# SUMMARY of math framework of ML (Last class)

- $P$: an unknown distribution on $X \times Y$ ;    ∀ sets

- We see training data: $(X_1, Y_1), \ldots, (X_n, Y_n) \sim P$ iid. Goal = oracle $X \to Y$    ↗ label

- Choose a **hypothesis class** $\mathcal{H}$    (functions $X \to Y$)

- $\forall h \in \mathcal{H}$, **Risk**, a.k.a. "test error":
    ↙ loss function, e.g. $\mathbb{E}(h(X) - Y)^2$ (quadratic loss)

$$R(h) := \mathbb{E}\, \ell\big(h(X), Y\big) \qquad h^* = \underset{h \in \mathcal{H}}{\arg\min}\, R(h). \qquad \underline{\text{Not computable}}$$

- **Empirical risk** a.k.a. **training error**

- $R_n(h) := \dfrac{1}{n} \sum\limits_{i=1}^{n} \ell\big(h(X_i), Y_i\big). \qquad h_n^* := \underset{h \in \mathcal{H}}{\arg\min}\, R_n(h). \qquad \underline{\text{Computable}}$

- **ERM** algorithm:
  - ① Training: for input data $(X_1, Y_1), \ldots, (X_n, Y_n)$; compute $h_n^*$.
  - ② Prediction: on query $X$, output $h_n^*(X)$    "oracle"

<u>Lem</u> (Generalization error)
$$R(h_n^*) \leq R(h^*) + 2 \sup_{h \in \mathcal{H}} \big| R_n(h) - R(h) \big| \qquad \ominus$$

↑ test error of ERM    ↖ Best possible error (with ∞ data)

$$\ominus \quad 2 \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^{n} \Big( \underbrace{\ell(h(X_i), Y_i) - \mathbb{E}\,\ell(h(X_i), Y_i)}_{Z_i(h) \text{ iid mean 0 r.v's}} \Big) \right| = 2 \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^{n} Z_i(h) \right|$$

$$= \text{stochastic process, called "empirical process"}$$

- For binary classification, $\ell(\cdot,\cdot) \in \{0,1\}$ $\Rightarrow$ $|z_i(h)| \leq 1$

- $P\left\{ \left| \frac{1}{n} \sum_{i=1}^{n} z_i(h) \right| > t \right\} = P\left\{ \left| \frac{1}{\sqrt{n}} \sum_{i=1}^{n} z_i(h) \right| > t\sqrt{n} \right\} \leq 2\exp\left( -\frac{(t\sqrt{n})^2}{2} \right)$ (*)

<span style="color:red">multiply both sides by $\sqrt{n}$ to scale like in CLT</span>

<span style="color:red">iid mean 0, bdd by 1</span>

<span style="color:red">Hoeffding's inequality (Lec. 7, Sep. 16)</span>

- Union Bound:

$$P\left\{ \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^{n} z_i(h) \right| > t \right\} \leq \sum_{h \in \mathcal{H}} P\left\{ \left| \frac{1}{n} \sum_{i=1}^{n} z_i(h) \right| > t \right\}$$

<span style="color:red">$\leftarrow$ assume $\mathcal{H}$ is finite</span>

<span style="color:red">"$\exists\, h \in \mathcal{H}$" = union</span>

$$\overset{(*)}{\leq} |\mathcal{H}| \cdot 2\exp\left( -t^2 n/2 \right) = 2\exp\left( \log|\mathcal{H}| - t^2 n/2 \right)$$

$$\leq 0.01 \quad \text{if} \quad t = C\sqrt{\frac{\log|\mathcal{H}|}{n}}.$$

$\Rightarrow$ We proved:

THM (Generalization Bound) If the hypothesis class $\mathcal{H}$ is finite,

$$R(h_n^*) \leq R(h^*) + C\sqrt{\frac{\log|\mathcal{H}|}{n}} \qquad \text{with prob.} \geq 0.99.$$

<span style="color:red">ERM's test error</span>   <span style="color:red">best possible error</span>   $\overset{\wedge}{0.01}$ if $n \geq C' \log|\mathcal{H}|$.

Hence ‖ the ERM algorithm generalizes well from

$$n \sim \log|\mathcal{H}|$$

training data points.

- Good: logarithmic in $|\mathcal{H}|$
- Bad: most hypothesis classes are infinite.

Can $\log|\mathcal{H}|$ be replaced by some "complexity" of an infinite $\mathcal{H}$?

Yes: VC dimension

— 2 —

# VC DIMENSION    (Вапник-Червоненкис)

- Heuristically:    $vc(\mathcal{H}) = $ largest #(data $\mathcal{H}$ overfits)

<span style="color:red">↙ i.e functions $h: X \to \{0,1\}$</span>

> **Def** Let $\mathcal{H}$ be any collection of Boolean functions on a set $X$.
> We say that $\mathcal{H}$ overfits, or "shatters" a subset $\{x_1, ..., x_d\} \subset X$
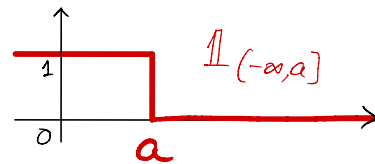> if $\forall$ labels $y_1, ..., y_d \in \{0, 1\}$   $\exists\, h \in \mathcal{H}$ such that
> $$h(x_i) = y_i \qquad \forall i = 1, ..., d.$$
> The vc dimension of $\mathcal{H}$, denoted $vc(\mathcal{H})$, is the maximal size $d$
> of a subset $\mathcal{H}$ shatters.

# Examples

$h(x)=1 \;\forall x$

**1.** $\mathcal{H} = \{ \mathbb{1} \}$  has  $\boxed{vc(\mathcal{H}) = 0}$ : it can't shatter even one point $x_i$
since   $h(x_i) = 1$.

**2.** Half-lines   $\mathcal{H} = \{ \mathbb{1}_{(-\infty, a]} : a \in \mathbb{R} \}$

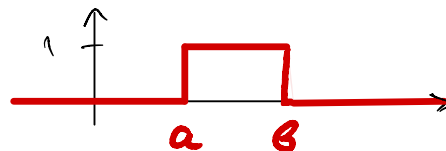$\boxed{vc(\mathcal{H}) = 1}$  Proof:


$\mathbb{1}_{(-\infty, a]}$

**HW:** $\{ \mathbb{1}_{(-\infty, a]} ; \mathbb{1}_{[b, +\infty)} \}$

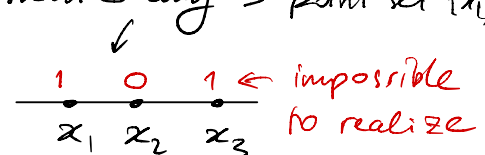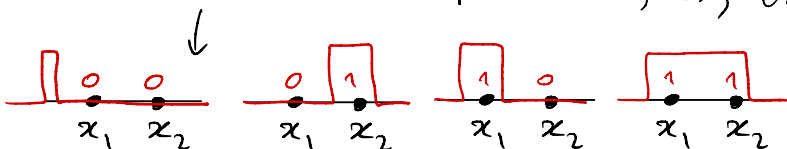$\mathcal{H}$ can shatter some 1-point set $\{x_1\}$, but can't shatter any 2-point set $\{x_1, x_2\}$



0  1  ← this label assignment
is impossible to realize
$x_1$  $x_2$

**3.** Intervals:   $\mathcal{H} = \{ \mathbb{1}_{[a, b]} : a \leq b \}$

$\boxed{vc(\mathcal{H}) = 2}$  Proof:



$\mathcal{H}$ can shatter some 2-pt set $\{x_1, x_2\}$, but can't shatter any 3-point set $\{x_1, x_2, x_3\}$



1  0  1 ← impossible
$x_1$  $x_2$  $x_3$   to realize

— 3 —

**4.** Half-planes in $\mathbb{R}^2$:

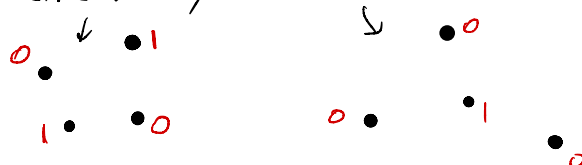$$\mathcal{H} = \left\{ \mathbb{1}_{\{a_1 x(1) + a_2 x(2) + b \geq 0\}} : a_1, a_2, b \in \mathbb{R} \right\}$$



$\boxed{vc(\mathcal{H}) = 3}$   Proof:

$\mathcal{H}$ can shatter some 3-pt set $\{x_1, x_2, x_3\}$,



but can't shatter any 4-point set $\{x_1, x_2, x_3, x_4\}$:

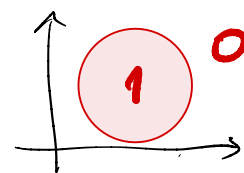$\forall$ 4-point set is like this, or like this



"convex position"

In either case, $\exists$ label assignment that is impossible to realize

**4.** Circles in $\mathbb{R}^2$:

$$\mathcal{H} = \left\{ \mathbb{1}_{\{(x(1) - a_1)^2 + (x(2) - a_2)^2 \leq r^2\}} : a_1, a_2, r \in \mathbb{R} \right\}$$
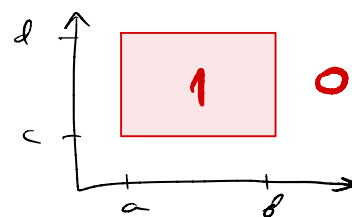
$\boxed{vc(\mathcal{H}) = 3}$   **HW ?**



**5.** Axis-aligned rectangles in $\mathbb{R}^2$:

$$\mathcal{H} = \left\{ \mathbb{1}_{(a,b) \times [c,d]} : a < b, c < d \right\}$$
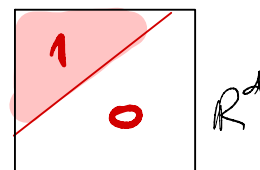
$\boxed{vc(\mathcal{H}) = 4}$



Generalizing Example 4:

**6.** Half-spaces in $\mathbb{R}^d$:

$$\mathcal{H} = \left\{ \mathbb{1}_{\{\langle w, x \rangle + b \geq 0\}} : w \in \mathbb{R}^d, b \in \mathbb{R} \right\}$$
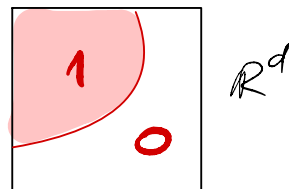
"linear classifier" (svm)

$\boxed{vc(\mathcal{H}) = d + 1}$

More generally:

**7.** Polynomial surfaces in $\mathbb{R}^d$: $\mathcal{H} = \{ \mathbb{1}_{\{p(x) \geq 0\}} : \deg(p) \leq r \}$

$$VC(\mathcal{H}) = \binom{d+r}{r} \quad \text{[Anthony 1995]}$$
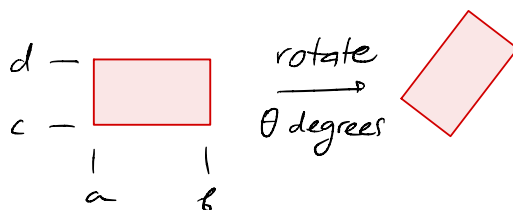
"Polynomial classifier"

$\mathbb{R}^d$

**Remark.** In all examples above, $VC(\mathcal{H}) = \#$ parameters that describe a function in $\mathcal{H}$.

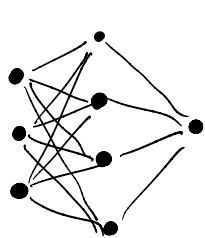- This is not true in general. For rectangles in $\mathbb{R}^2$, not necessarily axis-aligned as in Ex.5,

$$VC(\mathcal{H}) = 7$$

while the $\#$ parameters is 5: $a, b, c, d, \theta$.

$d$ — $c$ — $a \quad b$ $\xrightarrow[\theta \text{ degrees}]{\text{rotate}}$

- But heuristically, and "approximately", this is often true:

**8.** $\mathcal{H} = \{$ functions a given neural network can compute $\}$

- If activation function $= \underline{\quad\rule{1em}{0.5pt}\text{⌐}}$
  network computes a composition of linear classifiers.

  Here $\boxed{VC(\mathcal{H}) \leq C \, W \log W}$ [Cover 69; Baum-Haussler 89]

  $\#$ connections $\left(= \#\text{weights,} \atop \text{parameters}\right)$

- $\exists$ networks for which this bound is tight [Maas 94]

- If activation function is piecewise linear (e.g. ReLU) $\Rightarrow$

  $\boxed{VC(\mathcal{H}) \leq C \, W L \log W}$ [Bartlett-Harvey-Liaw-Mehrabian 2017]

  $\#$ layers

  and $\exists$ examples showing tightness.