function
red is Sigmoid function
green is Sign "

Sign function is better in terms of class.fication
better than ReLu 2 Sigmoid since the start $[0, (00)]$
$\Rightarrow$ all for output such Softmax

# Repetition: Feedforward neural networks (FNN)

- Recall a FNN consists of $L$ layers of $n_k$ neurons processing and passing information from layer to layer

- Each neuron $v_{k,i}$, $k = 1, \ldots, L$, $i = 1, \ldots, n_k$ is a linear hypothesis

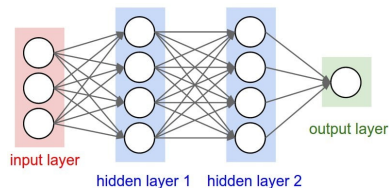$$y_{k,i} = v_{k,i}(\mathbf{y}_{k-1}) := \phi \left( \sum_{j=1}^{n} w_j \; y_{k-1,j} + b \right)$$

*layers* (handwritten annotation)

with activation function $\phi \colon \mathbb{R} \to \mathbb{R}$.



collects of wieghs in layers (handwritten annotation)

- The output of the $k$th layer $V_k = \{v_{k,1}, \ldots, v_{k,n_k}\}$ can then be written by

$$\mathbf{y}_k = \phi \circ f_{\mathbf{W}_k, \mathbf{b}_k}(\mathbf{y}_{k-1}), \qquad f_{\mathbf{W}_k, \mathbf{b}_k}(\mathbf{y}) := \mathbf{W}_k \; \mathbf{y} + \mathbf{b}_k$$

where $\phi$ is applied componentwise and we introduced the layerwise weight matrices $\mathbf{W}_k \in \mathbb{R}^{n_k \times n_{k-1}}$ and bias vectors $\mathbf{b}_k \in \mathbb{R}^{n_k}$

- The whole neural network is then a hypothesis $h\colon \mathcal{X} \to \mathcal{Y}$ of the form

$$h(\mathbf{x}) = \rho \circ f_{\mathbf{W}_L,\mathbf{b}_L} \circ \phi \circ f_{\mathbf{W}_{L-1},\mathbf{b}_{L-1}} \circ \phi \circ \cdots \circ \phi \circ f_{\mathbf{W}_1,\mathbf{b}_1}(\mathbf{x}),$$

where $\phi\colon \mathbb{R} \to \mathbb{R}$ as well as $\rho\colon \mathbb{R} \to \mathcal{Y}$ are the chosen activation functions

- Given a certain architecture $(V, E)$ with $V = (V_0, \ldots, V_L)$ consisting of
    - collection of layers $V = (V_0, \ldots, V_L)$ where $V_k = \{v_{k,1}, \ldots, v_{k,n_k}\}$ and $|V_0| = d$, $|V_L| = 1$,
    - collection of communication edges between adjacent layers:

$$E \subseteq \{(v_{k,i}, v_{k+1,j})\colon v_{k,i} \in V_k \text{ and } v_{k+1,j} \in V_{k+1}\}.$$

and chosen activation functions $\phi\colon \mathbb{R} \to \mathbb{R}$, $\rho\colon \mathbb{R} \to \mathcal{Y}$ we introduce the class of all FNN $h\colon \mathcal{X} \to \mathcal{Y}$ with just this architecture

$$\mathcal{H}_{V,E,\phi,\rho} = \big\{[\rho \circ f_{\mathbf{W}_L,\mathbf{b}_L}] \circ [\phi \circ f_{\mathbf{W}_{L-1},\mathbf{b}_{L-1}}] \circ \cdots \circ [\phi \circ f_{\mathbf{W}_1,\mathbf{b}_1}]\colon \mathbf{W}_k \in \mathbb{R}^{n_k \times n_{k-1}}, \mathbf{b} \in \mathbb{R}^{n_k},$$
$$[\mathbf{W}_k]_{i,j} \neq 0 \text{ iff } (v_{k,i}, v_{k+1,j}) \in E\big\}$$

*Output activation fonction* (handwritten annotation pointing to $\rho$)

- If $\phi = \rho$, we write only $\mathcal{H}_{V,E,\phi}$.

# The VC dimension of neural networks

> **Theorem 3.36:**
>
> Let $p_{V,E} = \sum_{k=1}^{L} n_k + |E|$ denote the number of parameters of the hypothesis class $\mathcal{H}_{V,E,\mathrm{sgn}}$. We have
>
> $$\mathrm{VCD}\left(\mathcal{H}_{V,E,\mathrm{sgn}}\right) \in \mathcal{O}(p_{V,E}\ \ln(p_{V,E})).$$

*which need to be trained*

*sum number of bias*

*we now use this in practice*

*larger might lead to overfitting*

- FNN with $\mathrm{sgn}$ as activation function are **PAC-learnable**, but the learnability decreases with the size of the network

- Lower bounds on $\mathrm{VCD}\left(\mathcal{H}_{V,E,\mathrm{sgn}}\right)$ can also be proved, as well as the VC dimension for other choices of $\sigma$ and $\rho = \mathrm{sgn}$

|  | $\mathrm{VCD}(\mathcal{H}_{V,E,\sigma,\mathrm{sgn}})$ | |
|---|---|---|
|  | Lower bound | Upper bound |
| sign | $\Omega(p\ln p)$ | $\mathcal{O}(p\ln p)$ |
| sigmoid | $\Omega(|E|^2)$ | $\mathcal{O}(p^2)$ |
| ReLU | $\Omega(L\ p\ln(p/L))$ | $\mathcal{O}(L\ p\ln p)$ |

- To prove Theorem 3.36 we exploit the special structure of $\mathcal{H} := \mathcal{H}_{V,E,\mathrm{sgn}}$

- Let $V = (V_0, \ldots, V_L)$ with $n_k = |V_k|$ then

$$\mathcal{H} = \mathcal{H}_L \circ \cdots \circ \mathcal{H}_1, \qquad \mathcal{H}_k = \mathcal{L}_{d_{k,1}} \times \cdots \times \mathcal{L}_{d_{k,n_k}},$$

$$VCD\left(\mathcal{L}_{d_{k,1}}\right) = d_k + 1$$

i.e., $\mathcal{H}_k \subset \{h \colon \mathbb{R}^{n_{k-1}} \to \{-1, +1\}^{n_k}\}$ and $d_{k,j} \leq n_{k-1}$ denotes the number of incoming edges at node $v_{k,j}$

- We then can use the following:

---

**Proposition 3.37:**

Let either

1. $\mathcal{H} := \mathcal{H}_2 \circ \mathcal{H}_1$ given $\mathcal{H}_1 \subseteq \mathcal{Y}^{\mathcal{X}}$ and $\mathcal{H}_2 \subseteq \mathcal{Z}^{\mathcal{Y}}$,
2. or $\mathcal{H} := \mathcal{H}_1 \times \mathcal{H}_2$ given $\mathcal{H}_i \subseteq \mathcal{Y}_i^{\mathcal{X}}$ for $i = 1, 2$.

Then we have

$$\tau_{\mathcal{H}}(m) \leq \tau_{\mathcal{H}_1}(m) \cdot \tau_{\mathcal{H}_2}(m), \qquad m \in \mathbb{N}.$$

growth function

---

Proof of Lemma 3,37:

let $H = H_1 \circ H_2$ then for any set $M = \{x_1, \ldots x_m\} \subset X$

we have $|H_m| = |\{[h(x_1), \ldots h[x_m]] \cdot h \in H\}|$

$= |\{[h_2(h_1(x_1)), \ldots, h_2(h_1(x_m))] : h_1 \in H_1, \; h_2 \in H_2\}|$

$= |\bigcup_{y \in H_{1,m}} [h_2(y_1), \ldots, h_2(y_m)] : h_2 \in H_2|$

$\underbrace{\qquad}_{\tau_{H_2}(m)}$

$\Rightarrow$

$\boxed{\begin{array}{l} \sup_{M \subset Z[m] = m} |H_m| \quad \text{or} \\[4pt] \leq \overline{\tau_{H_1}(m)} \circ \overline{\tau_{H_2}(m)} \end{array}}$

$\leq |H_{1,m}| \cdot \overline{\tau_{H_2(m)}} \leq \overline{\tau_{H_1}(m)} \circ \overline{\tau_{H_2(m)}}$  $\#$

Same reasoning can be applied to :

$$\mathcal{H} = H_1 \propto H_2 = \left\{ h = \begin{pmatrix} h_1 \\ h_2 \end{pmatrix} ; h_1 \in H_1, h_2 \in H_2 \right)$$