

# PAC, Generalization and SRM

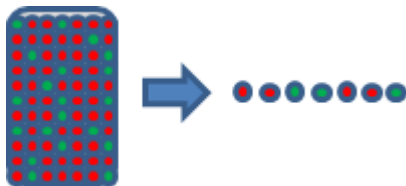
Machine Learning, A.Y. 2022/23, Padova



Fabio Aioli

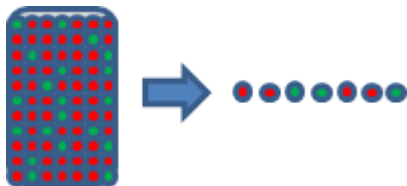
October 12th, 2022

# A simple experiment



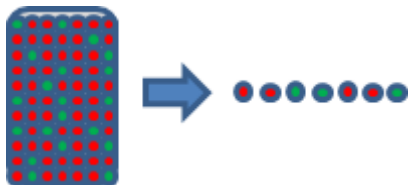
- $P(\text{red}) = \pi$
- $P(\text{green}) = 1 - \pi$
- $\pi$  is unknown
- Pick  $N$  marbles (the *sample*) from the bin, independently
- $\sigma$  = fraction of red marbles in the sample

# A simple experiment



- Does  $\sigma$  say anything about  $\pi$ ?
- Short answer... NO
- Ans: The sample can be mostly green, while the bin is mostly red
- Long answer... YES
- Ans: The sample frequency  $\sigma$  is likely close to the bin frequency  $\pi$

# What does $\sigma$ say about $\pi$

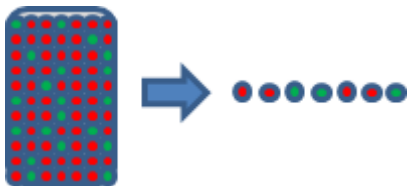


In a large sample (large  $N$ ), the value  $\sigma$  is likely close to  $\pi$  (within  $\epsilon$ )  
More formally (Hoeffding's Inequality),

$$P(\underbrace{|\sigma - \pi| > \epsilon}_{\text{bad event}}) \leq 2e^{-2\epsilon^2 N}$$

That is,  $\sigma = \pi$  is P.A.C. (Probably Approximately Correct)

# What does $\sigma$ say about $\pi$ ?



$$P(\underbrace{|\sigma - \pi| > \epsilon}_{\text{bad event}}) \leq 2e^{-2\epsilon^2 N}$$

- Valid for all  $N$  and  $\epsilon$
- Bound does not depend on  $\pi$
- Tradeoff:  $N$ ,  $\epsilon$ , and the bound
- $\sigma \approx \pi \Rightarrow \pi \approx \sigma$ , that is "  $\pi$  tends to be close to  $\sigma$  "



- In the Bin example, the unknown is  $\pi$
- In the Learning example the unknown is  $f : \mathcal{X} \rightarrow \mathcal{Y}$
- The bin is the input space  $\mathcal{X}$
- Given an hypothesis  $h$ , **green** marbles correspond to examples where the hypothesis is right, i.e.  $h(\mathbf{x}) = f(\mathbf{x})$
- Given an hypothesis  $h$ , **red** marbles correspond to examples where the hypothesis is wrong, i.e.  $h(\mathbf{x}) \neq f(\mathbf{x})$

So, for *this*  $h$ ,  $\sigma$  (empirical error) **actually generalizes** to  $\pi$  (ideal error) but... this is **verification**, not learning!

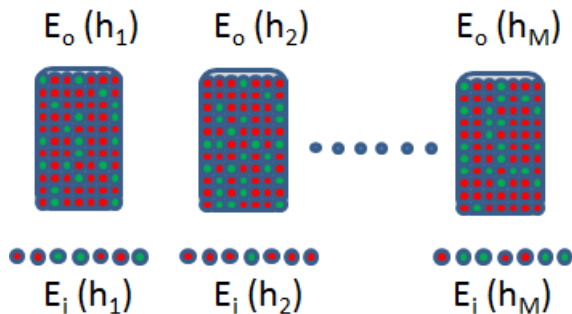


We need to choose from **multiple** hypotheses!  
 $\pi$  and  $\sigma$  depend on which  $h$  we choose

Change of notation

- in-sample error  $\sigma \rightarrow E_i(h)$
- out-of-sample error  $\pi \rightarrow E_o(h)$
- then,  $P(|E_i(h) - E_o(h)| > \epsilon) \leq 2e^{-2\epsilon^2 N}$

# Multiple Bins



Hoeffding's inequality does not directly apply here!



## Analogy: Head and Cross



- If you toss a (fair) coin 10 times, what is the probability that you will get 10 heads?
- $(0.5)^{10} = 0.0009765625 \approx 0.1\%$
- If you toss 1000 (fair) coins 10 times each, what is the probability that *some coin* will get 10 heads?
- $(1 - (1 - 0.001)^{1000}) = 0.6323045752290363 \approx 63\%$



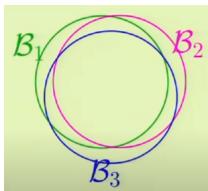
We resort to the so called **Union Bound**:

$$\begin{aligned} P[|E_i(g) - E_o(g)| > \epsilon] &\leq P[|E_i(h_1) - E_o(h_1)| > \epsilon \\ &\quad \text{or } |E_i(h_2) - E_o(h_2)| > \epsilon \\ &\quad \dots \\ &\quad \text{or } |E_i(h_M) - E_o(h_M)| > \epsilon] \\ &\leq \sum_{m=1}^M P[|E_i(h_m) - E_o(h_m)| > \epsilon] \leq 2Me^{-2\epsilon^2 N} \end{aligned}$$

Remember,  $M$  is generally very big (can be also infinite)!!

## Back to the learning problem

- **Testing:**  $P(|E_i(g) - E_o(g)| > \epsilon) \leq 2e^{-2\epsilon^2 N}$
- **Training:**  $P(|E_i(g) - E_o(g)| > \epsilon) \leq 2Me^{-2\epsilon^2 N}$

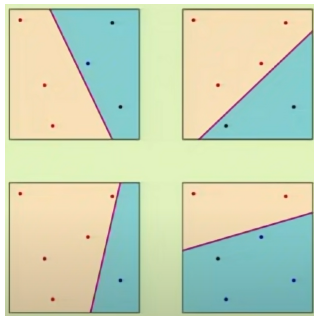


In fact  $M$  can be substituted by  $m_{\mathcal{H}}(N) \leq 2^N$  which is related to the *complexity* of the hypothesis space!

Remember that  $P(E \cup F) = P(E) + P(F) - P(E \cap F)$ .

So, when the bad events overlaps a lot (low complexity of the hypothesis space), then the value  $m_{\mathcal{H}}(N) \ll 2^N$ . What happens if only poly( $N$ )?

# Overlapping of hypotheses and Growth Function



$$m_{\mathcal{H}}(N) = \max_{\mathbf{x}_1, \dots, \mathbf{x}_N} |\mathcal{H}(\mathbf{x}_1, \dots, \mathbf{x}_N)| \leq 2^N$$

where  $|\mathcal{H}(\mathbf{x}_1, \dots, \mathbf{x}_N)|$  is the number of dichotomies we can have on  $\mathbf{x}_1, \dots, \mathbf{x}_N$  using hypotheses from  $\mathcal{H}$ .

# Measuring the complexity of a hypothesis space

## Shattering



**Shattering:** Given  $S \subset X$ ,  $S$  is shattered by the hypothesis space  $\mathcal{H}$  iff

$$\forall S' \subseteq S, \exists h \in \mathcal{H}, \text{ such that } \forall x \in S, h(x) = 1 \Leftrightarrow x \in S'$$

( $\mathcal{H}$  is able to implement all possible dichotomies of  $S$ )

# Measuring the complexity of the hypothesis space

## VC-dimension



**VC-dimension:** The VC-dimension of a hypothesis space  $\mathcal{H}$  defined over an instance space  $X$  is the size of the largest finite subset of  $X$  shattered by  $\mathcal{H}$ :

$$VC(\mathcal{H}) = \max_{S \subseteq X} |S| : S \text{ is shattered by } \mathcal{H}$$

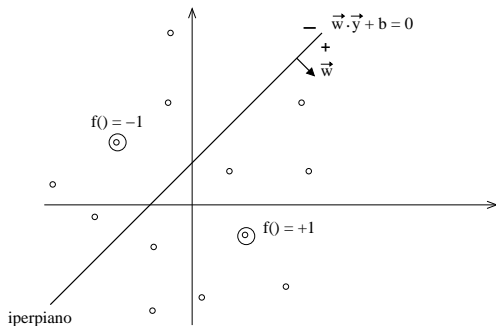
If arbitrarily large finite sets of  $X$  can be shattered by  $\mathcal{H}$ , then  $VC(\mathcal{H}) = \infty$ .



# VC-dimension: Example (1)

What is the VC-dimension of  $\mathcal{H}_1$  ?

$$\mathcal{H}_1 = \{f_{(\vec{w}, b)}(\vec{y}) \mid f_{(\vec{w}, b)}(\vec{y}) = \text{sign}(\vec{w} \cdot \vec{y} + b), \vec{w} \in \mathbb{R}^2, b \in \mathbb{R}\}$$

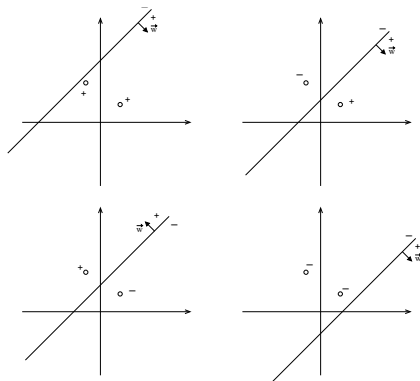




## VC-dimension: Example (2)

What is the VC-dimension of  $\mathcal{H}_1$  ?

$VC(\mathcal{H}) \geq 1$  trivial. Let consider 2 points:



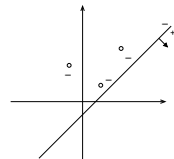
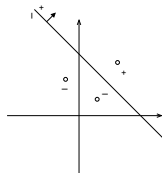
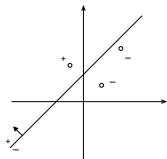
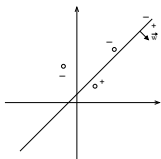
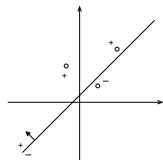
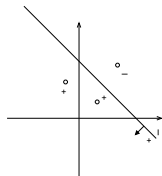
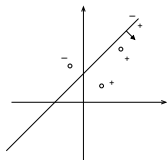
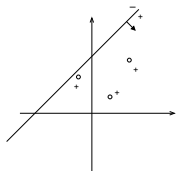




# VC-dimension: Example (3)

What is the VC-dimension of  $\mathcal{H}_1$  ?

Thus  $VC(\mathcal{H}) \geq 2$ . Let consider 3 points:



## VC-dimension: Example (4)



What is the VC-dimension of  $\mathcal{H}_1$  ?

Thus  $VC(\mathcal{H}) \geq 3$ . What happens with 4 points ?

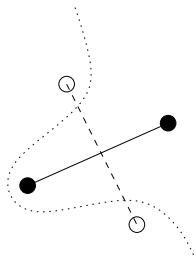


## VC-dimension: Example (5)

What is the VC-dimension of  $\mathcal{H}_1$  ?

Thus  $VC(\mathcal{H}) \geq 3$ . What happens with 4 points ? It is impossible to shatter 4 points!!

In fact there always exist two pairs of points such that if we connect the two members by a segment, the two resulting segments will intersect. So, if we label the points of each pair with a different class, a curve is necessary to separate them! Thus  $VC(\mathcal{H}) = 3$



What if  $n > 2$  ?



# Generalization Error: the VC bound

Consider a binary classification learning problem with:

- Training set  $\mathcal{S} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$
- Hypothesis space  $\mathcal{H} = \{h_\theta(\mathbf{x})\}$
- Learning algorithm  $\mathcal{L}$ , returning the hypothesis  $g = h_\theta^*$  minimizing the empirical error on  $\mathcal{S}$ , that is  $g = \arg \min_{h \in \mathcal{H}} \text{error}_{\mathcal{S}}(h)$ .

It is possible to derive an upper bound of the ideal error which is valid with probability  $(1 - \delta)$ ,  $\delta$  being arbitrarily small, of the form:

$$\text{error}(g) \leq \text{error}_{\mathcal{S}}(g) + F\left(\frac{\text{VC}(\mathcal{H})}{n}, \delta\right)$$



Let's take the two terms of the bound

- $A = \text{error}_S(g)$
- $B = F(\text{VC}(\mathcal{H})/n, \delta)$
- The term  $A$  depends on the hypothesis returned by the learning algorithm  $\mathcal{L}$ .
- The term  $B$  (often called **VC-confidence**) does not depend on  $\mathcal{L}$ . It only depends on:
  - the training size  $n$  (inversely),
  - the VC dimension of the hypothesis space  $\text{VC}(\mathcal{H})$  (proportionally)
  - the confidence  $\delta$  (inversely).



# Structural Risk Minimization

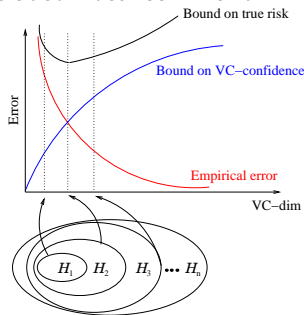
Problem: as the VC-dimension grows, the empirical risk (A) decreases, however the VC confidence (B) increases !

Because of that, Vapnik and Chervonenkis proposed a **new inductive principle**, i.e. **Structural Risk Minimization (SRM)**, which aims to minimizing the right hand of the confidence bound, so to get a tradeoff between **A** and **B**:

Consider  $\mathcal{H}_i$  such that

- $\mathcal{H}_1 \subseteq \mathcal{H}_2 \subseteq \dots \subseteq \mathcal{H}_n$
- $VC(\mathcal{H}_1) \leq \dots \leq VC(\mathcal{H}_n)$
- select the hypothesis with the smallest bound on the true risk

Example: Neural networks with an increasing number of hidden units





## Notions

- Hoeffding's Inequality
- Connection to learning
- Measuring the complexity of the hypotheses space (VC-Dimension)
- VC-Dimension of hyperplanes
- Structural Risk Minimization

## Exercises

- VC-Dimension of other hypothesis spaces, e.g. intervals in  $\mathbb{R}$  :

$$h(x) = +1 \text{ if } a \leq x \leq b, h(x) = -1 \text{ otherwise.}$$