

Discussion 7

Machine Learning Theory

VC dimension

田鹏超

tianpch@shanghaitech.edu.cn

Machine Learning Theory

1. Binary classification
2. Noise-free labeling

$$err_D(h) \leq err_S(h) + \varepsilon$$

?

Finite H

Infinite H

Realizable
 $c^* \in H$

Agnostic
 $c^* \notin H$

Realizable
 $c^* \in H$

Agnostic
 $c^* \notin H$

$$err_D(h) \leq err_S(h) + \sqrt{\frac{1}{2m} \ln \frac{2|H|}{\delta}}$$

$$err_D(h) \leq err_S(h) + \mathcal{O}\left(\sqrt{\frac{1}{m} \left(d + \ln \frac{1}{\delta}\right)}\right)$$

$$err_D(h) \leq \frac{1}{m} \left(\ln |H| + \ln \frac{1}{\delta} \right)$$

$$err_D(h) \leq \mathcal{O}\left(\frac{1}{m} \left(d \ln \frac{m}{d} + \ln \frac{1}{\delta} \right)\right)$$

Some concepts:

1. Consistent learner
2. Version space (VS)
3. ε -exhausted VS
4. PAC-learnable
5. Hoeffding inequality
6. Dichotomy
7. Shattering
8. Shattering coefficient (growth function)
9. VC-dimension
10. Sauer's lemma

Effective number of hypotheses

- $H[S]$ - the set of splittings of dataset S using concepts from H .
- $H[m]$ - max number of ways to split m points using concepts in H

$$H[m] = \max_{\substack{|S|=m \\ \forall S \subseteq X}} |H[S]| \quad H[m] \leq 2^m$$

Definition: H shatters S if $|H[S]| = 2^{|S|} = 2^m$
specific

Shattering, VC-dimension

Definition: H shatters S if $|H[S]| = 2^{|S|} = 2^m$

A set of points S is shattered by H if there are hypotheses in H that split S in all of the $2^{|S|}$ possible ways, all possible ways of classifying points in S are achievable using concepts in H .

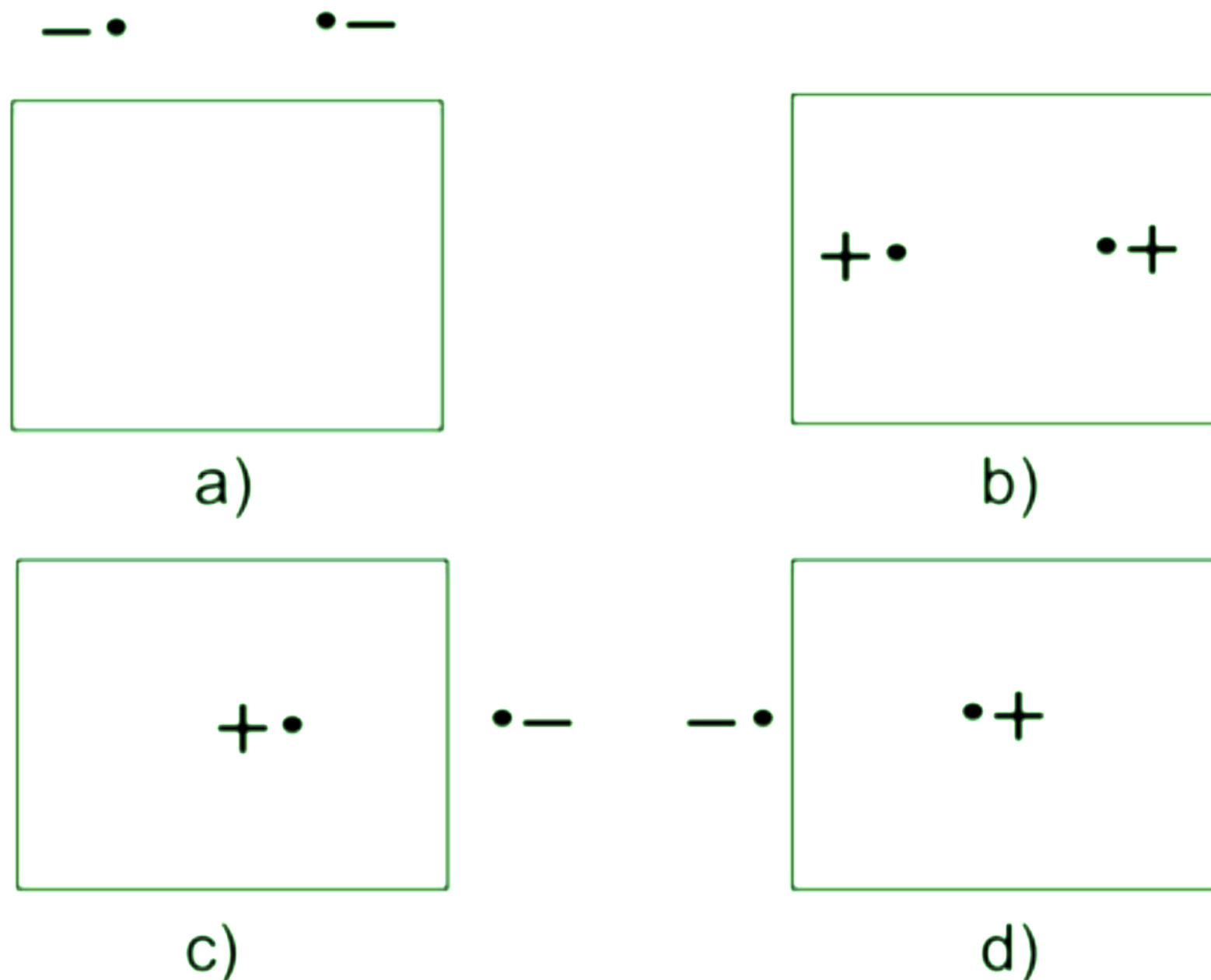
Definition: VC-dimension (Vapnik-Chervonenkis dimension)

The **VC-dimension** of a hypothesis space H is the cardinality of the largest set S that can be shattered by H . $|S| = m$

If arbitrarily large finite sets can be shattered by H , then VCdim(H) = ∞

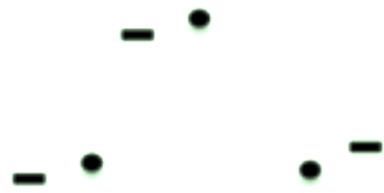
VC dimension: rectangles

Now we look at another example, where the hypothesis labels the point inside the rectangle decided by the two points (x_1, y_1) , (x_2, y_2) positive, and otherwise negative. The case for two points:

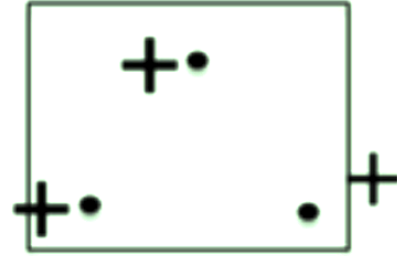


VC dimension: rectangles

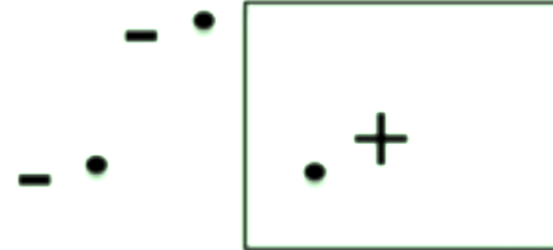
For three points:



a)



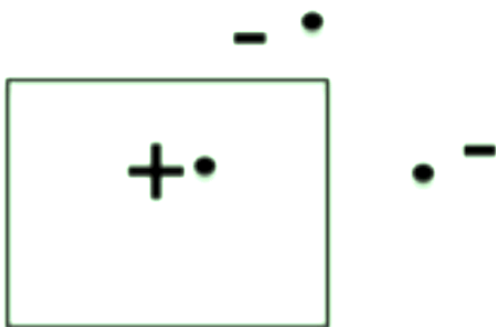
b)



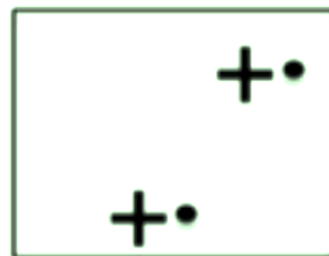
c)



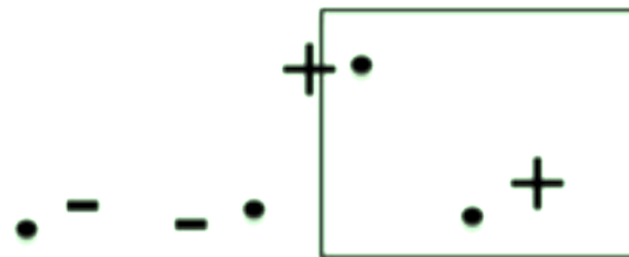
d)



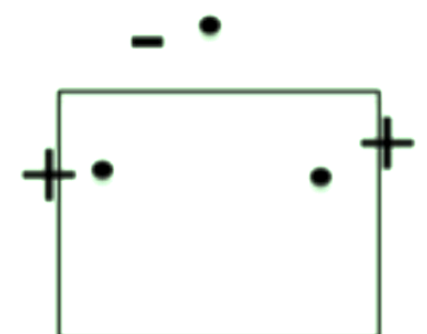
e)



f)



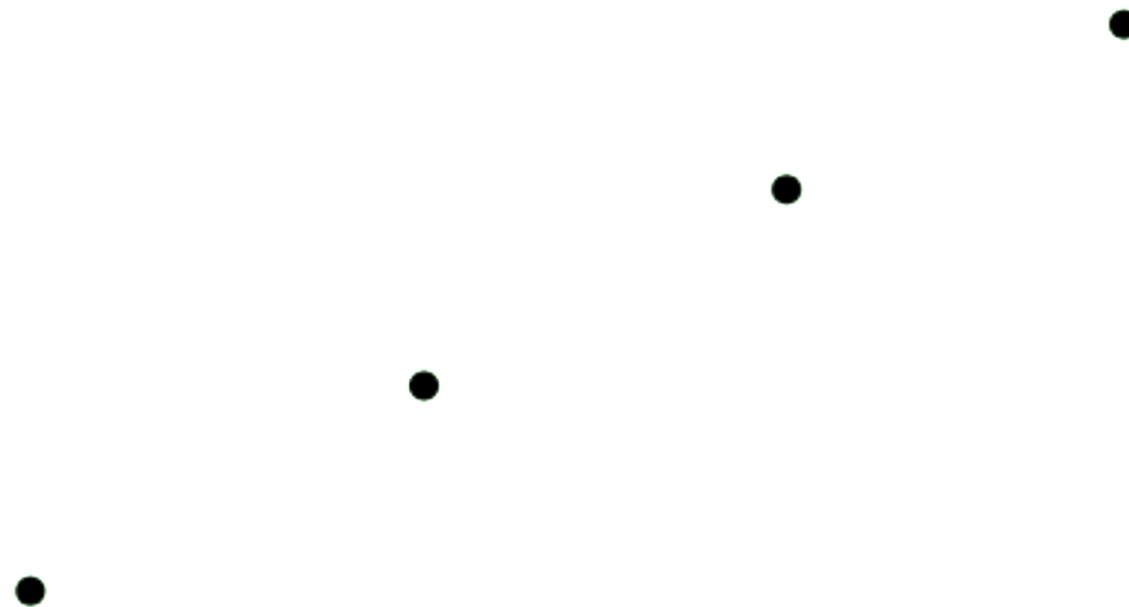
g)



h)

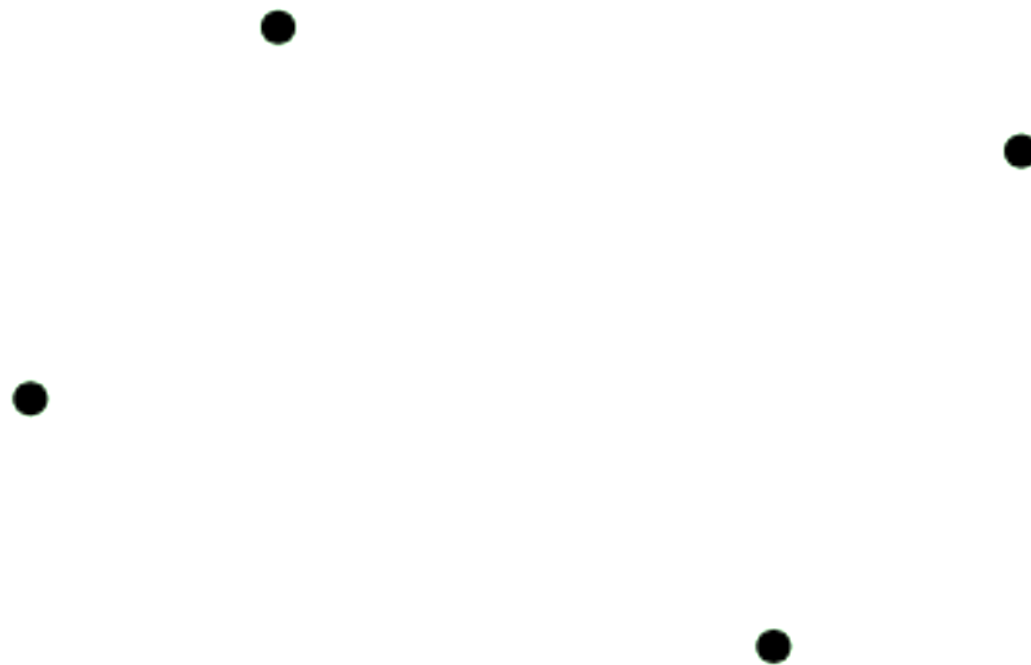
VC dimension: rectangles

The case for four points is a little different; it is not possible to produce all the dichotomies for certain situations, one of them is presented below:



VC dimension: rectangles

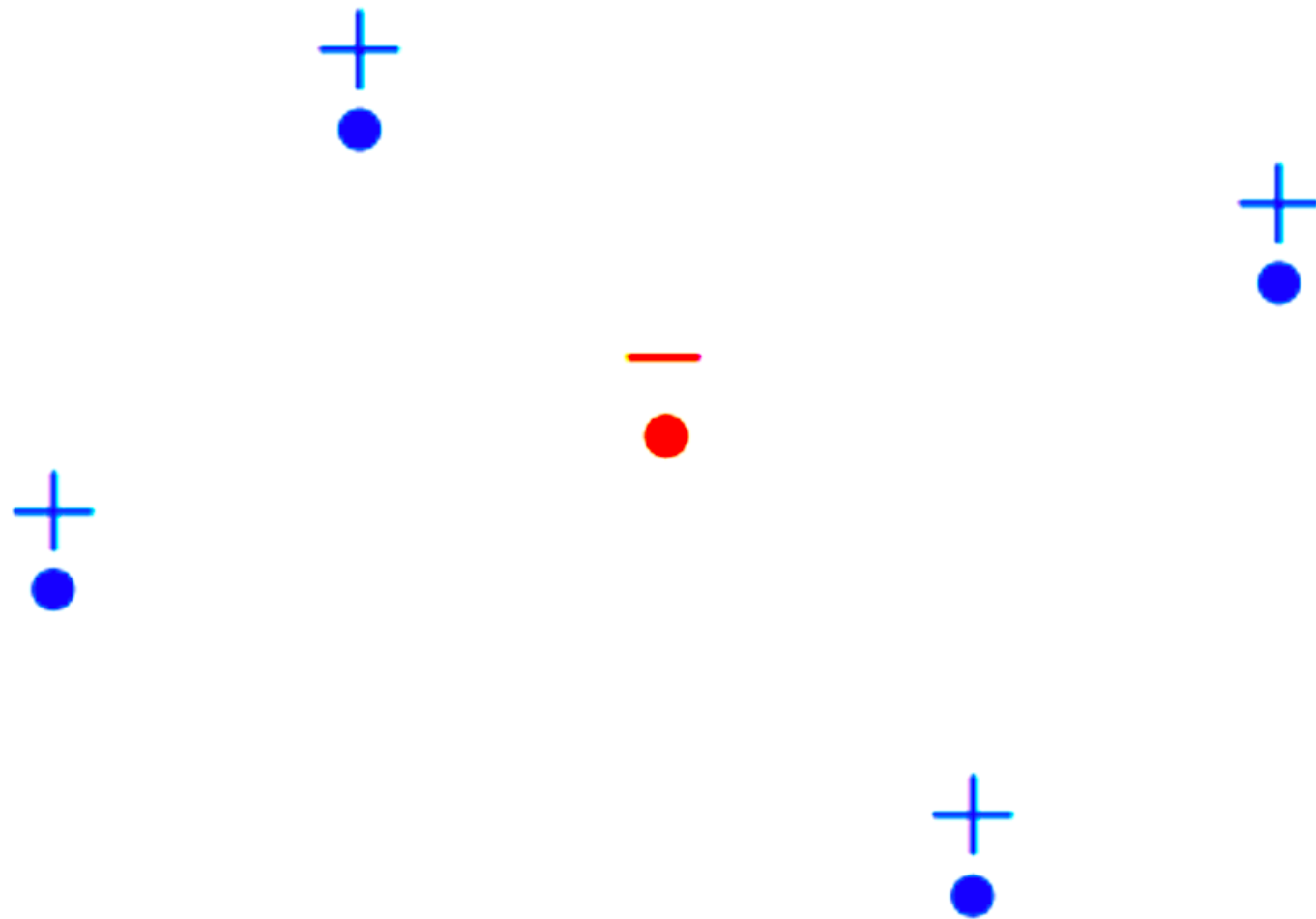
However, this configuration can be shattered!



Therefore, the VC dimension is at least 4

VC dimension: rectangles

But not this one:



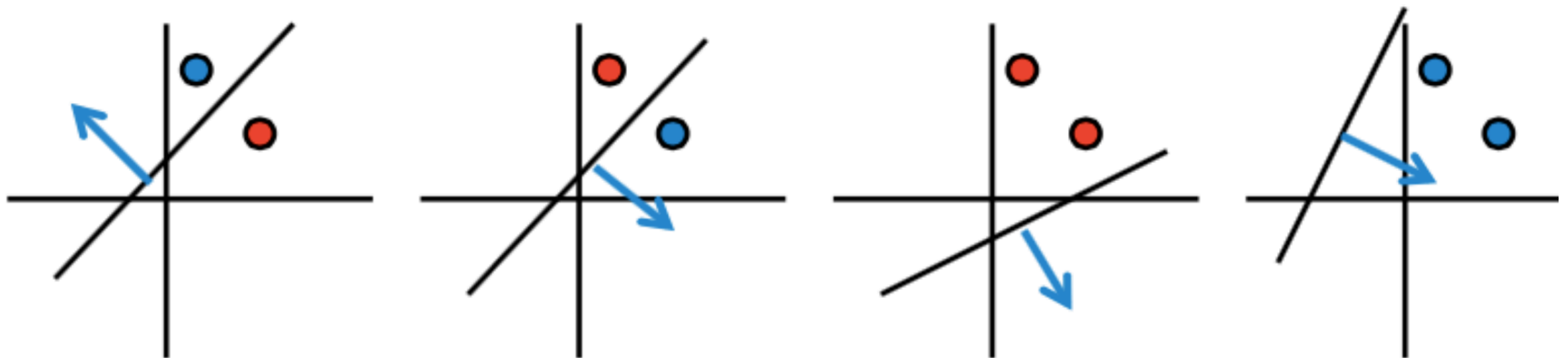
VC dimension: rectangles

The VC dimension of rectangles is the cardinality of the **maximum** set of points that can be shattered by a rectangle

The VC dimension of rectangles is 4 because there **exists** a set of 4 points that can be shattered by a rectangle and **any** set of 5 points cannot be shattered by a rectangle

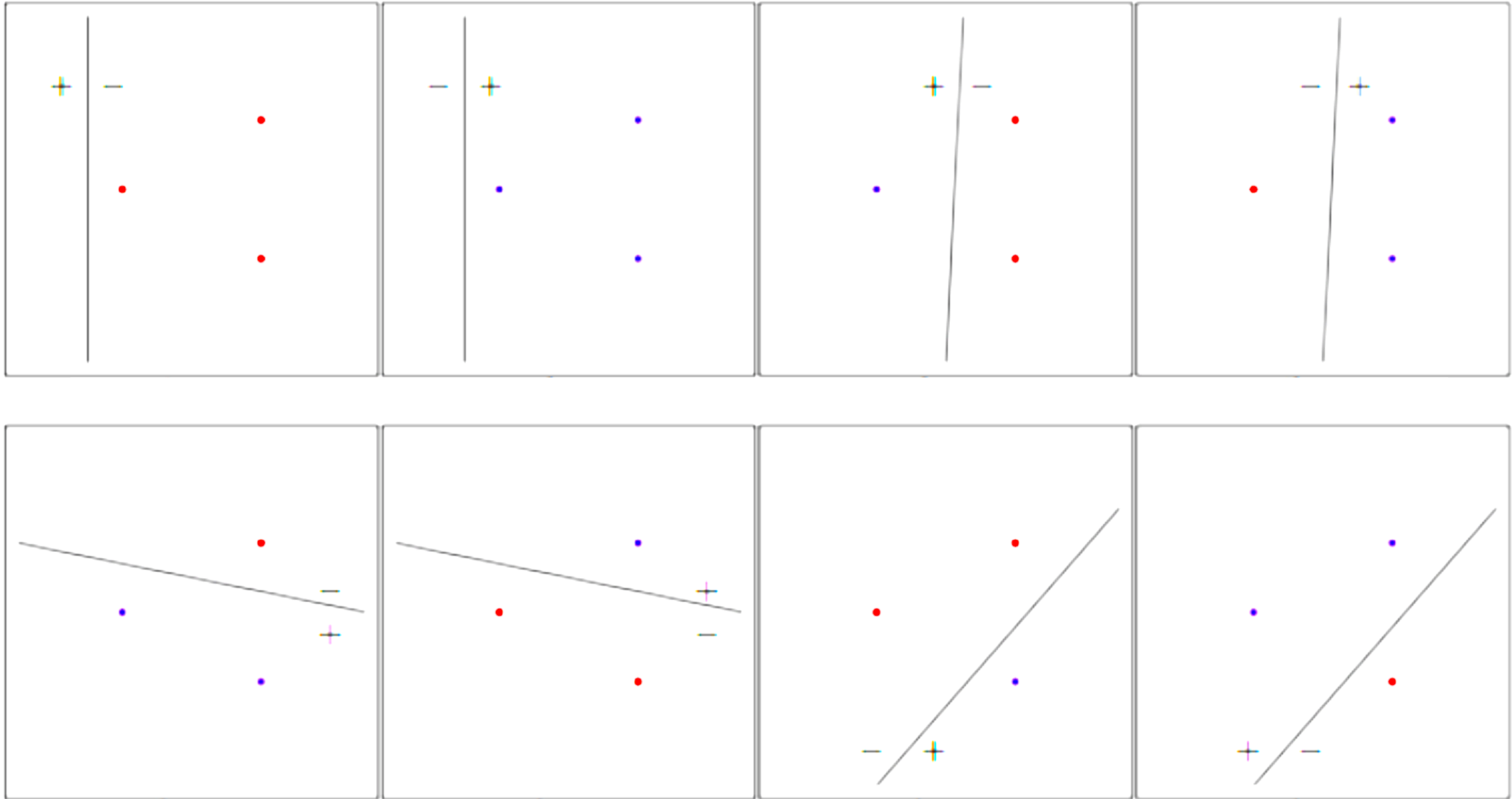
VC dimension: linear separator

Can $h_{\theta}(\mathbf{x}) = \text{sign}(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$ shatter these points?



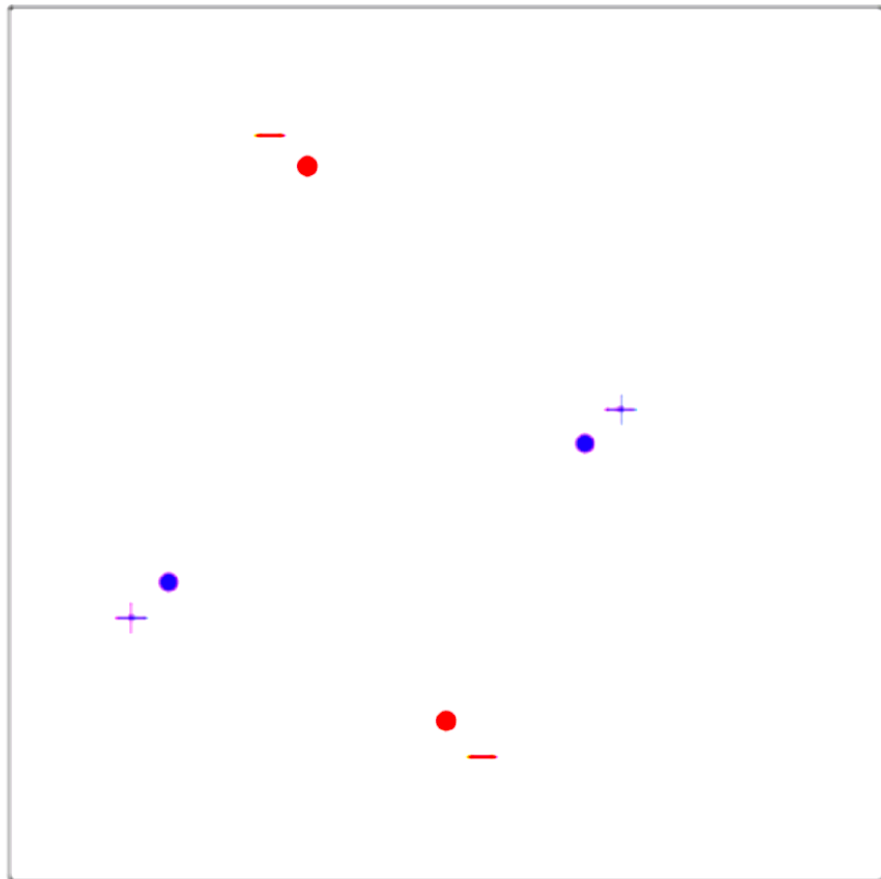
- $d_{vc} = 2?$, $d_{vc} \leq 2?$, $d_{vc} \geq 2?$

VC dimension: linear separator



VC dimension

However, things are a little different with the case of 4 points. For the case of 4 points, there are $2^4 - 2 = 14$ kinds of labeling. As the usual 2^m number of labelings, this time there are two labeling that is not achievable by linear classifiers. Below presents one of them:



- In \mathbb{R}^2 , linear separator has $d_{vc} = 3$

VC dimension

In general, linear classifier (perceptron) in d dimensions with a constant term

$$d_{\text{VC}} = d + 1$$

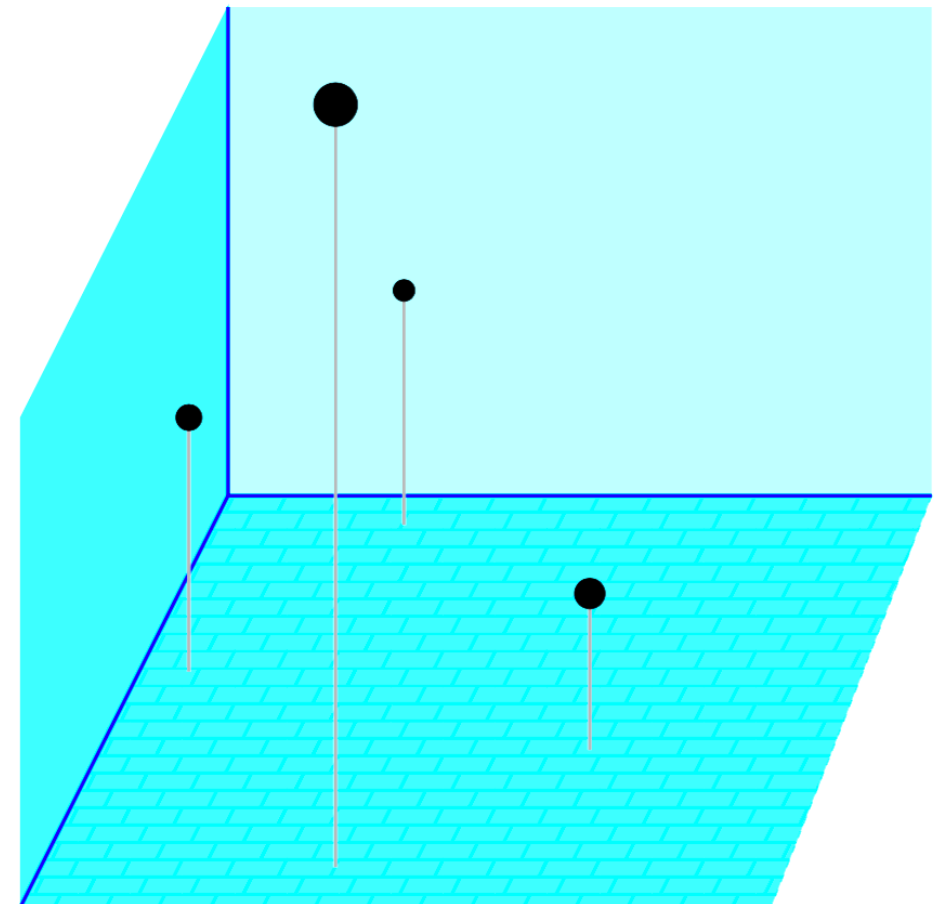
For $d = 2$, $d_{\text{VC}} = 3$

In general, $d_{\text{VC}} = d + 1$

We will prove two directions:

$$d_{\text{VC}} \leq d + 1$$

$$d_{\text{VC}} \geq d + 1$$



数据是在 d 维空间里的，但是分离平面的参数要加上常数项，一共是 $d+1$ 个参数。 $\text{sign}(w_0 \cdot 1 + w_1x_1 + \dots + w_dx_d)$

Here is one direction

A set of $N = d + 1$ points in \mathbb{R}^d shattered by the perceptron:

$$\mathbf{x} = [1, x_1, x_2, \dots, x_d]^T$$

$$X = \begin{bmatrix} \text{---} \mathbf{x}_1^\top \text{---} \\ \text{---} \mathbf{x}_2^\top \text{---} \\ \text{---} \mathbf{x}_3^\top \text{---} \\ \vdots \\ \text{---} \mathbf{x}_{d+1}^\top \text{---} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & & 0 \\ & \vdots & & \ddots & 0 \\ 1 & 0 & \dots & 0 & 1 \end{bmatrix} \left. \vphantom{\begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & & 0 \\ & \vdots & & \ddots & 0 \\ 1 & 0 & \dots & 0 & 1 \end{bmatrix}} \right\} d+1$$

$\underbrace{\hspace{10em}}_{d+1}$

X is invertible

Can we shatter this data set?

For any $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{d+1} \end{bmatrix} = \begin{bmatrix} \pm 1 \\ \pm 1 \\ \vdots \\ \pm 1 \end{bmatrix}$, can we find a vector \mathbf{w} satisfying

$$\text{sign}(\mathbf{X}\mathbf{w}) = \mathbf{y}$$

Easy! Just make $\mathbf{X}\mathbf{w} = \mathbf{y}$

which means $\mathbf{w} = \mathbf{X}^{-1}\mathbf{y}$

我们对于一个特定的包含 $d + 1$ 个数据点的数据集，可以产生所有的 2^d 个 dichotomies。这意味着我们可以“粉碎”某个 $d + 1$ 样本容量的数据集。所以“断点”肯定不是 $d + 1$ 。

We can shatter these $d + 1$ points

This implies what?

[a] $d_{\text{VC}} = d + 1$

[b] $d_{\text{VC}} \geq d + 1$ ✓

[c] $d_{\text{VC}} \leq d + 1$

[d] No conclusion

Now, to show that $d_{vc} \leq d + 1$

We need to show that:

- [a] There are $d + 1$ points we cannot shatter
- [b] There are $d + 2$ points we cannot shatter
- [c] We cannot shatter *any* set of $d + 1$ points
- [d] We cannot shatter *any* set of $d + 2$ points ✓



Prove for "ANY"!!!

Here is the other direction

Take **any** $d + 2$ points in \mathbb{R}^d !!

For any $d + 2$ points in \mathbb{R}^d : $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{d+1}, \mathbf{x}_{d+2}$

More points than dimensions \implies we must have

$$\mathbf{x}_j = \sum_{i \neq j} a_i \mathbf{x}_i$$

where not all the a_i s are zeros

Our purpose is then to design a dichotomy that any linear separator cannot generate on $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{d+1}, \mathbf{x}_{d+2}$!!

$$\mathbf{x}_j = \sum_{i \neq j} a_i \mathbf{x}_i$$

- Consider the following dichotomy:

$$\begin{aligned} y_i &= \text{sign}(a_i) && \text{for } \mathbf{x}_i \text{'s with non-zero } a_i \\ y_j &= -1 && \text{for } \mathbf{x}_j \end{aligned}$$

- No perceptron can implement such dichotomy!

- The dichotomy we construct ($j = 1$)

$$\begin{array}{ccccccc} \mathbf{x}_i & \mathbf{x}_1 & = & a_2 \mathbf{x}_2 & + \dots & + 0 \cdot \mathbf{x}_{d+1} & + a_{d+2} \mathbf{x}_{d+2} \\ & \downarrow & & \downarrow & & \downarrow & \downarrow \\ y_i & -1 & & \text{sign}(a_2) & & \text{whatever} & \text{sign}(a_{d+2}) \end{array}$$

- Show for any $\mathbf{w} \in \mathbb{R}^{d+1}$, this dichotomy cannot appear!

- The dichotomy we construct ($j = 1$)

$$\begin{array}{ccccccc}
 \mathbf{x}_i & \mathbf{x}_1 & = & 1 \cdot \mathbf{x}_2 & + & \dots & + 0 \cdot \mathbf{x}_{d+1} - 1 \cdot \mathbf{x}_{d+2} \\
 & \downarrow & & \downarrow & & & \downarrow & \downarrow \\
 y_i & -1 & & +1 & & & +1 & -1
 \end{array}$$

- Show for any $\mathbf{w} \in \mathbb{R}^{d+1}$, this dichotomy cannot appear!
- Notice that $y_i = \text{sign}(\mathbf{w}^T \mathbf{x}_i)$

$$\mathbf{x}_j = \sum_{i \neq j} a_i \mathbf{x}_i \implies \mathbf{w}^T \mathbf{x}_j = \sum_{i \neq j} a_i \boxed{\mathbf{w}^T \mathbf{x}_i} \longrightarrow$$

这项的sign就是 y_i ，假设我们可以选 \mathbf{w} 使得这项的sign和 y_i 匹配。则由于我们在设定dichotomy的时候，把 y_i 选的和 a_i 的sign选的一样！！到处矛盾！

- Since $\text{sign}(\mathbf{w}^T \mathbf{x}_i) = y_i = \text{sign}(a_i)$, then $a_i \mathbf{w}^T \mathbf{x}_i > 0$
- This forces $\mathbf{w}^T \mathbf{x}_j = \sum_{i \neq j} a_i \mathbf{w}^T \mathbf{x}_i > 0$
- Therefore, $y_i = \text{sign}(\mathbf{w}^T \mathbf{x}_j) = +1$, contradiction!!!

Putting it together...

We proved $d_{vc} \leq d + 1$ and $d_{vc} \geq d + 1$

$$d_{vc} = d + 1$$