

机器学习基石第5、6周笔记

机器学习基石第5、6周笔记

5.1 Recap and Preview

5.2 Effective Number of Line

5.3 Effective Number of Hypotheses

5.4 Break Point

6.1

6.2 Bounding Function: Basic Cases

6.3 Bounding Function: Inductive Cases

6.4 A Pictorial Proof

5.1 Recap and Preview

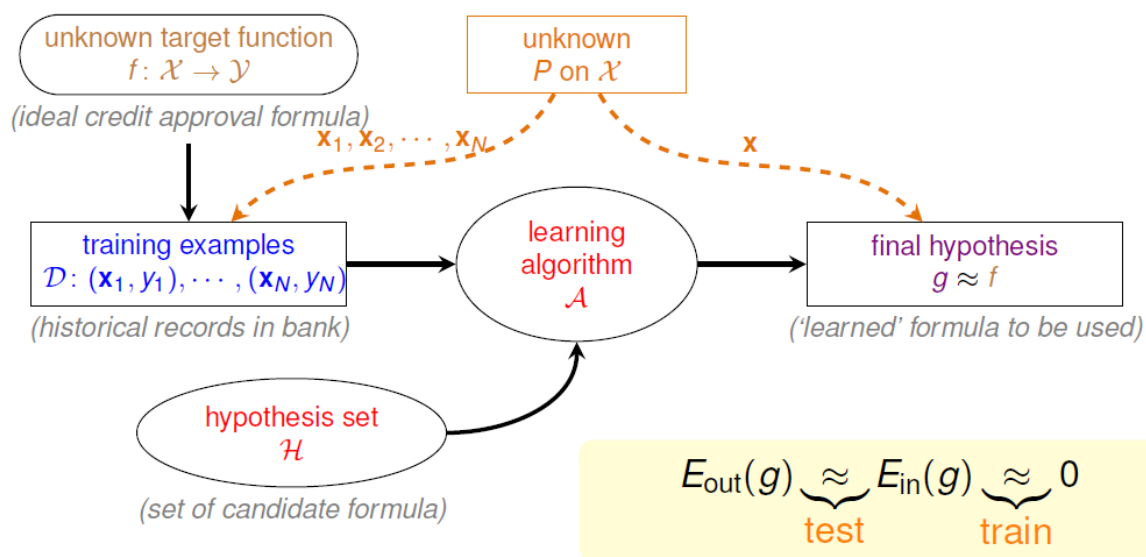
Recap: the 'Statistical' Learning Flow

if $|\mathcal{H}| = M$ finite, N large enough,

for whatever g picked by \mathcal{A} , $E_{\text{out}}(g) \approx E_{\text{in}}(g)$

if \mathcal{A} finds one g with $E_{\text{in}}(g) \approx 0$,

PAC guarantee for $E_{\text{out}}(g) \approx 0 \implies$ **learning possible :-)**



机器学习能够学习的前提：

训练数据 \mathcal{D} 与测试hypothesis的数据都来自同样的一个样本分布

如果**假设集合(hypothesis set)**不太大，是有限的，而且**数据量(N)**够大的话，根据霍夫丁不等式，可以确保对每一个hypothesis而言， $E_{\text{in}} \approx E_{\text{out}}$ ，所以干脆选一个 E_{in} 最小的，如果这个 $E_{\text{in}} \approx 0$ ，那么 E_{out} 也接近于0，这就达到了学习的效果。

Trade-off on M

- 1 can we make sure that $E_{\text{out}}(g)$ is close enough to $E_{\text{in}}(g)$?
- 2 can we make $E_{\text{in}}(g)$ small enough?

small M

- 1 Yes!,
 $\mathbb{P}[\text{BAD}] \leq 2 \cdot M \cdot \exp(\dots)$
- 2 No!, too few choices

large M

- 1 No!,
 $\mathbb{P}[\text{BAD}] \leq 2 \cdot M \cdot \exp(\dots)$
- 2 Yes!, many choices

using the right M (or \mathcal{H}) is important
 $M = \infty$ **doomed?**

这里将学习拆为两个问题:

- E_{out} 与 E_{in} 是不是接近
- E_{in} 是不是足够小

在讨论这两个问题的时候, 涉及到一个参数 M , M 是 hypothesis set 的大小, 我们现在暂时说 hypothesis set 是有限大的, 那它到底到大, 我们用这个 M 来代表.

- 当 M 很小的时候, 根据公式可知 Bad data 发生的概率很小, 那么 E_{out} 与 E_{in} 就很接近, 但是算法可供选择的假设少, 不一定能够找到一个很小的 E_{in} , 使 $E_{\text{in}} \approx 0$
- 当 M 很大的时候, 算法的选择足够多, 可以找到一个好的 E_{in} , 使 $E_{\text{in}} \approx 0$, 但是对第一个问题而言, Bad data 发生的概率大大增加,

Known

$$\mathbb{P} [|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq 2 \cdot M \cdot \exp(-2\epsilon^2 N)$$

Todo

- establish a **finite quantity** that replaces M

$$\mathbb{P} [|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \stackrel{?}{\leq} 2 \cdot m_{\mathcal{H}} \cdot \exp(-2\epsilon^2 N)$$

- justify the feasibility of learning for infinite M
- study $m_{\mathcal{H}}$ to understand its trade-off for 'right' \mathcal{H} , just like M

用 $m_{\mathcal{H}}$ 替换 M , 用 m 代表它可能比原来的大 M 来得小, 然后用 \mathcal{H} 代表它跟 hypothesis set 的一些性质可能是有关系的

Data size: how large do we need?

One way to use the inequality

$$\mathbb{P} \left[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon \right] \leq \underbrace{2 \cdot M \cdot \exp(-2\epsilon^2 N)}_{\delta}$$

is to pick a tolerable difference ϵ as well as a tolerable **BAD** probability δ , and then gather data with size (N) large enough to achieve those tolerance criteria. Let $\epsilon = 0.1$, $\delta = 0.05$, and $M = 100$. What is the data size needed?

① 215

② 415

③ 615

④ 815

$$\delta = 2M * e^{-2\epsilon^2 N} \Rightarrow N = \frac{1}{2\epsilon^2} \ln \frac{2M}{\delta}$$

N=415

5.2 Effective Number of Line

$$\mathbb{P} \left[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon \right] \leq 2 \cdot M \cdot \exp(-2\epsilon^2 N)$$

- **BAD events** \mathcal{B}_m : $|E_{\text{in}}(h_m) - E_{\text{out}}(h_m)| > \epsilon$
- to give \mathcal{A} freedom of choice: bound $\mathbb{P}[\mathcal{B}_1 \text{ or } \mathcal{B}_2 \text{ or } \dots \mathcal{B}_M]$
- worst case: all \mathcal{B}_m non-overlapping

$$\mathbb{P}[\mathcal{B}_1 \text{ or } \mathcal{B}_2 \text{ or } \dots \mathcal{B}_M] \leq \mathbb{P}[\mathcal{B}_1] + \mathbb{P}[\mathcal{B}_2] + \dots + \mathbb{P}[\mathcal{B}_M]$$

union bound

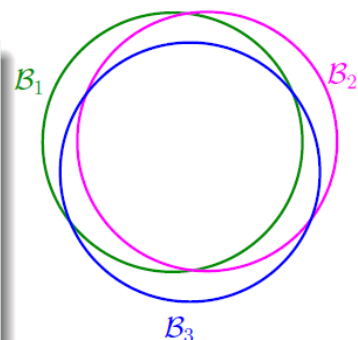
我们用联集的方式把m个不好的事情联集起来:计算每一个不好的事情发生的几率, 然后把它加加起来。但是当m无限大的时候, 每次这个概率不为0的话, 这个union bound 的值会大于1, 这样union bound的值就没有了意义。

$$\text{union bound } \mathbb{P}[\mathcal{B}_1] + \mathbb{P}[\mathcal{B}_2] + \dots + \mathbb{P}[\mathcal{B}_M]$$

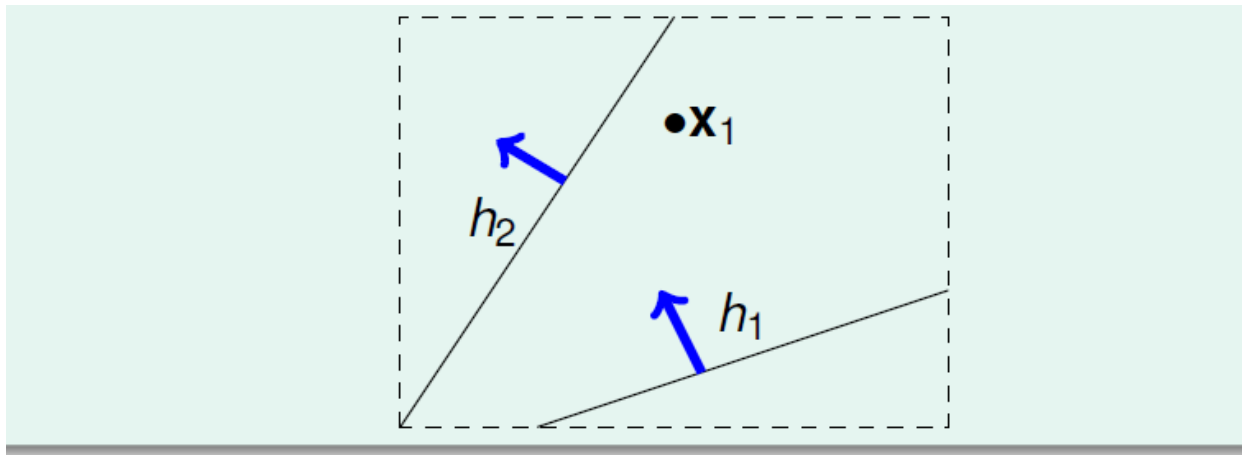
- **BAD events** \mathcal{B}_m : $|E_{\text{in}}(h_m) - E_{\text{out}}(h_m)| > \epsilon$

overlapping for similar hypotheses $h_1 \approx h_2$

- why? ① $E_{\text{out}}(h_1) \approx E_{\text{out}}(h_2)$
 ② for most \mathcal{D} , $E_{\text{in}}(h_1) = E_{\text{in}}(h_2)$
- union bound **over-estimating**

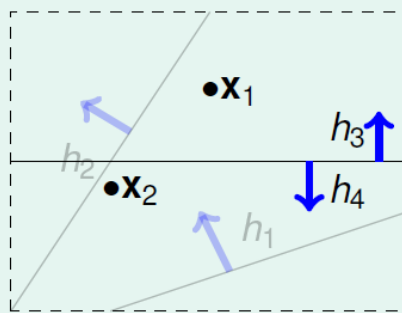


坏事情是重叠的，单纯的加和就会造成过分的估计，所以我们的目的就是将无数个hypothesis分成有限个不同的类，去除掉重叠部分

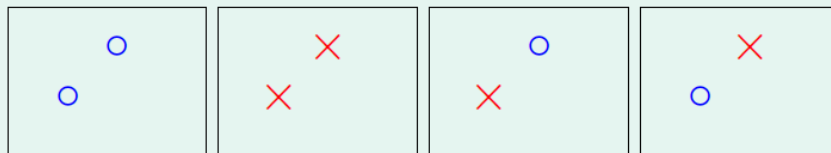


2 kinds: $h_1\text{-like}(\mathbf{x}_1) = \circ$ or $h_2\text{-like}(\mathbf{x}_1) = \times$

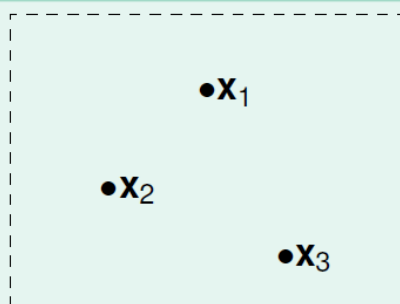
- how many **kinds of** lines if viewed from two inputs $\mathbf{x}_1, \mathbf{x}_2$?



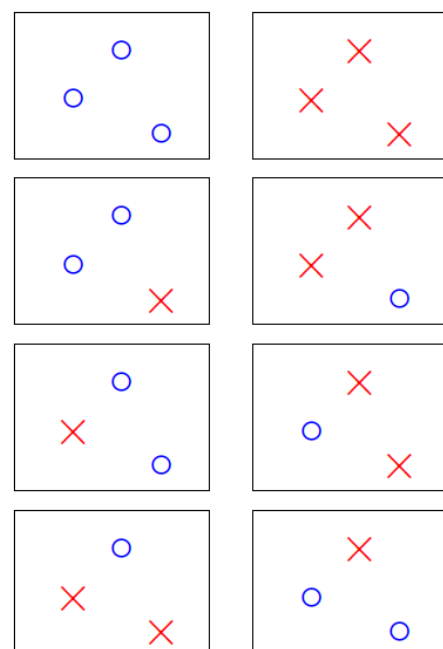
4:



for three inputs $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3$

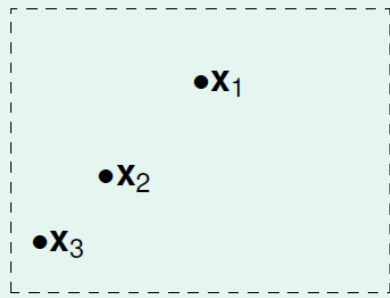


8:



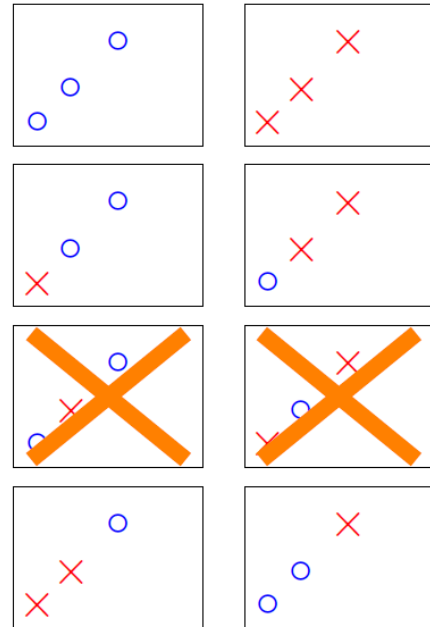
always **8 for three inputs?**

for **another** three inputs
 x_1, x_2, x_3

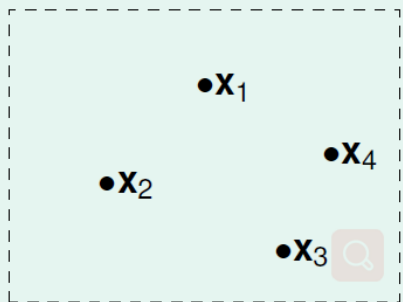


'fewer than 8' when degenerate
 (e.g. collinear or same inputs)

6:

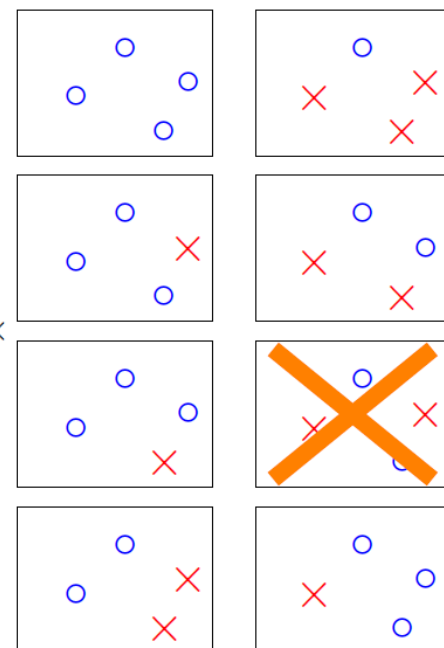


for four inputs x_1, x_2, x_3, x_4



for any four inputs
at most 14

14: 2x



Effective Number of Lines

maximum kinds of lines with respect to N inputs $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$
 \iff **effective number of lines**

- must be $\leq 2^N$ (why?)
- finite 'grouping' of infinitely-many lines $\in \mathcal{H}$
- wish:

$$\begin{aligned} & \mathbb{P} [|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \\ & \leq 2 \cdot \text{effective}(N) \cdot \exp(-2\epsilon^2 N) \end{aligned}$$

lines in 2D

N	effective(N)
1	2
2	4
3	8
4	$14 < 2^4$

if ① effective(N) can replace M and

② effective(N) $\ll 2^N$

learning possible with infinite lines :-)

Fun Time

What is the effective number of lines for five inputs $\in \mathbb{R}^2$?

① 14

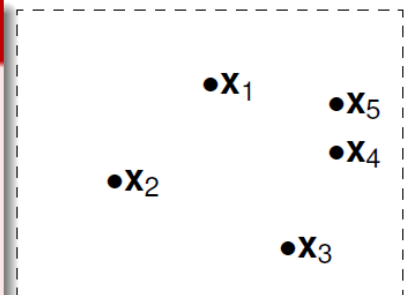
② 16

③ 22

④ 32

Reference Answer: ③

If you put the inputs roughly around a circle, you can then pick any consecutive inputs to be on one side of the line, and the other inputs to be on the other side. The procedure leads to effectively 22 kinds of lines, which is **much smaller than** $2^5 = 32$. You shall find it difficult to generate more kinds by varying the inputs, and we will give a formal proof in future lectures.



0个点: 1种

1个点: 5种

2个点: 5种

5.3 Effective Number of Hypotheses

Dichotomies: Mini-hypotheses

$$\mathcal{H} = \{\text{hypothesis } h: \mathcal{X} \rightarrow \{\times, \circ\}\}$$

- call

$$h(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) = (h(\mathbf{x}_1), h(\mathbf{x}_2), \dots, h(\mathbf{x}_N)) \in \{\times, \circ\}^N$$

a **dichotomy**: hypothesis 'limited' to the eyes of $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$

- $\mathcal{H}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$:

all dichotomies 'implemented' by \mathcal{H} on $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$

	hypotheses \mathcal{H}	dichotomies $\mathcal{H}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$
e.g.	all lines in \mathbb{R}^2	$\{\circ\circ\circ\circ, \circ\circ\circ\times, \circ\circ\times\times, \dots\}$
size	possibly infinite	upper bounded by 2^N

$|\mathcal{H}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)|$: candidate for **replacing M**

dichotomy: 二分意思是我们手上的点 分成两堆，圈圈的一堆，叉叉的一堆

dichotomy只对N个特定的点来做取值，代表着N个点有几种不一样的组合，所以它的上限是 2^N

hypothesis是对我们的输入空间 \mathcal{X} 里面所有的点取值，所以用线来区分的话，可能有无限多条线

Growth Function

- $|\mathcal{H}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)|$: depend on inputs $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$
- growth function:
remove dependence by **taking max of all possible $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$**

$$m_{\mathcal{H}}(N) = \max_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \in \mathcal{X}} |\mathcal{H}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)|$$

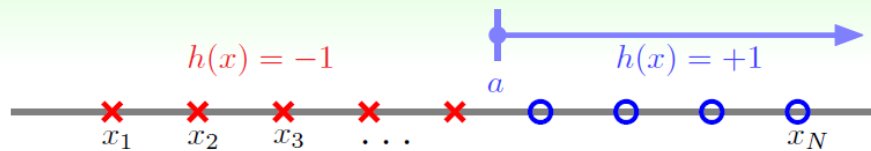
- finite, upper-bounded by 2^N

lines in 2D

N	$m_{\mathcal{H}}(N)$
1	2
2	4
3	$\max(\dots, 6, 8) = 8$
4	$14 < 2^4$

成长函数(Growth Function): 记为 $m_{\mathcal{H}}(N)$,对于由N个点组成的不同集合中, 某集合对应的 dichotomy最大, 那么这个dichotomy值就是 $m_{\mathcal{H}}(N)$

Growth Function for Positive Rays



- $\mathcal{X} = \mathbb{R}$ (one dimensional)
- \mathcal{H} contains h , where **each** $h(x) = \text{sign}(x - a)$ **for threshold** a
- 'positive half' of 1D perceptrons

one dichotomy for $a \in$ each spot (x_n, x_{n+1}) :

$$m_{\mathcal{H}}(N) = N + 1$$

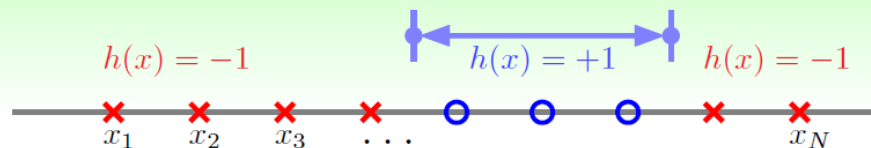
$$(N + 1) \ll 2^N \text{ when } N \text{ large!}$$

x_1	x_2	x_3	x_4
○	○	○	○
×	○	○	○
×	×	○	○
×	×	×	○
×	×	×	×

Training versus Testing

Effective Number of Hypotheses

Growth Function for Positive Intervals



- $\mathcal{X} = \mathbb{R}$ (one dimensional)
- \mathcal{H} contains h , where **each** $h(x) = +1$ **iff** $x \in [\ell, r]$, **-1 otherwise**

one dichotomy for each 'interval kind'

$$m_{\mathcal{H}}(N) = \underbrace{\binom{N+1}{2}}_{\text{interval ends in } N+1 \text{ spots}} + \underbrace{1}_{\text{all } \times}$$

$$= \frac{1}{2}N^2 + \frac{1}{2}N + 1$$

$$\left(\frac{1}{2}N^2 + \frac{1}{2}N + 1\right) \ll 2^N \text{ when } N \text{ large!}$$

x_1	x_2	x_3	x_4
○	×	×	×
○	○	×	×
○	○	○	×
○	○	○	○
×	○	×	×
×	○	○	×
×	○	○	○
×	×	○	×
×	×	×	○
×	×	×	×

$$C_{N+1}^2 + 1 = \frac{N^2}{2} + \frac{N}{2} + 1$$

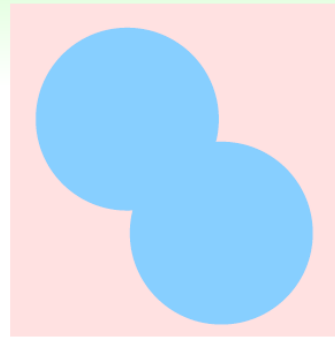
C_{N+1}^2 : 从 $N+1$ 个空中取两个空, 放置两个端点

1: 两个端点放置在同一空中, 得到全为叉的情况

Growth Function for Convex Sets (1/2)



convex region in blue



non-convex region

- $\mathcal{X} = \mathbb{R}^2$ (two dimensional)
- \mathcal{H} contains h , where $h(\mathbf{x}) = +1$ iff \mathbf{x} in a convex region, -1 otherwise

what is $m_{\mathcal{H}}(N)$?

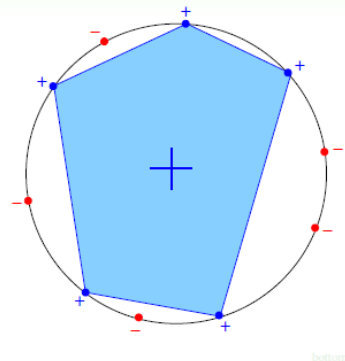
二维集合凸平面

Growth Function for Convex Sets (2/2)

- one possible set of N inputs: $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ on a big circle
- **every dichotomy can be implemented** by \mathcal{H} using a convex region slightly extended from contour of positive inputs

$$m_{\mathcal{H}}(N) = 2^N$$

- call those N inputs '**shattered**' by \mathcal{H}



$m_{\mathcal{H}}(N) = 2^N \iff$
exists N inputs that can be shattered

shattered: N 个点所有可能的分类情况都能够被 hypotheses set 覆盖

Consider positive **and negative** rays as \mathcal{H} , which is equivalent to the perceptron hypothesis set in 1D. The hypothesis set is often called '**decision stump**' to describe the shape of its hypotheses. What is the growth function $m_{\mathcal{H}}(N)$?

1 N

2 $N + 1$

3 $2N$

4 2^N

Reference Answer: ③

Two dichotomies when threshold in each of the $N - 1$ 'internal' spots; two dichotomies for the all-○ and all-× cases.

$$2(N - 1) + 2 = 2N$$

N 个点中有 $N-1$ 个空，射线有正有负，所以有 $2(N - 1)$ 种，除此之外还有 N 个点之外的两个空，得到全正全负两种情况，所以一共有 $2(N - 1) + 2 = 2N$ 种分类

5.4 Break Point

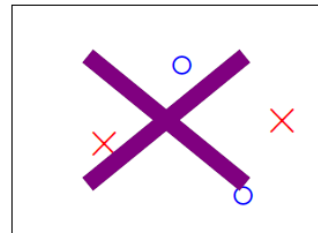
Break Point of \mathcal{H}

what do we know about 2D perceptrons now?

three inputs: 'exists' shatter;
four inputs, 'for all' no shatter

if no k inputs can be shattered by \mathcal{H} ,
call k a **break point** for \mathcal{H}

- $m_{\mathcal{H}}(k) < 2^k$
- $k + 1, k + 2, k + 3, \dots$ also break points!
- will study **minimum break point** k



2D perceptrons: **break point at 4**

2D perceptrons, 我们之前分析了3个点，可以做出8种所有的dichotomy，而4个点，就无法做出所有16个点的dichotomy了。所以，我们就把4称为2D perceptrons的break point（5、6、7等都是break point）。令有 k 个点，如果 k 大于等于break point时，它的成长函数一定小于2的 k 次方。

The Four Break Points

- positive rays: $m_{\mathcal{H}}(N) = N + 1 = O(N)$
break point at 2
- positive intervals: $m_{\mathcal{H}}(N) = \frac{1}{2}N^2 + \frac{1}{2}N + 1 = O(N^2)$
break point at 3
- convex sets: $m_{\mathcal{H}}(N) = 2^N$
no break point
- 2D perceptrons: $m_{\mathcal{H}}(N) < 2^N$ in some cases
break point at 4

conjecture:

- no break point: $m_{\mathcal{H}}(N) = 2^N$ (sure!)
- break point k : $m_{\mathcal{H}}(N) = O(N^{k-1})$

excited? wait for next lecture :-)

Fun Time

Consider positive **and negative** rays as \mathcal{H} , which is equivalent to the perceptron hypothesis set in 1D. As discussed in an earlier quiz question, the growth function $m_{\mathcal{H}}(N) = 2N$. What is the minimum break point for \mathcal{H} ?

① 1

② 2

③ 3

④ 4

Reference Answer: ③

At $k = 3$, $m_{\mathcal{H}}(k) = 6$ while $2^k = 8$.

- positive rays: $m_{\mathcal{H}}(N) = N + 1$
 $\circ \times$ $m_{\mathcal{H}}(2) = 3 < 2^2$: break point at 2
- positive intervals: $m_{\mathcal{H}}(N) = \frac{1}{2}N^2 + \frac{1}{2}N + 1$
 $\circ \times \circ$ $m_{\mathcal{H}}(3) = 7 < 2^3$: break point at 3
- convex sets: $m_{\mathcal{H}}(N) = 2^N$
 $\circ \quad \times$
 $\times \quad \circ$ $m_{\mathcal{H}}(N) = 2^N$ always: no break point
- 2D perceptrons: $m_{\mathcal{H}}(N) < 2^N$ in some cases
 $\times \quad \circ \quad \times$ $m_{\mathcal{H}}(4) = 14 < 2^4$: break point at 4

what 'must be true' when **minimum break point $k = 2$**

- $N = 1$: every $m_{\mathcal{H}}(N) = 2$ by definition
- $N = 2$: every $m_{\mathcal{H}}(N) < 4$ by definition
 (so **maximum possible = 3**)

如果 $k=2$ ，那么当 N 取不同值的时候，计算其成长函数 $m_{\mathcal{H}}(N)$ 是多少。很明显，当 $N=1$ 时， $m_{\mathcal{H}}(N) = 2$ ；当 $N=2$ 时，由break point为2可知，任意两点都不能被shattered（shatter的意思是对 N 个点，能够分解为种dichotomies）； $m_{\mathcal{H}}(N)$ 最大值只能是3；当 $N=3$ 时，简单绘图分析可得其 $m_{\mathcal{H}}(N) = 4$ ，即最多只有4种dichotomies。

maximum possible $m_{\mathcal{H}}(N)$ when $N = 3$ and $k = 2$?

4 dichotomies, shatter any two points? **yes**

x_1	x_2	x_3
\circ	\circ	\circ
\circ	\circ	\times
\circ	\times	\circ
\circ	\times	\times

maximum possible $m_{\mathcal{H}}(N)$ when $N = 3$ and $k = 2$?

4 dichotomies, shatter any two points? **no**

x_1	x_2	x_3
\circ	\circ	\circ
\circ	\circ	\times
\circ	\times	\circ
\times	\circ	\circ

maximum possible $m_{\mathcal{H}}(N)$ when $N = 3$ and $k = 2$?

5 dichotomies, shatter any two points? **yes**

x_1	x_2	x_3
○	○	○
○	○	×
○	×	○
×	○	○
×	○	×

maximum possible $m_{\mathcal{H}}(N)$ when $N = 3$ and $k = 2$?

5 dichotomies, shatter any two points? **yes**

x_1	x_2	x_3
○	○	○
○	○	×
○	×	○
×	○	○
×	×	○

maximum possible $m_{\mathcal{H}}(N)$ when $N = 3$ and $k = 2$?

5 dichotomies, shatter any two points? **yes**

x_1	x_2	x_3
○	○	○
○	○	×
○	×	○
×	○	○
×	×	×

Restriction of Break Point (2/2)

what 'must be true' when **minimum break point** $k = 2$

- $N = 1$: every $m_{\mathcal{H}}(N) = 2$ by definition
- $N = 2$: every $m_{\mathcal{H}}(N) < 4$ by definition
(so **maximum possible** = 3)
- $N = 3$: **maximum possible** = 4 $\ll 2^3$

—break point k **restricts maximum possible** $m_{\mathcal{H}}(N)$ **a lot** for $N > k$

idea: $m_{\mathcal{H}}(N)$
 \leq maximum possible $m_{\mathcal{H}}(N)$ given k
 $\leq poly(N)$

当 $N > K$ 时, break point K 限制了 $m_{\mathcal{H}}(N)$ 的大小, 所以如果给定了 N 和 K , 只要证明 $m_{\mathcal{H}}(N)$ 的最大的上界是多项式, 则根据霍夫丁不等式, 就能用 $m_{\mathcal{H}}(N)$ 代替 M , 得到机器学习是可行的。

When minimum break point $k = 1$, what is the maximum possible $m_{\mathcal{H}}(N)$ when $N = 3$?

① 1

② 2

③ 4

④ 8

Reference Answer: ①

Because $k = 1$, the hypothesis set cannot even shatter one point. Thus, every 'column' of the table cannot contain both \circ and \times . Then, after including the first dichotomy, it is not possible to include any other different dichotomy. Thus, the maximum possible $m_{\mathcal{H}}(N)$ is 1.

x_1	x_2	x_3
\circ	\times	\circ
\circ	\times	\times

6.2 Bounding Function: Basic Cases

Bounding Function

bounding function $B(N, k)$:

maximum possible $m_{\mathcal{H}}(N)$ when break point = k

- combinatorial quantity:
maximum number of length- N vectors with (\circ, \times)
while '**no shatter**' any **length- k** subvectors
- irrelevant of the details of \mathcal{H}
e.g. $B(N, 3)$ bounds both
 - positive intervals ($k = 3$)
 - 1D perceptrons ($k = 3$)

new goal: $B(N, k) \leq \text{poly}(N)$?

Bounding Function($B(N, k)$):当break point= k 时, 成长函数 $m_{\mathcal{H}}(N)$ 最多有多少种Dichotomy的可能. 这里我们不去管成长函数到底长什么样子, 而只需要关心的是在排列组合上, 到底可以做出多少种排列组合。

Table of Bounding Function (1/4)

$B(N, k)$		k						
		1	2	3	4	5	6	...
N	1							
	2		3					
	3		4					
	4							
	5							
	6							
	\vdots							

Known

- $B(2, 2) = 3$ (maximum < 4)
- $B(3, 2) = 4$ ('pictorial' proof previously)

Table of Bounding Function (2/4)

		k						
$B(N, k)$		1	2	3	4	5	6	...
N	1	1						
	2	1	3					
	3	1	4					
	4	1						
	5	1						
	6	1						
	⋮	⋮						

Known

- $B(N, 1) = 1$ (see previous quiz)

Table of Bounding Function (3/4)

		k						
$B(N, k)$		1	2	3	4	5	6	...
N	1	1	2	2	2	2	2	...
	2	1	3	4	4	4	4	...
	3	1	4		8	8	8	...
	4	1				16	16	...
	5	1					32	...
	6	1						...
	⋮	⋮						

Known

- $B(N, k) = 2^N$ for $N < k$
—including all dichotomies not violating ‘breaking condition’

Table of Bounding Function (4/4)

		k						
$B(N, k)$		1	2	3	4	5	6	...
N	1	1	2	2	2	2	2	...
	2	1	3	4	4	4	4	...
	3	1	4	7	8	8	8	...
	4	1			15	16	16	...
	5	1				31	32	...
	6	1					63	...
	\vdots	\vdots						\ddots

Known

- $B(N, k) = 2^N - 1$ for $N = k$
 —removing a single dichotomy satisfies 'breaking condition'

6.3 Bounding Function: Inductive Cases

'Achieving' Dichotomies of $B(4, 3)$

after checking all 2^{2^4} sets of dichotomies, the winner is ...

	x_1	x_2	x_3	x_4
01	○	○	○	○
02	×	○	○	○
03	○	×	○	○
04	○	○	×	○
05	○	○	○	×
06	×	×	○	×
07	×	○	×	○
08	×	○	○	×
09	○	×	×	○
10	○	×	○	×
11	○	○	×	×

		k					
$B(N, k)$		1	2	3	4	5	6
N	1	1	2	2	2	2	2
	2	1	3	4	4	4	4
	3	1	4	7	8	8	8
	4	1		11	15	16	16
	5	1				31	32
	6	1					63

how to reduce $B(4, 3)$ to $B(3, ?)$ cases?

首先，把 $B(4, 3)$ 所有情况写下来，共有11组。也就是说再加一种dichotomy，任意三点都能被shattered，11是极限。

Reorganized Dichotomies of $B(4, 3)$

after checking all 2^{2^4} sets of dichotomies, **the winner is ...**

	x_1	x_2	x_3	x_4		x_1	x_2	x_3	x_4
01	○	○	○	○	01	○	○	○	○
02	×	○	○	○	05	○	○	○	×
03	○	×	○	○	02	×	○	○	○
04	○	○	×	○	08	×	○	○	×
05	○	○	○	×	03	○	×	○	○
06	×	×	○	×	10	○	×	○	×
07	×	○	×	○	04	○	○	×	○
08	×	○	○	×	11	○	○	×	×
09	○	×	×	○	06	×	×	○	×
10	○	×	○	×	07	×	○	×	○
11	○	○	×	×	09	○	×	×	○

orange: pair; purple: single

对这11种dichotomy分组，目前分成两组，分别是orange和purple，orange的特点是， x_1, x_2 和 x_3 是一致的， x_4 不同并成对，例如1和5，2和8等，purple则是单一的， x_1, x_2, x_3 都不同，如6,7,9三组。

将orange去掉 x_4 后去重得到4个不同vector，记为 α ，相应的purple记为 β ， $B(4, 3) = 2\alpha + \beta$ ，因为 $B(4, 3)$ 要求任意三点是不能shatter的，所以，所以由 α, β 构成的三点组合也不能shatter，即 $\alpha + \beta \leq B(3, 3)$

Estimating Part of $B(4, 3)$ (1/2)

$$B(4, 3) = 11 = 2\alpha + \beta$$

	x_1	x_2	x_3
α	○	○	○
	×	○	○
	○	×	○
	○	○	×
β	×	×	○
	×	○	×
	○	×	×

- $\alpha + \beta$: dichotomies on (x_1, x_2, x_3)
- $B(4, 3)$ 'no shatter' any 3 inputs
 $\Rightarrow \alpha + \beta$ 'no shatter' any 3

	x_1	x_2	x_3	x_4
2α	○	○	○	○
	○	○	○	×
	×	○	○	○
	×	○	○	×
	○	×	○	○
	○	×	○	×
	○	○	×	○
	○	○	×	×
β	×	×	○	×
	×	○	×	○
	○	×	×	○

$$\alpha + \beta \leq B(3, 3)$$

另一方面，由于中 x_4 是成对存在的，且是不能被任意三点shatter的，则能推导出 α 是不能被任意两点shatter的。这是因为，如果 α 是不能被任意两点shatter，而 x_4 又是成对存在的，那么 x_1, x_2, x_3, x_4 组成的必然能被三个点shatter。这就违背了条件的设定，所以 $\alpha + \beta \leq B(3, 2)$

Estimating Part of $B(4, 3)$ (2/2)

$$B(4, 3) = 11 = 2\alpha + \beta$$

	x_1	x_2	x_3
α	○	○	○
	×	○	○
	○	×	○
	○	○	×

- α : dichotomies on (x_1, x_2, x_3)
with x_4 **paired**
- $B(4, 3)$ 'no shatter' any 3 inputs
 $\Rightarrow \alpha$ 'no shatter' any 2

	x_1	x_2	x_3	x_4
2α	○	○	○	○
	○	○	○	×
	×	○	○	○
	×	○	○	×
	○	×	○	○
	○	×	○	×
	○	○	×	○
	○	○	×	×
β	×	×	○	×
	×	○	×	○
	○	×	×	○

$$\alpha \leq B(3, 2)$$

由此得出 $B(4, 3)$ 与 $B(3, x)$ 的关系：

$$\begin{aligned}
 B(4, 3) &= 2\alpha + \beta \\
 \alpha + \beta &\leq B(3, 3) \\
 \alpha &\leq B(3, 2) \\
 \Rightarrow B(4, 3) &\leq B(3, 3) + B(3, 2)
 \end{aligned}$$

归纳为公式如下：

$$\begin{aligned}
 B(N, k) &= 2\alpha + \beta \\
 \alpha + \beta &\leq B(N-1, k) \\
 \alpha &\leq B(N-1, k-1) \\
 \Rightarrow B(N, k) &\leq B(N-1, k) + B(N-1, k-1)
 \end{aligned}$$

根据公式填表：

$B(N, k)$		k					
		1	2	3	4	5	6
N	1	1	2	2	2	2	2
	2	1	3	4	4	4	4
	3	1	4	7	8	8	8
	4	1	≤ 5	11	15	16	16
	5	1	≤ 6	≤ 16	≤ 26	31	32
	6	1	≤ 7	≤ 22	≤ 42	≤ 57	63

根据递推公式，推导出 $B(N, K)$ 满足下列不等式：

$$B(N, k) \leq \underbrace{\sum_{i=0}^{k-1} \binom{N}{i}}_{\text{highest term } N^{k-1}}$$

证明：

(1)当 $N=0$ 时

$$B(0, k) \leq \sum_{i=0}^{k-1} C_0^i$$

$$B(0, k-1) \leq \sum_{i=0}^{k-2} C_0^i$$

(2)假设下列两式成立：

$$B(N-1, k) = \sum_{i=0}^{k-1} C_{N-1}^i$$

$$B(N-1, k-1) = \sum_{i=0}^{k-2} C_{N-1}^i$$

则:

$$B(N, k) \leq B(N-1, k) + B(N-1, k-1)$$

$$B(N-1, k) \leq B(N-2, k) + B(N-2, k-1)$$

$$B(N-2, k) \leq B(N-3, k) + B(N-3, k-1)$$

\vdots

$$B(2, k) \leq B(1, k) + B(1, k-1)$$

$$B(1, k) \leq B(0, k) + B(0, k-1)$$

证明 $B(1, k)$ 的值:

$$C_n^m = C_{n-1}^m + C_{n-1}^{m-1} \quad (6-1)$$

$$B(0, k) = C_0^0 + C_0^1 + C_0^2 + C_0^3 \dots + C_0^{k-2} + C_0^{k-1}$$

$$B(0, k-1) = C_0^0 + C_0^1 + C_0^2 \dots + C_0^{k-3} + C_0^{k-2}$$

由公式6-1可得:

$$C_0^1 + C_0^0 = C_1^1$$

$$C_0^2 + C_0^1 = C_1^2$$

\vdots

$$C_0^{k-2} + C_0^{k-3} = C_1^{k-2}$$

$$C_0^{k-1} + C_0^{k-2} = C_1^{k-1}$$

所以

$$B(0, k) + B(0, k-1) = C_0^0 + C_1^1 + C_1^2 + \dots + C_1^{k-2} + C_1^{k-1} \quad (6-2)$$

$$= C_1^0 + C_1^1 + C_1^2 + \dots + C_1^{k-2} + C_1^{k-1} \quad (6-3)$$

$$\text{所以 } B(1, k) = C_1^0 + C_1^1 + C_1^2 + \dots + C_1^{k-2} + C_1^{k-1}$$

$$\text{同理: } B(1, k-1) = C_1^0 + C_1^1 + C_1^2 + \dots + C_1^{k-2}$$

$$\text{根据公式6-1可得: } B(2, k) = C_2^0 + C_2^1 + C_2^2 + \dots + C_2^{k-2} + C_2^{k-1}$$

$$\text{以此类推 } B(N, k) = C_N^0 + C_N^1 + C_N^2 + C_N^3 + \dots + C_N^{k-2} + C_N^{k-1}$$

(3)由(1)、(2)可得

$$B(N, k) = \sum_{i=0}^{k-1} C_N^i$$

成立

上述不等式的右边是最高阶为k-1的N项多项式，也就是说成长函数的上界B(N,K)的上界满足多项式分布poly(N),

6.4 A Pictorial Proof

want:

$$\mathbb{P}\left[\exists h \in \mathcal{H} \text{ s.t. } |E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon\right] \leq 2 m_{\mathcal{H}}(N) \cdot \exp\left(-2 \epsilon^2 N\right)$$

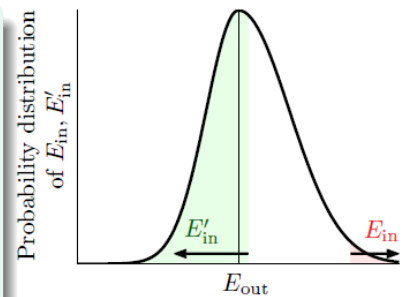
actually, when N large enough,

$$\mathbb{P}\left[\exists h \in \mathcal{H} \text{ s.t. } |E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon\right] \leq 2 \cdot 2 m_{\mathcal{H}}(2N) \cdot \exp\left(-2 \cdot \frac{1}{16} \epsilon^2 N\right)$$

Step 1: Replace E_{out} by E'_{in}

$$\begin{aligned} & \frac{1}{2} \mathbb{P}\left[\exists h \in \mathcal{H} \text{ s.t. } |E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon\right] \\ & \leq \mathbb{P}\left[\exists h \in \mathcal{H} \text{ s.t. } |E_{\text{in}}(h) - E'_{\text{in}}(h)| > \frac{\epsilon}{2}\right] \end{aligned}$$

- $E_{\text{in}}(h)$ finitely many, $E_{\text{out}}(h)$ infinitely many
—replace the evil E_{out} first
- how? sample verification set \mathcal{D}' of size N to calculate E'_{in}
- BAD h of $E_{\text{in}} - E_{\text{out}}$
 $\xRightarrow{\text{probably}}$ BAD h of $E_{\text{in}} - E'_{\text{in}}$



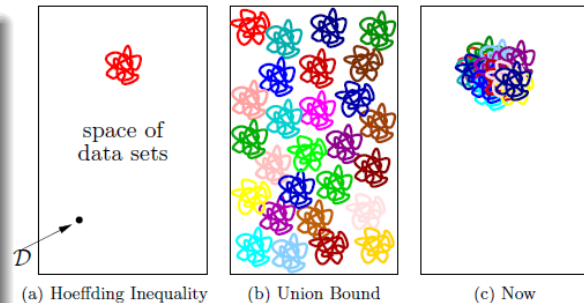
evil E_{out} removed by
verification with 'ghost data'

取另外的 N 个点拿来验证，用来估计 E_{out} 的值 我把这另外 N 个点叫 \mathcal{D}' ，OK，那这个用 \mathcal{D}' 估计出来的值我叫 E'_{in} ，当发生bad data时， E_{out} 与 E_{in} 差很远，如果我有很大的机会抽到 E'_{in} ， E'_{in} 与 E_{in} 也相差很远，那么便可以换一种坏事情： E'_{in} 与 E_{in} 相差很远。
式子左边的 $\frac{1}{2}$ 对应有很大的几率抽中 E'_{in} 与 E_{in} 相差很远

Step 2: Decompose \mathcal{H} by Kind

$$\begin{aligned} \text{BAD} &\leq 2\mathbb{P}\left[\exists h \in \mathcal{H} \text{ s.t. } |E_{\text{in}}(h) - E'_{\text{in}}(h)| > \frac{\epsilon}{2}\right] \\ &\leq 2m_{\mathcal{H}}(2N)\mathbb{P}\left[\text{fixed } h \text{ s.t. } |E_{\text{in}}(h) - E'_{\text{in}}(h)| > \frac{\epsilon}{2}\right] \end{aligned}$$

- E_{in} with \mathcal{D} , E'_{in} with \mathcal{D}'
—now $m_{\mathcal{H}}$ comes to play
- how? infinite \mathcal{H} becomes
 $|\mathcal{H}(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{x}'_1, \dots, \mathbf{x}'_N)|$
kinds
- union bound on $m_{\mathcal{H}}(2N)$ kinds

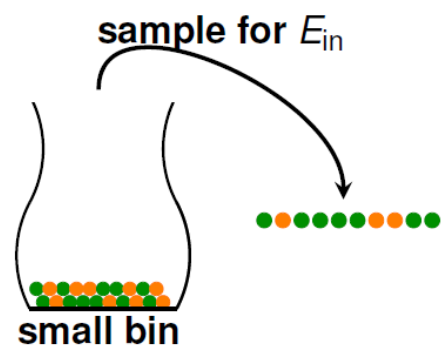


use $m_{\mathcal{H}}(2N)$ to calculate BAD-overlap properly

Step 3: Use Hoeffding without Replacement

$$\begin{aligned} \text{BAD} &\leq 2m_{\mathcal{H}}(2N)\mathbb{P}\left[\text{fixed } h \text{ s.t. } |E_{\text{in}}(h) - E'_{\text{in}}(h)| > \frac{\epsilon}{2}\right] \\ &\leq 2m_{\mathcal{H}}(2N) \cdot 2 \exp\left(-2\left(\frac{\epsilon}{4}\right)^2 N\right) \end{aligned}$$

- consider bin of $2N$ examples,
choose N for E_{in} , leave others for E'_{in}
 $|E_{\text{in}} - E'_{\text{in}}| > \frac{\epsilon}{2} \Leftrightarrow \left|E_{\text{in}} - \frac{E_{\text{in}} + E'_{\text{in}}}{2}\right| > \frac{\epsilon}{4}$
- so? just 'smaller bin', 'smaller ϵ ', and
Hoeffding without replacement



use Hoeffding after zooming to fixed h

Vapnik-Chervonenkis (VC) bound:

$$\begin{aligned} & \mathbb{P}\left[\exists h \in \mathcal{H} \text{ s.t. } |E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon\right] \\ & \leq 4m_{\mathcal{H}}(2N) \exp\left(-\frac{1}{8}\epsilon^2 N\right) \end{aligned}$$

- replace E_{out} by E'_{in}
- decompose \mathcal{H} by kind
- use Hoeffding without replacement