

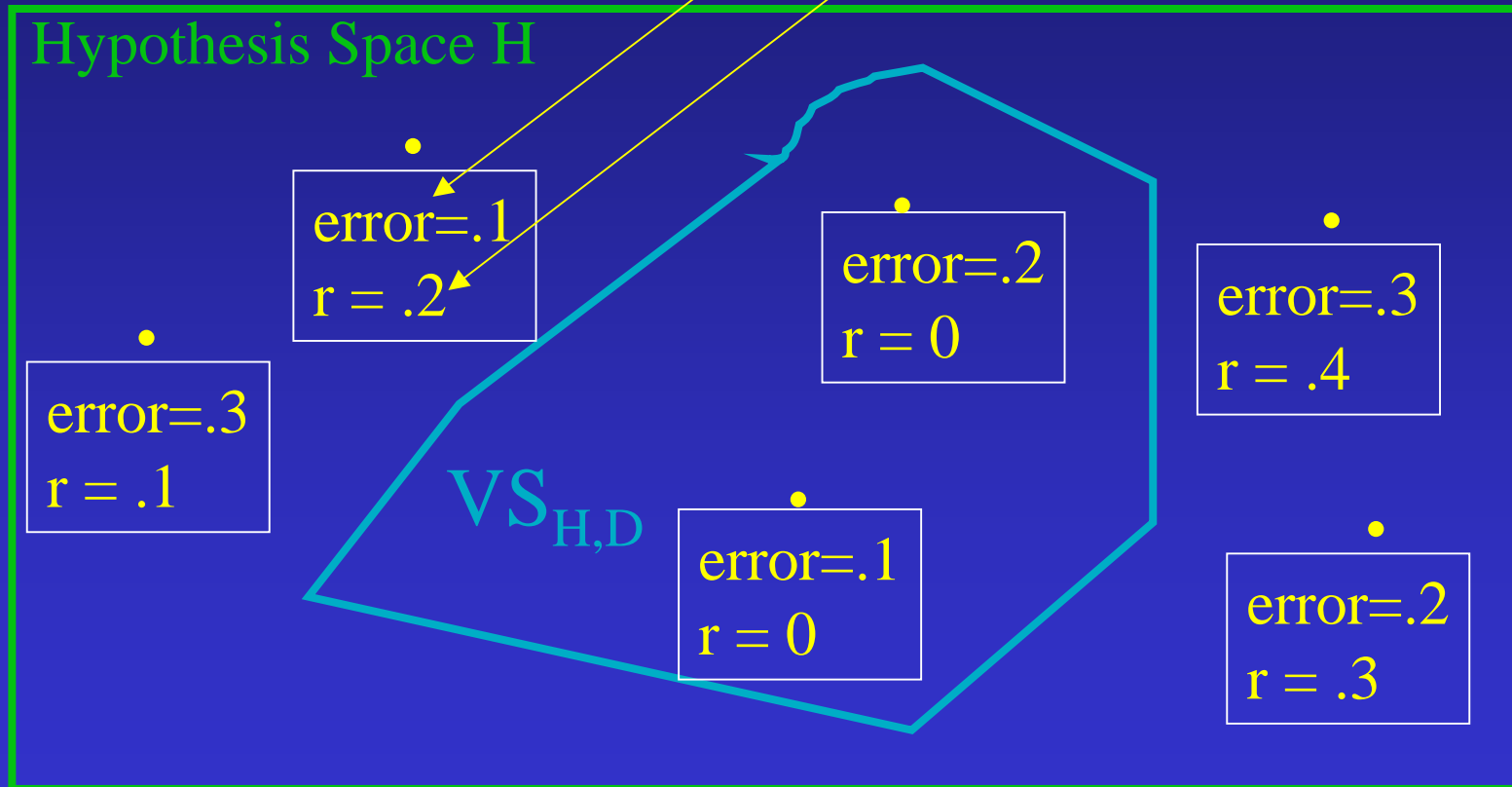
---

# Computational Learning Theory (VC Dimension)

1. Difficulty of machine learning problems
2. Capabilities of machine learning algorithms

# Version Space with associated errors

error is the true error,  
r is the training error



## Number of training samples required

$$m \geq \frac{1}{\varepsilon} (\ln |H| + \ln(1/\delta))$$

- With probability  $1-\delta$ ,
- every hypothesis in  $H$  having zero training error will have a true error of at most  $\varepsilon$
- Sample complexity for PAC learning grows as the logarithm of the size of the hypothesis space

## Disadvantage of Sample Complexity for finite hypothesis spaces

- Bound in terms of  $|H|$  has two disadvantages
  - Weak bounds

$$\delta \geq |H| e^{-\epsilon m}$$

- Bound grows linearly with  $|H|$
- For large  $H$ ,  $\delta > 1$ , or Probability at least  $(1-\delta)$  is negative!

$$m \geq \frac{1}{\epsilon} (\ln |H| + \ln(1/\delta))$$

- Can overestimate number of samples required
- Cannot be applied in case of infinite hypotheses

## Sample Complexity for infinite hypothesis spaces

- Another measure of the complexity of  $H$  called Vapnik-Chervonenkis dimension, or  $VC(H)$ 
  - We will use  $VC(H)$  instead of  $|H|$
  - Results in tighter bounds
  - Allows characterizing sample complexity of infinite hypothesis spaces and is fairly tight

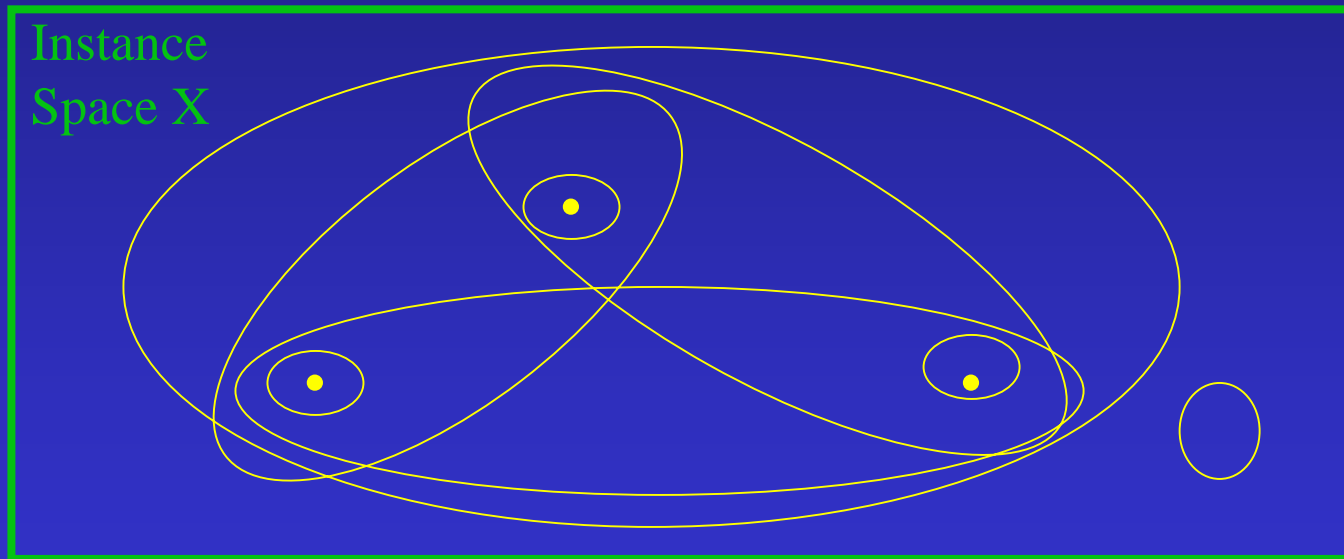
## VC Dimension

- VC Dimension is a property of a set of functions  $\{ f(\alpha) \}$
- It can be defined for various classes of functions
- Bounds that relate the capacity of a learning machine and its performance

## Shattering a Set of Instances

- Complexity of hypothesis space is measured
  - not by no. of distinct hypotheses  $|H|$
  - but by no. of distinct instances from  $X$  that can be completely discriminated using  $H$
- Definition:
  - 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

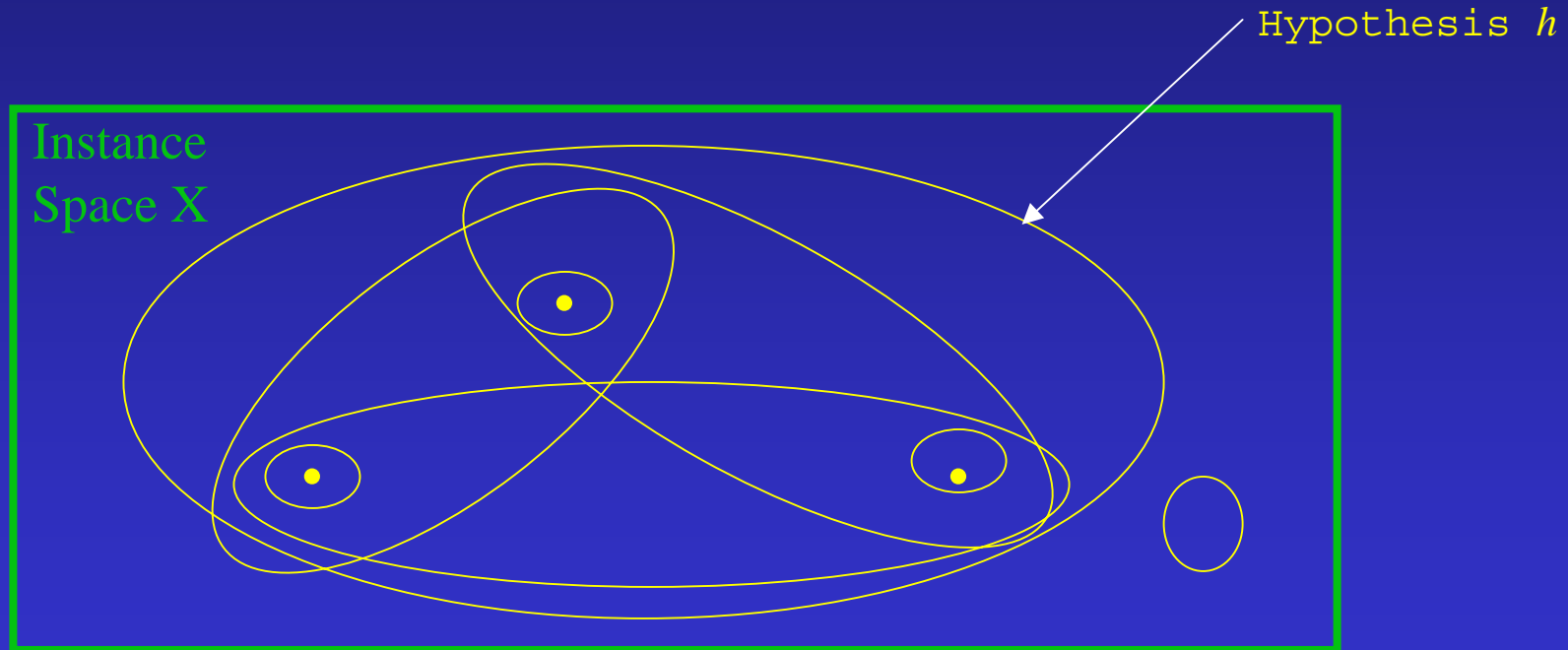
# Shattering a Set of A set of three instances by eight hypotheses





## Shattering a Set of Instances

- $S$  is a subset of  $X$  ( three instances in example below )
- into two subsets



- Ability of  $H$  to shatter a set of instances is its capacity to represent target concepts defined over these instances

## Vapnik-Chervonenkis Dimension

- Definition:  
 $VC(H)$  of hypothesis space  $H$  defined over instance space  $X$  is the size of the largest finite subset of  $X$  shattered by  $H$
- If arbitrarily large finite sets of  $X$  can be shattered by  $H$ , then  $VC(H) = \text{infinity}$

## VC Dimension

**Definition:** The VC-dimension of  $H$  is equal to the maximal number  $d$  of examples that can be split into two sets in all  $2^d$  ways using functions from  $H$  (shattering).

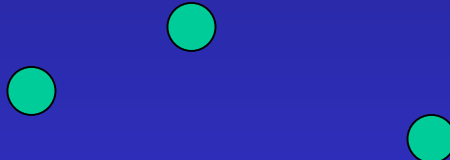
	$x_1$	$x_2$	$x_3$	...	$x_d$
$h_1$	+	+	+	...	+
$h_2$	-	+	+	...	+
$h_3$	+	-	+	...	+
$h_4$	-	-	+	...	+
...	...	...	...	...	...
$h_N$	-	-	-	...	-

## Examples to illustrate VC Dimension

1. Instance space:  $X =$  set of real numbers  $\mathbb{R}$   
 $H$  is the set of intervals on the real number line



2. Instance space:  $X =$  points on the  $x$ - $y$  plane  
 $H$  is the set of all linear decision surfaces in the plane



3. Instance space:  $X =$  three Boolean literals  
 $H$  is conjunction

a    b    c

## Example 1: VC dimension of 1-dimensional intervals

- $X = R$  (e.g., heights of people)
- $H$  is the set of hypotheses of the form  $a < x < b$
- Subset containing two instances  $S = \{3.1, 5.7\}$

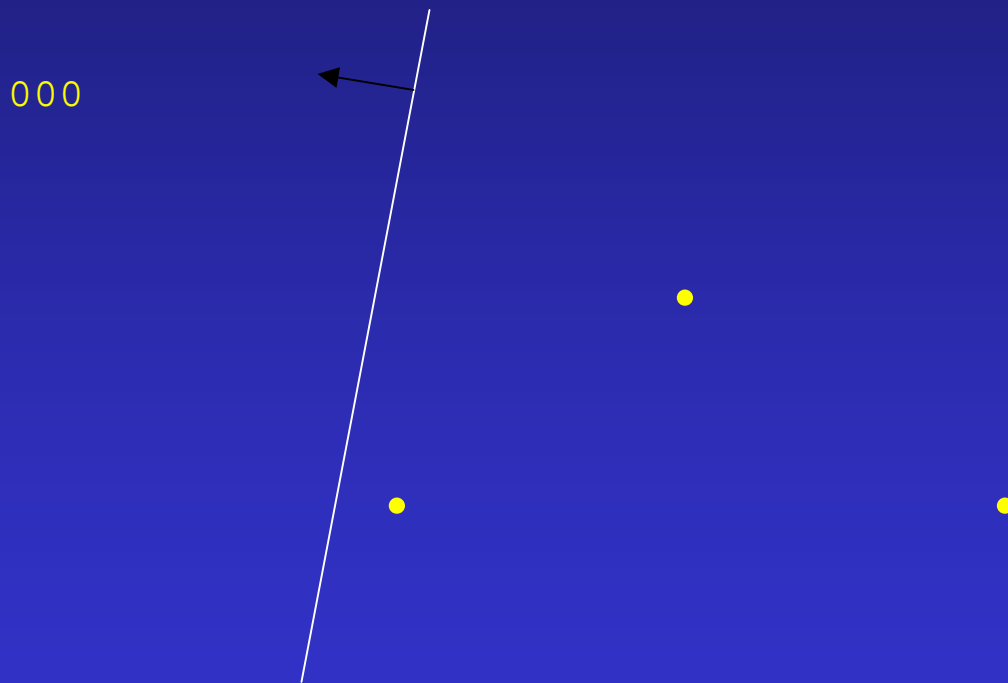


- Can  $S$  be shattered by  $H$ ?
- Yes, e.g.,  $(1 < x < 2)$ ,  $(1 < x < 4)$ ,  $(4 < x < 7)$ ,  $(1 < x < 7)$
- Since we have found a set of two that can be shattered,  $VC(H)$  is at least two
- However, no subset of size three can be shattered



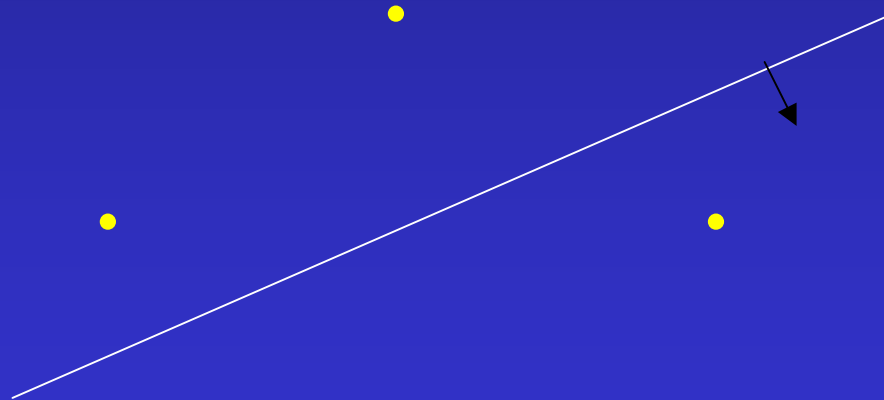
- Therefore  $VC(H) = 2$
- Here  $|H|$  is infinite but  $VC(H)$  is finite

## Example 2: VC Dimension of linear discriminants on a plane



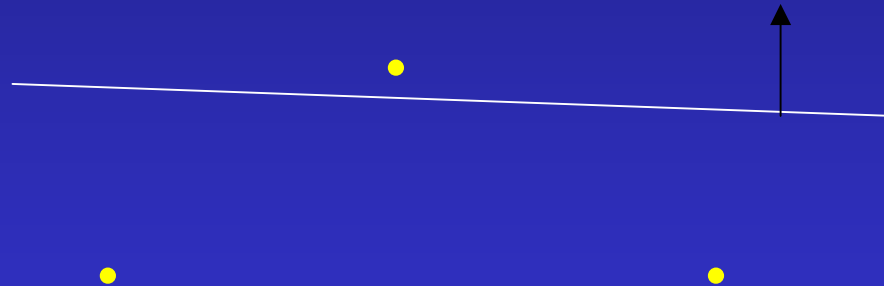
# VC Dimension of linear discriminants on a plane

001



# VC Dimension of linear discriminants on a plane

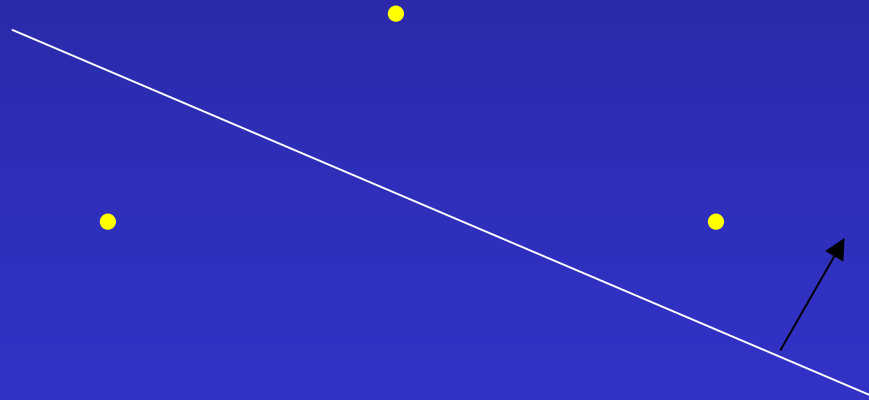
010





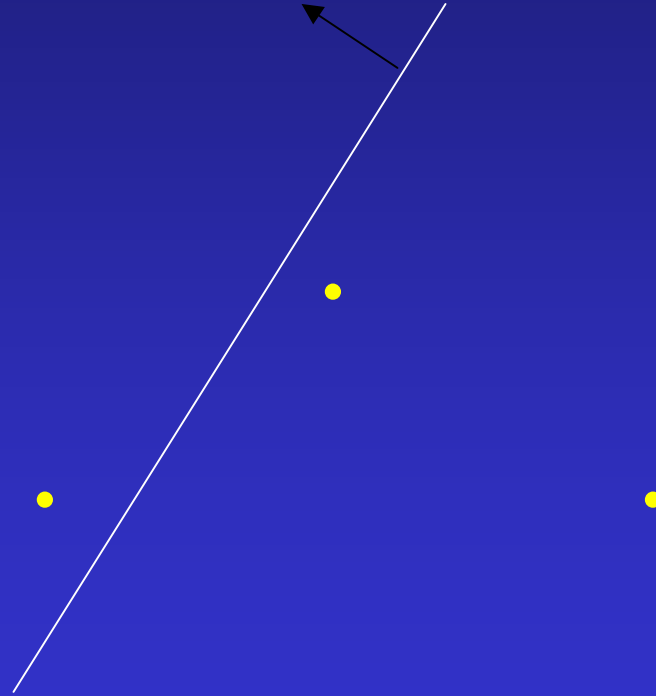
# VC Dimension of linear discriminants on a plane

011



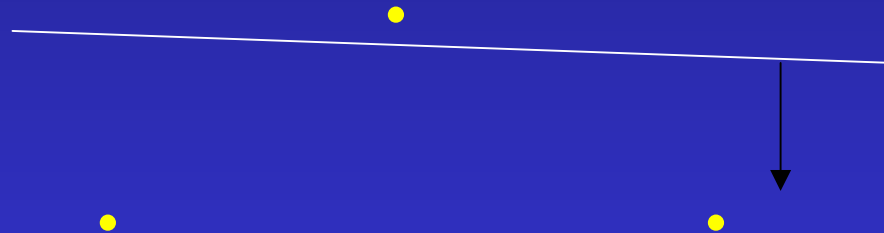
# VC Dimension of linear discriminants on a plane

100



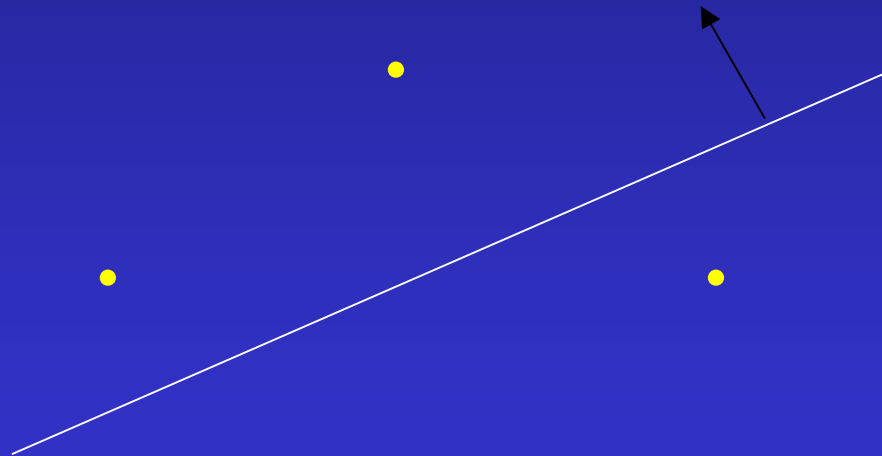
# VC Dimension of linear discriminants on a plane

101



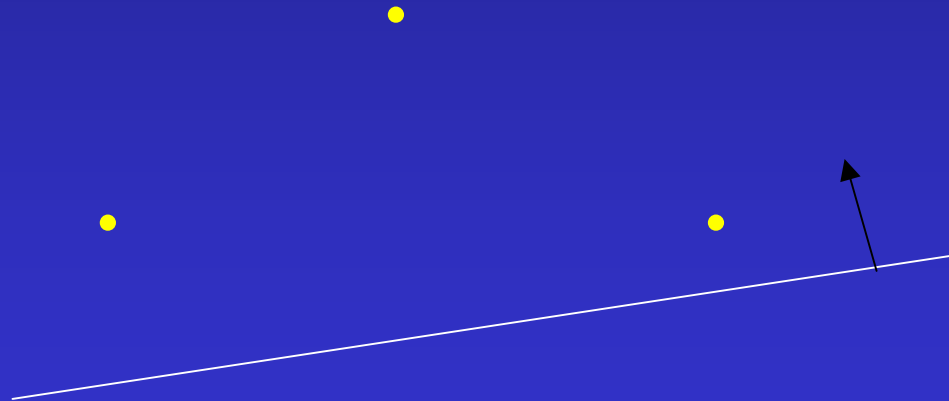
# VC Dimension of linear discriminants on a plane

110

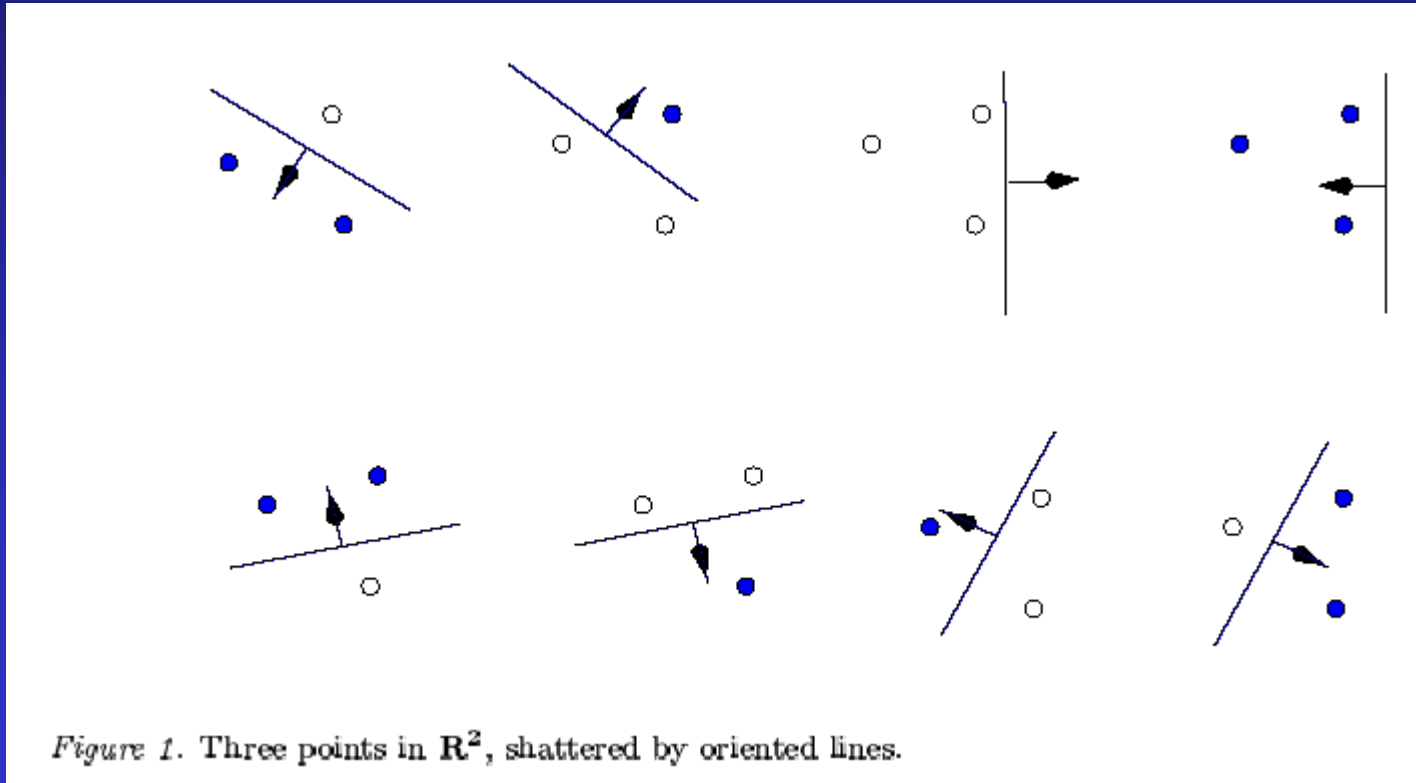


# VC Dimension of linear discriminants on a plane

111



# VC Dimension of points in 2-d space and perceptron

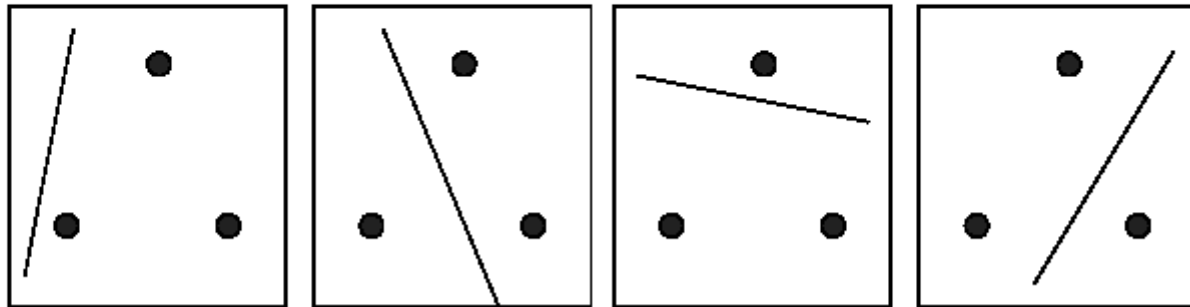


## VC Dimension of single perceptron with 2 input units (or points are in 2-d space)

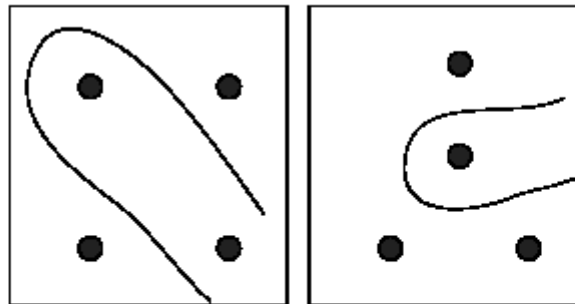
- $VC(H)$  is at least 3 since 3 non-collinear points can be shattered
- It is not 4 since we cannot find a set of four points that can be shattered
- Therefore  $VC(H) = 3$
- More generally, VC dimension of linear decision surfaces in  $r$ -dimensional space is  $r+1$

## VC-Dimension of Hyperplanes in $\mathbb{R}^2$

- Three points in  $\mathbb{R}^2$  can be shattered with hyperplanes.



- Four points cannot be shattered.



$\Rightarrow$  Hyperplanes in  $\mathbb{R}^2 \rightarrow VCdim=3$

**General:** Hyperplanes in  $\mathbb{R}^N \rightarrow VCdim=N+1$



## Capacity of a hyperplane

Fraction of dichotomies of  $n$  points in  $d$  dimensions that are linearly separable

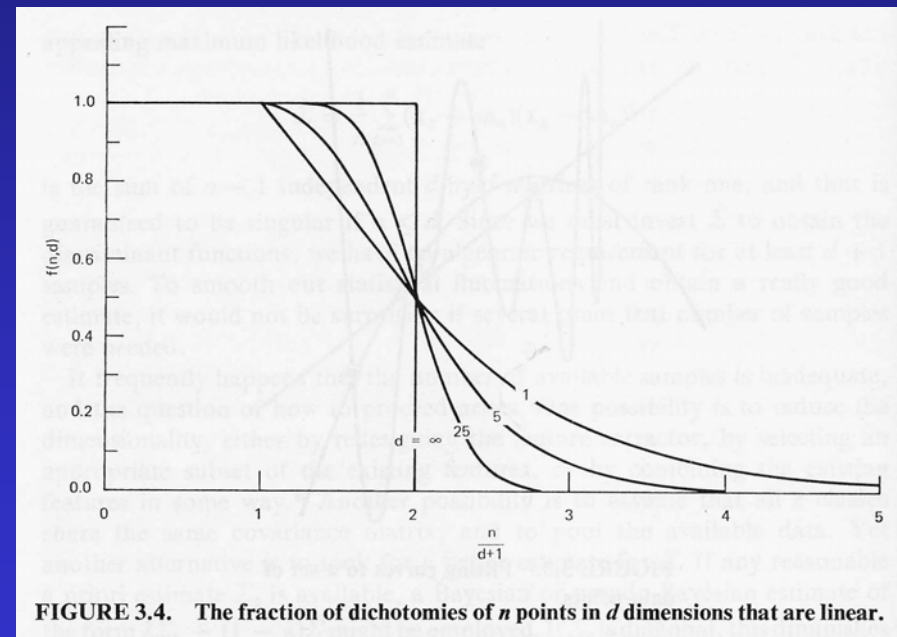
$$f(n, d) = \begin{cases} 1 & \dots \dots \dots n \leq d+1 \\ \frac{1}{2^n \sum_{i=0}^d \binom{n-1}{i}} & \dots \dots \dots n > d+1 \end{cases}$$

At  $n=2(d+1)$ , called the capacity of the hyperplane nearly one half of the dichotomies

are still linearly separable

Hyperplane is not overdetermined

Until number of samples is several times the dimensionality



## Example 3: VC Dimension of three Boolean literals and each hypothesis in $H$ is conjunction

- instance 1 = 100
- instance 2 = 010
- instance 3 = 001
  
- Exclude instance  $i$ :  $\sim l_i$
- Example: include instance 2 but exclude instances 1 and 3: use hypothesis  $\sim l_1 \wedge \sim l_3$
- *VC dimension is at least 3*
- VC dimension of conjunction of  $n$  Boolean literals is exactly  $n$  (proof is more difficult)

## Sample Complexity and the VC dimension

- How many randomly drawn samples are *sufficient* to PAC-learn any target concept  $C$  ?
- Using  $VC(H)$  as the measure of complexity of  $H$

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

- Analogous to the finite hypothesis case:

$$m \geq \frac{1}{\epsilon} (\ln |H| + \ln(1/\delta))$$

## Lower bound on sample complexity

If  $C$  is a concept class such that  $VC(C) > 2$  then observing fewer than

$$\max \left[ \frac{1}{\varepsilon} \log(1/\delta), \frac{VC(C) - 1}{32\varepsilon} \right]$$

samples, then  $L$  outputs a hypothesis  $h$  having  $error_D(h) > \varepsilon$

1. Provides number of samples *necessary* to PAC learn
2. Defined in terms of concept class  $C$  rather than  $H$

## VC Dimension for Neural Networks

- G-composition of C is the class of all functions that can be implemented on the network G, i.e., it is the hypothesis space representable by network G
- For acyclic layered networks containing s perceptrons, each with r inputs

$$\mathbf{VC}(\mathbf{C}_G^{\text{perceptrons}}) \leq 2(\mathbf{r} + 1)\mathbf{s} \log(\mathbf{e}\mathbf{s})$$

## VC Dimension for Neural Networks

- Bound on number of samples sufficient to learn with probability at least  $(1-d)$  any target concept from  $C_g$  to within error  $\epsilon$

$$\mathbf{m} \geq \frac{1}{\epsilon} (4 \log_2(2/\delta) + 16(\mathbf{r} + 1)\mathbf{s} \log(\mathbf{e}\mathbf{s}) \log(13/\epsilon))$$

- Does not apply to Backprop since the units are not sigmoid
- VC dimension of sigmoid units will at least as great as that of perceptrons

## VC Dimension of SVMs

### 6.2. The VC Dimension for Radial Basis Function Kernels

**THEOREM 5** Consider the class of Mercer kernels for which  $K(\mathbf{x}_1, \mathbf{x}_2) \rightarrow 0$  as  $\|\mathbf{x}_1 - \mathbf{x}_2\| \rightarrow \infty$ , and for which  $K(\mathbf{x}, \mathbf{x})$  is  $O(1)$ , and assume that the data can be chosen arbitrarily from  $\mathbf{R}^d$ . Then the family of classifiers consisting of support vector machines using these kernels, and for which the error penalty is allowed to take all values, has infinite VC dimension.

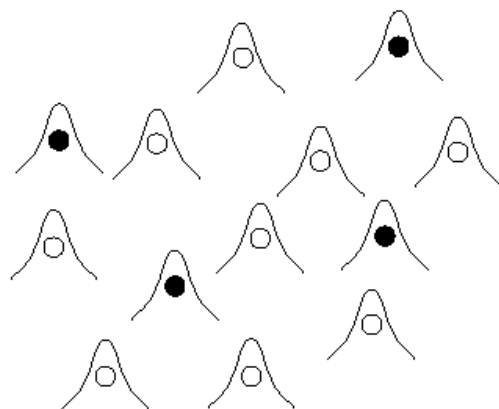


Figure 11. Gaussian RBF SVMs of sufficiently small width can classify an arbitrarily large number of training points correctly, and thus have infinite VC dimension.

## Rate of Convergence

**Question:** After  $n$  training examples, how close is the training error to the true error?

With probability  $\eta$  it hold for all  $h \in H$ :

$$Err_P(h) - Err_S(h) > \Phi(d, n, \eta)$$

$$\Phi(d, n) = \frac{1}{2} \sqrt{4 \frac{d \left( \ln \frac{2n}{d} + 1 \right) - \ln \frac{\eta}{4}}{n}}$$

- $n$  number of training examples
- $d$  VC-dimension of hypothesis space  $H$

$\implies$

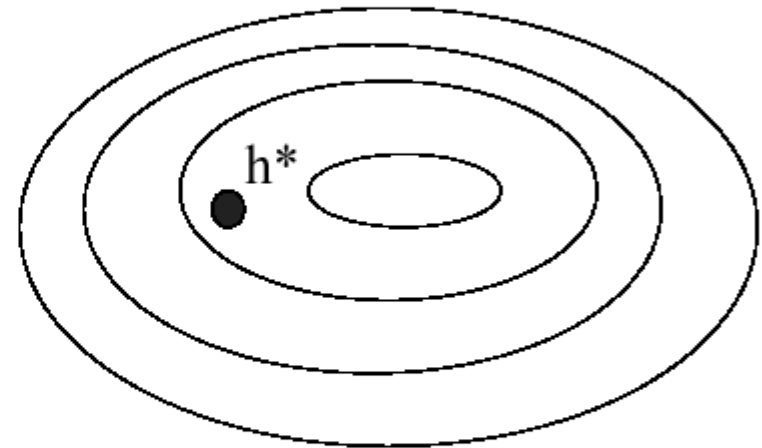
$$Err_P(h) \leq Err_S(h) + \Phi(d, n, \eta)$$



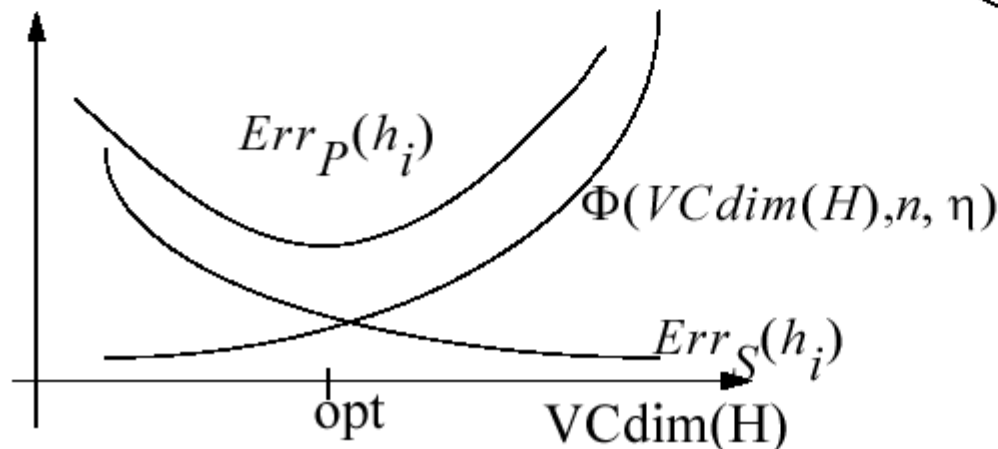
## SVM Motivation: Structural Risk Minimization

$$Err_P(h_i) \leq Err_S(h_i) + \Phi(VCdim(H), n, \eta)$$

**Idea:** Structure on hypothesis space.



**Goal:** Minimize upper bound on true error rate.



---

## MISTAKE BOUND MODEL OF LEARNING

- How many mistakes will learner make in its predictions before it learns the target concept
- Significant in practical settings where learning must be done while the system is in actual use, rather than in an off-line training stage
- Example, system to learn to approve credit card purchases based on data collected during use

## Mistake bound for Find-S algorithm

- Find-S
- Initialize  $h$  to the most specific hypothesis

$$\mathbf{l}_1 \wedge \neg \mathbf{l}_1 \wedge \mathbf{l}_2 \wedge \neg \mathbf{l}_2 \cdots \wedge \mathbf{l}_n \wedge \neg \mathbf{l}_n$$

- For each positive training instance  $x$ 
  - Remove from  $h$  any literal that is not satisfied by  $x$
- Output hypothesis  $h$
- Total no. of mistakes can be at most  $n+1$

## Mistake bound for Halving Algorithm

- Candidate Elimination Algorithm maintains a description of the version space, incrementally refining the version space as each new sample is encountered
- Assume majority vote is used among all hypotheses in the current version space
- Total no. of mistakes can be at most  $\log_2|H|$