

### Outline and References

- Outline
  - ▶ Basics [ML 7.1, 7.2]
  - ▶ PAC-Learning [ML 7.3]
  - ▶ VC Dimension [ML 7.4]
  
- References
  - ▶ Machine Learning, Mitchell [ML]



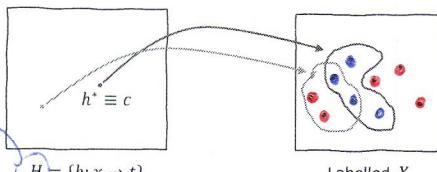
### Machine Learning - Daniele Loriacono

- What is computational learning theory? — Something that allows us to bound the test error based on the training error. (obviously under some assumptions)
- It aims at studying the general laws of inductive learning, by modeling:
    - ▶ Complexity of hypothesis space
    - ▶ Bound on training samples
    - ▶ Bound on accuracy
    - ▶ Probability of successful learning
    - ▶ ...
  - This allows to answer to questions like:
    - ▶ How many training samples do a learner need to converge (with some probability) to a successful (with some minimum accuracy) hypothesis?
    - ▶ How many training samples will be misclassified by the learner before converging to a successful hypothesis?
    - ▶ ...

### Machine Learning - Daniele Loriacono

#### (Let's Go Back to) The Big Picture

hypothesis space  
(all the points represent a function)  
(each point corresponds to a different function)



In this case each point of the hypothesis space corresponds to a decision boundary. Some points lead to good decision boundaries, others not.

- A learner ( $L$ ) wants to learn a concept ( $c$ ) that maps the data in the input space ( $X$ ) to a target ( $t$ )
- Let assume that  $L$  found an hypothesis  $h^*$  with no errors on the training data
- How many training samples of  $X$  are necessary to be sure that  $L$  actually learned the true concept, i.e.,  $h^* \equiv c$ ?



Supposing that we have no noise, if we don't have any particular assumption and we consider all the possible concepts that we can learn, unless we observe all the points, we cannot tell if we're learning the true concept or not. This is the basis of the "No Free Lunch" theorem.

### Machine Learning - Daniele Loriacono

#### «No Free Lunch» Theorems

- Let  $ACC_G(L)$  be the generalization accuracy of learner  $L$ , i.e., the accuracy of  $L$  on samples that are not in the training set.
- Let  $\mathcal{F}$  be the set of all the possible concepts  $y = f(x)$  (concepts = functions)
- For any learner  $L$  and any possible training set:

$$\frac{1}{|\mathcal{F}|} \sum_{\mathcal{F}} ACC_G(L) = \frac{1}{2}$$

▶ Proof Sketch: for every concept  $f$  where  $ACC_G(f) = 0.5 + \delta$ , exists a concept  $f'$  where  $ACC_G(f') = 0.5 - \delta$ :  $\forall x \in D, f'(x) = f(x); \forall x \in D, f'(x) \neq f(x)$

▶ Corollary: for any two learners,  $L_1$  and  $L_2$ , if  $\exists f$  where  $ACC_G(f) > ACC_G(L_1)$  then  $\exists f'$  where  $ACC_G(f') > ACC_G(L_2)$

### Machine Learning - Daniele Loriacono

## «No Free Lunch» Theorems

□ Let  $ACC_G(L)$  be the generalization accuracy of learner  $L$ , i.e., the accuracy of  $L$  on samples that are not in the training set

□ Let

□ For

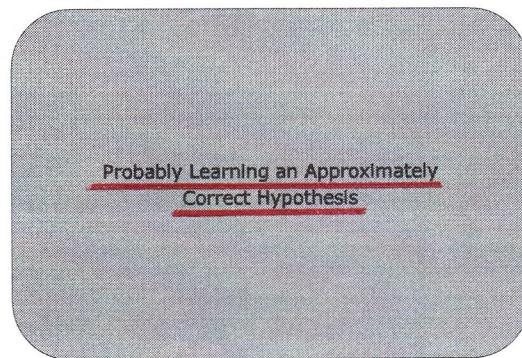
What does this mean in practice?

There is no such thing as a winner-takes-all in ML!

In ML we always operate under some assumptions!

$\exists f'$  where  $ACC_G(L_2) > ACC_G(L_1)$  except  $f'$

for example we assume that the training data is distributed like the test data (strong assumption!) We need to be aware of the assumptions we're using. We assume that what we see (training set) is representative. A violation of this assumption leads to non-negligible errors.  
ML is useless if the training set is not representative.



## Basics

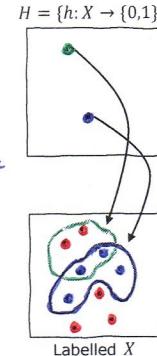
(for now we focus on binary classification)

□ Problem setting

- ▶ Let  $X$  be the instance space (input space)
- ▶ Let  $H = \{h: X \rightarrow \{0,1\}\}$  be the hypothesis space of  $L$
- ▶ Let  $C = \{c: X \rightarrow \{0,1\}\}$  be the set of all the possible target functions (concepts) we might want to learn
- ▶ Let  $D$  be training data drawn from a stationary distribution  $P(X)$  and labeled (without noise) according to a concept  $c$  we want to learn

□ A learner  $L$  outputs a hypothesis  $h \in H$  such that

$$h^* = \arg \min_{h \in H} error_{train}(h)$$



Assumptions: (big!) we want to learn something that is correlated with our inputs and it's well represented by our training set. This is formally specified by \*

$H$ : space of all the possible functions that map a point in the input space into a target (all the possible models) that we can produce

$C$ : space of all the possible functions that map a point in the input space into a target THAT WE MAY WANT TO LEARN  
(≠ from the ones we can produce / obtain in  $H$ )  
We have no guarantee that our learner  $L$  is able to learn any function

## How do we compute the error?

□ We compute the error of an hypothesis as the probability of misclassifying a sample:

$$error_D(h) = \Pr_{x \in D} [h(x) \neq c(x)] = \frac{1}{|D|} \sum_{x \in D} I(h(x) \neq c(x))$$

$D$  is the training data

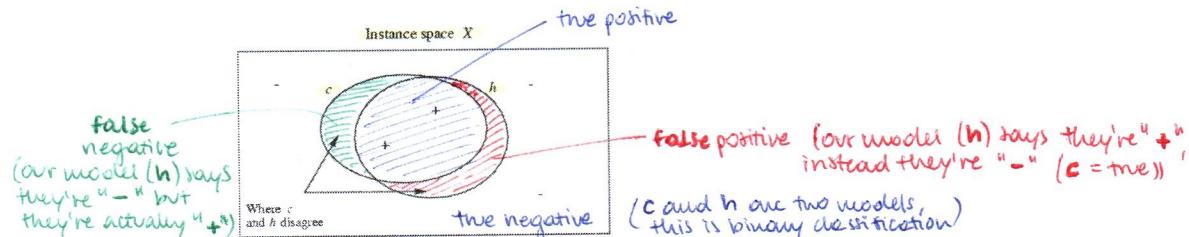
□ This is the training error, instead we are interested in the true error of  $h$ :

$$error_{true}(h) = \Pr_{x \sim P(X)} [h(x) \neq c(x)]$$

$P(X)$  is the input space distribution

Our goal is to find a model that is accurate on points that we never saw

## How do we compute the error?



□ But we have to remember that...

$$error_{true}(h) = \Pr_{x \sim P(X)} [h(x) \neq c(x)]$$

## What now?

- We say that  $h$  overfits the training data if  $\text{error}_{\text{true}} > \text{error}_{\mathcal{D}}$  ...
- ... but can we bound  $\text{error}_{\text{true}}$  given  $\text{error}_{\mathcal{D}}$ ?
- Let assume...
  - ▶  $\text{error}_{\text{true}}$  is the probability of making a mistake on a sample
  - ▶ we can compute  $\text{error}_{\mathcal{D}}$  that is the average error probability on  $\mathcal{D}$
  - ▶ assuming a Bernoulli distribution for the error probability, the 95% CI is:
$$\text{error}_{\text{true}}(h) = \text{error}_{\mathcal{D}}(h) \pm 1.96 \sqrt{\frac{\text{error}_{\mathcal{D}}(h)(1 - \text{error}_{\mathcal{D}}(h))}{n}}$$
- Is this correct? No! Because  $\mathcal{D}$  is the training data and not independent of  $h$
- So, we need to bound the error under more strict assumptions

It could be cool, but

what we are observing ( $\text{error}_{\mathcal{D}}$ ) is the empirical average on the training set  $\mathcal{D}$

let's start in a simpler setting :

because of this we're introducing a bias in the estimate of  $\text{error}_{\mathcal{D}}$  (the probability of making an error on a point of  $\mathcal{D}$  is not equal to the probability of making an error on a point  $\notin \mathcal{D}$ )

## Machine Learning - Daniela Lolaco

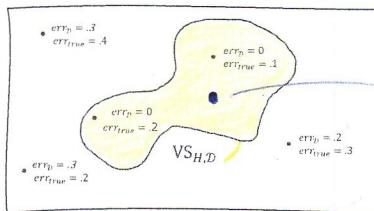
### Version Space

- A hypothesis  $h$  is **consistent** with a training dataset  $\mathcal{D}$  of the concept  $c$  if and only if  $h(x) = c(x)$  for each training sample in  $\mathcal{D}$
- $\text{Consistent}(h, \mathcal{D}) \stackrel{\text{def}}{=} \forall \langle x, c(x) \rangle \in \mathcal{D}, h(x) = c(x)$
- The version space,  $VS_{H, \mathcal{D}}$  with respect to hypothesis space  $H$  and labelled dataset  $\mathcal{D}$ , is the subset of hypotheses in  $H$  consistent with  $\mathcal{D}$
- $VS_{H, \mathcal{D}} \stackrel{\text{def}}{=} \{h \in H | \text{Consistent}(h, \mathcal{D})\}$
- From now on, we consider only **consistent learners**, that always output a **consistent hypothesis**, i.e., an hypothesis in  $VS_{H, \mathcal{D}}$ , assuming it is not empty
- Can we bound the  $\text{error}_{\text{true}}$  of a consistent learner?

## Machine Learning - Daniela Lolaco

### Version Space (2)

#### Hypothesis space $H$



- If we wish to bound the  $\text{error}_{\text{true}}$  of a consistent learner, we need to find a bound for all the hypotheses in  $VS_{H, \mathcal{D}}$

every point here has:  $\begin{cases} \text{error}_{\mathcal{D}} = 0 \\ \text{error}_{\text{true}} = ? \end{cases}$   
we want to bound  $\text{error}_{\text{true}}$  for all these points that have  $\text{error}_{\mathcal{D}} = 0$ .  
Clearly we cannot rely on  $\text{error}_{\mathcal{D}}$  to bound  $\text{error}_{\text{true}}$ .

## Machine Learning - Daniela Lolaco

### Bound for Consistent Learners

If the hypothesis space  $H$  is **finite** and  $\mathcal{D}$  is a sequence of  $N \geq 1$  independent random examples of some target concept  $c$ , then for any  $0 \leq \varepsilon \leq 1$ , the probability that  $VS_{H, \mathcal{D}}$  contains a hypothesis error greater than  $\varepsilon$  is less than  $|H|e^{-\varepsilon N}$

$$\Pr(\exists h \in H : \text{error}_{\mathcal{D}}(h) = 0 \wedge \text{error}_{\text{true}}(h) \geq \varepsilon) \leq |H|e^{-\varepsilon N}$$

#### Proof

$$\begin{aligned} &\Pr((\text{error}_{\mathcal{D}}(h_1) = 0 \wedge \text{error}_{\text{true}}(h_1) \geq \varepsilon) \vee \dots \vee (\text{error}_{\mathcal{D}}(h_{|VS_{H, \mathcal{D}}|}) = 0 \wedge \text{error}_{\text{true}}(h_{|VS_{H, \mathcal{D}}|}) \geq \varepsilon)) \\ &\leq \sum_{h \in VS_{H, \mathcal{D}}} \Pr(\text{error}_{\mathcal{D}}(h) = 0 \wedge \text{error}_{\text{true}}(h) \geq \varepsilon) \quad (\text{Union bound}) \\ &\leq \sum_{h \in VS_{H, \mathcal{D}}} \Pr(\text{error}_{\mathcal{D}}(h) = 0 | \text{error}_{\text{true}}(h) \geq \varepsilon) \quad (\text{Bound using Bayes' rule}) \\ &\leq \sum_{h \in VS_{H, \mathcal{D}}} (1 - \varepsilon)^N \quad (\text{Bound on individual } h) \\ &\leq |H|(1 - \varepsilon)^N \quad (|VS_{H, \mathcal{D}}| \leq |H|) \\ &\leq |H|e^{-\varepsilon N} \quad (1 - \varepsilon \leq e^{-\varepsilon}, \text{ for } 0 \leq \varepsilon \leq 1) \end{aligned}$$

We're bounding the probability that the version space contains an hypothesis with an error (true) greater than  $\varepsilon$ . In this way, if our learner is inside the version space also its error is bounded.

## Machine Learning - Daniela Lolaco

### What does it mean in practice?

- Let say that  $\delta$  is the probability to have  $\text{error}_{\text{true}} > \varepsilon$  for a consistent hypothesis:

$$|H|e^{-\varepsilon N} \leq \delta \quad \text{confidence level that we want}$$

- We can bound  $N$  after setting  $\varepsilon$  and  $\delta$ :

$$N \geq \frac{1}{\varepsilon} \left( \ln |H| + \ln \left( \frac{1}{\delta} \right) \right)$$

- We can bound  $\varepsilon$  after setting  $N$  and  $\delta$ :

$$\varepsilon \geq \frac{1}{N} \left( \ln |H| + \ln \left( \frac{1}{\delta} \right) \right)$$

## Machine Learning - Daniela Lolaco

the probability of correct classification is  $1 - \text{prob. of incorrect classification}$   
 $\rightarrow \leq 1 - \varepsilon$  (on  $\mathcal{D}$ )  
We need this to happen to every sample:  
 $\rightarrow \leq (1 - \varepsilon)^N$

$$\begin{aligned} &\Pr(\text{error}_{\mathcal{D}}(h) = 0 | \text{error}_{\text{true}}(h) \geq \varepsilon) = \\ &= \prod_{j=1}^N \Pr(\text{error}_{x_j}(h) = 0 | \text{error}_{\text{true}}(h) \geq \varepsilon) \\ &= \prod_{j=1}^N [1 - \Pr(\text{error}_{x_j}(h) \neq 0 | \text{error}_{\text{true}}(h) \geq \varepsilon)] \\ &\Pr(\text{error}_{x_j}(h) \neq 0 | \text{error}_{\text{true}}(h) \geq \varepsilon) \geq \varepsilon \\ &1 - \Pr(\cdot) \leq 1 - \varepsilon \\ &\prod(1 - \Pr(\cdot)) \leq \prod(1 - \varepsilon) = (1 - \varepsilon)^N \end{aligned}$$

## What does it mean in practice?

- Let say that  $\delta$  is the probability to have  $\text{error}_{\text{true}} > \varepsilon$  for a consistent hypothesis:

$$|H|e^{-\varepsilon N} \leq \delta$$

- We can bound  $N$  after setting  $\varepsilon$  and  $\delta$ :

$$N \geq \frac{1}{\varepsilon} \left( \ln |H| + \ln \left( \frac{1}{\delta} \right) \right)$$

- We can bound  $\varepsilon$  after setting  $N$  and  $\delta$ :

$$\varepsilon \geq \frac{1}{N} \left( \ln |H| + \ln \left( \frac{1}{\delta} \right) \right)$$

Can be exponential in  
#features

## Example: Conjunction of up to $N$ Boolean Literals

- Consider a classification problem

- Instance space is  $X = \langle x_1, x_2, x_3, x_4 \rangle$ , where  $x_i$  is a boolean variable
- Each hypothesis  $h$  is a rule like this:  
$$\text{if } (x_1 = 1, x_2 = ?, x_3 = 0, x_4 = 1) \text{ then } y = 1, \text{ otherwise } y = 0$$

- How many samples  $N$  are necessary to guarantee that, with a probability at least of 0.99, the error of a consistent hypothesis is not greater than 0.05?

$$N \geq \frac{1}{\varepsilon} \left( \ln |H| + \ln \left( \frac{1}{\delta} \right) \right) \Rightarrow N \geq 180$$

0.05      3<sup>4</sup>      0.01

How big is the hypothesis space?  
We have 3 options for all the 4 variables  $\Rightarrow 3^4$   
(All the rules are achieved by combining: 0, 1, ?)

$$\begin{aligned} x_1 &= 0, 1, ? \\ x_2 &= 0, 1, ? \\ x_3 &= 0, 1, ? \\ x_4 &= 0, 1, ? \end{aligned}$$

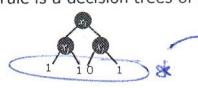
- How does it scale with respect to the number of variables ( $M$ )?

- $M=10 \rightarrow N \geq 312$  ( $|H| = 3^{10}$ )
- $M=100 \rightarrow N \geq 2290$  ( $|H| = 3^{100}$ )

## Example: Decision Tree (depth=2)

- Consider a classification problem

- Instance space is  $X = \langle x_1, \dots, x_M \rangle$ , where  $x_i$  is a boolean variable
- Each hypothesis  $h$  is a rule is a decision tree of depth 2 using only two variables:



for each we have  $\frac{M(M-1)}{2}$  possibility of choosing  $x_i$  and  $x_j$ . Moreover, we have the possibility of  $*$   $\Rightarrow 2^4 = 16$

- How many samples  $N$  are necessary to guarantee that, with a probability at least of 0.99, the error of a consistent hypothesis is not greater than 0.05?

$$N \geq \frac{1}{\varepsilon} \left( \ln |H| + \ln \left( \frac{1}{\delta} \right) \right)$$

0.05      0.01

$\frac{M(M-1)}{2} = 16$

## Probably Learning an Approximately Correct Hypothesis

- Considering a class  $C$  of possible target concepts defined over an instance space  $X$  with an encoding length  $M$ , and a learner  $L$  using an hypothesis space  $H$  we define:

$C$  is PAC-learnable by  $L$  using  $H$  if for all  $c \in C$ , for any distribution  $P(X)$ ,  $\varepsilon$  (such that  $0 < \varepsilon < 1/2$ ), and  $\delta$  (such that  $0 < \delta < 1/2$ ), learner  $L$  will with a probability at least  $(1 - \delta)$  output a hypothesis  $h \in H$  such that  $\text{error}_{\text{true}}(h) \leq \varepsilon$ , in time that is polynomial in  $1/\varepsilon$ ,  $1/\delta$ ,  $M$ , and  $\text{size}(c)$ .

- So, PAC-learnability is only about computational complexity? What about the complexity with respect to the number of training samples  $N$ ?

- A sufficient condition to prove PAC-learnability is proving that a learner  $L$  requires only a polynomial number of training examples, and processing per example is polynomial.

## Agnostic Learning

- So far, we assumed that  $c \in H$ , or at least that  $\text{VS}_{H,D}$  is not empty, and the learner  $L$  will always output a hypothesis  $h$  such that  $\text{error}_D(h) = 0$
- But in general (agnostic) learner will output a hypothesis  $h$  such that  $\text{error}_D(h) > 0$
- Can we bound  $\text{error}_{\text{true}}(h)$  given  $\text{error}_D(h)$ ?

If the hypothesis space  $H$  is finite and  $D$  is a sequence of  $N \geq 1$  i.i.d. examples of some target concept  $c$ , then for any  $0 \leq \varepsilon \leq 1$ , and for any learned hypothesis  $h$ , the probability that  $\text{error}_{\text{true}}(h) - \text{error}_D(h) > \varepsilon$  is less than  $|H|e^{-2N\varepsilon^2}$

$$\Pr(\exists h \in H : \text{error}_{\text{true}}(h) > \text{error}_D(h) + \varepsilon) \leq |H|e^{-2N\varepsilon^2}$$

We try to bound not anymore the true error but the difference between the true error and the target error  
 $\Rightarrow$  we're trying to compute the IP that exists an hypothesis in  $H$  that has a true error that is greater than the training error of a given threshold  $\varepsilon$

- Additive Hoeffding Bound: let  $\hat{\theta}$  be the empirical mean of  $N$  i.i.d. Bernoulli random variables with mean  $\theta$ :

$$\Pr(\theta > \hat{\theta} + \varepsilon) \leq e^{-2N\varepsilon^2}$$

$\Pr(\text{true mean} > \text{empirical mean} + \varepsilon) \leq e^{-2N\varepsilon^2}$

- So for any single hypothesis  $h$ :

$$\Pr(\text{error}_{\text{true}}(h) > \text{error}_{\mathcal{D}}(h) + \varepsilon) \leq e^{-2N\varepsilon^2}$$

translated in our case

- As we want this to be true for all the hypothesis in  $H$ :

$$\Pr(\exists h \in H : \text{error}_{\text{true}}(h) > \text{error}_{\mathcal{D}}(h) + \varepsilon) \leq |H|e^{-2N\varepsilon^2}$$

## Bounds for Agnostic Learning

- Similarly to what done before, we can bound the sample complexity:

$$N \geq \frac{1}{2\varepsilon^2} \left( \ln |H| + \ln \left( \frac{1}{\delta} \right) \right)$$

- We can also bound the true error of the hypothesis as:

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\mathcal{D}}(h) + \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2N}}$$

- We found the bias and variance decomposition we previously saw in the course!

we assume no noise  $\Rightarrow$  no irreducible error. The role of the bias is played by  $\text{error}_{\mathcal{D}}(h)$  and the variance is represented by  $\sqrt{\frac{1}{2N} (\ln |H| + \ln \frac{1}{\delta})}$ .

## PAC-Learning with Infinite Hypotheses Spaces

- Previously we found this PAC-Learning bound for the number of samples:

$$N \geq \frac{1}{\varepsilon} \left( \ln |H| + \ln \left( \frac{1}{\delta} \right) \right)$$

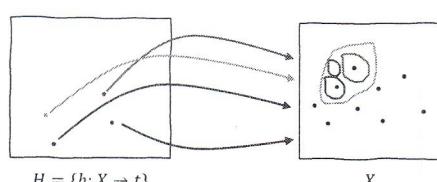
- If  $|H|$  is infinite, what does this mean? What can we use instead of  $|H|$ ?

- The answer is the largest subset of  $X$  for which  $|H|$  can guarantee a zero training error (regardless of the target function  $c$ )

- We call VC dimension the size of this subset

We see from the variance term that in order to afford a large hypothesis space ( $|H|$ ) we need a large amount of data ( $N$ ), so that variance doesn't explode. The assumption (strong) that we're making is that  $|H| < \infty$ . However, in most of the cases the hypothesis space is not finite.

## Intuition Behind Using VC Dimension



→ We have  $N$  samples: how many concepts we can learn?  $2^N$ , because each point can be classified as +1/-1. (or 0/1). Assuming  $|H| = 2^N$  = same dimension as the number of possible concepts, we can always find  $h \in H$  s.t.  $\text{error}_{\mathcal{D}}(h) = 0$ .

- Let assume that  $|X|=N$ , how big is  $|C|$ ?  
 □ Assuming  $|H| = 2^N$ , we can always find  $h \in H$  with  $\text{error}_{\mathcal{D}}(h)=0$   
 □ Does  $\text{error}_{\mathcal{D}}(h)$  tells something more on the error on other samples in  $X$ ?  
 □ What happens instead if with  $H$  we can classify correctly no more than 2 training samples?

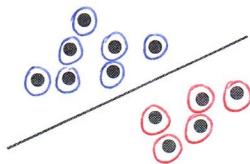
## VC Dimension

- We define a **dichotomy** of a set  $S$  of instances as a partition of  $S$  into two disjoint subsets, i.e., labeling each instance in  $S$  as positive or negative  
 □ We say that a set of instances  $S$  is **shattered** by hypothesis space  $H$  if and only if for every dichotomy of  $S$  there exists some hypothesis in  $H$  consistent with this dichotomy  
 □ The **Vapnik-Chervonenkis dimension**,  $\text{VC}(H)$ , of hypothesis space  $H$  over instance space  $X$ , is the largest finite subset of  $X$  shattered by  $H$

] if and only if any labelling can be learned in  $H$

### Example: VC dimension of linear classifier

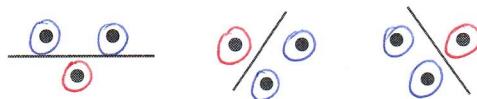
- What about a linear classifier in 2D input space?



What is the VC dimension of this linear classifier?

### Example: VC dimension of linear classifier

- What about a linear classifier in 2D input space?

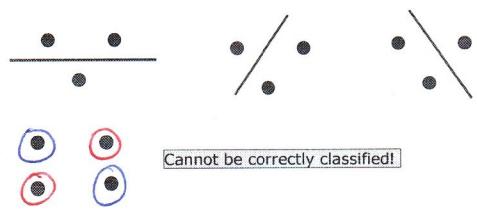


→ Also extreme case: (3 points equally labelled)

If we consider 3 points we're sure that any choice of labelling (so, any dichotomy) allows to find (through a linear classifier) a consistent hypothesis.  
Is this possible with an arbitrary number of points?

### Example: VC dimension of linear classifier

- What about a linear classifier in 2D input space?

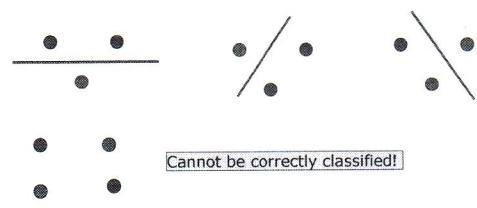


Cannot be correctly classified!

[ Not any subset of 4 points can guarantee a linear model that classifies all the points correctly ]

### Example: VC dimension of linear classifier

- What about a linear classifier in 2D input space?



Cannot be correctly classified!

- A linear classifier in a 2D input space has  $VC(h)=3$

- We can prove that a linear classifier in M-D input space has  $VC(h)=M+1$  !

For this reason, a linear classifier in a 2D space has  $VC(h)=3$ . If we choose 4 points (any) we won't be able to always classify all of them correctly.

### VC Dimension

- We define a **dichotomy** of a set  $S$  of instances as a partition of  $S$  into two disjoint subsets, i.e., labeling each instance in  $S$  as positive or negative
- We say that a set of instances  $S$  is **shattered** by hypothesis space  $H$  if and only if for every dichotomy of  $S$  there exists some hypothesis in  $H$  consistent with this dichotomy
- The **Vapnik-Chervonenkis dimension**,  $VC(H)$ , of hypothesis space  $H$  over instance space  $X$ , is the largest finite subset of  $X$  shattered by  $H$
- If an arbitrarily large set of  $X$  can be shattered by  $H$ ,  $VC(H)=\infty$
- If  $|H| < \infty$  then  $VC(H) \leq \log_2(|H|)$ 
  - If  $VC(H) = d$  it means there are in  $H$  at least  $2^d$  hypotheses to label  $d$  instances
  - Thus,  $|H| \geq 2^d$

- How many randomly drawn examples suffice to guarantee that any hypothesis that perfectly fits the training data is probably  $(1 - \delta)$  approximately  $(\epsilon)$  correct ?

$$N \geq \frac{1}{\epsilon} \left( \ln |H| + \ln \left( \frac{1}{\delta} \right) \right)$$



$$N \geq \frac{1}{\epsilon} (8VC(H) \log_2(13/\epsilon) + 4 \log_2(2/\delta))$$

## Agnostic Learning: VC Bounds

- With probability at least  $(1 - \delta)$  every  $h \in H$  satisfies the following inequality:

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\mathcal{D}}(h) + \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2N}}$$



$$\text{error}_{\text{true}}(h) \leq \text{error}_{\mathcal{D}}(h) + \sqrt{\frac{VC(H)(\ln \frac{2N}{VC(H)} + 1) + \ln \frac{4}{\delta}}{N}}$$