

Lecture 7. Generalisation with Finite VC Dimension

COMP90051 Statistical Machine Learning

Lecturer: Jean Honorio



THE UNIVERSITY OF
MELBOURNE

This lecture

- Motivation
- Growth function
 - * Considering patterns of labels possible on a data set
 - * Gives good generalisation bounds provided possible patterns don't grow too fast in the data set size
- Vapnik-Chervonenkis (VC) dimension
 - * Max number of points that can be labelled in all ways
 - * Beyond this point, growth function is polynomial in data set size
 - * Leads to famous VC generalisation theorem

Motivation

...from last lecture

A Countably Finite Model Class

- Consider we have 2 features and a countably finite set \mathcal{F} of classifiers, containing:

$$f(x) = \text{sgn}(x_1 + x_2) = \begin{cases} +1, & \text{if } x_1 + x_2 > 0 \\ -1, & \text{if } x_1 + x_2 \leq 0 \end{cases}$$

$$f(x) = \text{sgn}(x_1 - x_2)$$

$$f(x) = \text{sgn}(-x_1 + x_2)$$

$$f(x) = \text{sgn}(-x_1 - x_2)$$

$$f(x) = \text{sgn}(x_1)$$

$$f(x) = \text{sgn}(-x_1)$$

$$f(x) = \text{sgn}(x_2)$$

$$f(x) = \text{sgn}(-x_2)$$

- Here $|\mathcal{F}| = 8$

Empirical Risk Minimisation

- Training data $\mathbf{D} = \{\mathbf{x}_1, y_1, \dots, \mathbf{x}_n, y_n\}$ is a random variable!
 - * (\mathbf{x}_i, y_i) i.i.d. with distribution P (unknown)

- The empirical risk of a classifier f for loss l is

$$\hat{R}_{\mathbf{D}}[f] = \frac{1}{n} \sum_{i=1}^n l(y_i, f(\mathbf{x}_i))$$

- ERM: $\hat{f}_{\mathbf{D}}$ minimises the empirical risk

$$\hat{f}_{\mathbf{D}} = \operatorname{argmin}_{f \in \mathcal{F}} \hat{R}_{\mathbf{D}}[f]$$

Go through all the $|\mathcal{F}| = 8$ classifiers and choose the best for data \mathbf{D}

- Given f and n samples in \mathbf{D} , we can compute $\hat{R}_{\mathbf{D}}[f]$

True Risk

- The true risk is the expected value of the loss l
 - * Intuitively speaking, the true risk is the empirical risk when using an infinite number of samples

- The true risk of a classifier f for loss l is

$$R[f] = \mathbb{E} l(Y, f(X)) = \int l(Y, f(X)) P(X, Y) dX dY$$

aka generalisation error
(expected test error) for

$$l(y, y') = \begin{cases} 1, & \text{if } y \neq y' \\ 0, & \text{if } y = y' \end{cases}$$

- Given f , we cannot compute $R[f]$ because the data distribution P is unknown

Generalisation Theorem

- For a finite model class \mathcal{F} , without knowing the data distribution P , with probability $\geq 1 - \delta$ over the choice of the training set D of n i.i.d. samples

$$R[\hat{f}_D] \leq \hat{R}_D[\hat{f}_D] + \sqrt{\frac{\log |\mathcal{F}| + \log(1/\delta)}{2n}}$$

move certain more pay

$\delta \downarrow \rightarrow \frac{1}{\delta} \uparrow \rightarrow \text{bound} \uparrow$

We cannot compute $R[f]$, but we can bound it!

- The proof-sketch required upper bounding

$$\max_{f \in \mathcal{F}} \varphi_D[f] = \max_{f \in \mathcal{F}} (R[f] - \hat{R}_D[f])$$

Non-(Countably Finite) Model Class?

- Finite model class
 - * Bounding uniform deviation with union bound and Hoeffding's inequality
- Consider we have 2 features and an uncountable set \mathcal{F} of classifiers, containing for all $w_1 \in \mathbb{R}$, $w_2 \in \mathbb{R}$:

$$f(x) = \text{sgn}(w_1 x_1 + w_2 x_2)$$
- As before, still requires upper bounding

$$\sup_{f \in \mathcal{F}} (R[f] - \hat{R}_D[f])$$

Handwritten notes:
 Above the sup: 不可数 infinite.
 Below the sup: max 可数

Mini Summary

- No good for general (countably infinite and uncountable) cases
- Need another fundamentally new idea

Next: Organising analysis around patterns of labels possible on any data set

Growth Function

*Focusing on the size of model families
on data samples*

Example: Decision stumps

- Consider a dataset of 6 samples, each with a single continuous feature (x) and label (y)

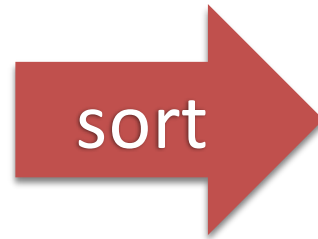
x	y
0	+1
4	-1
-2	+1
1	+1
-3	-1
2	-1

- We would like to find a threshold β , and then classify all samples with feature value x above β as +1, and feature value x below β as -1 (or viceversa)

Example: Decision stumps

- Lets sort with respect to x

x	y
0	+1
4	-1
-2	+1
1	+1
-3	-1
2	-1



x	y
-3	-1
-2	+1
0	+1
1	+1
2	-1
4	-1

- Lets use the classifier:

$$f(x) = \text{sgn}(x - \beta) = \begin{cases} +1, & \text{if } x > \beta \\ -1, & \text{if } x \leq \beta \end{cases}$$

- How to find the threshold β ? Try all midpoints of x

Example: Decision stumps

- Lets use the classifier:

$$f(x) = \text{sgn}(x - \beta) = \begin{cases} +1, & \text{if } x > \beta \\ -1, & \text{if } x \leq \beta \end{cases}$$

- Count the number of mistakes for all thresholds β

x	y	<i>best threshold</i> $f(x)$ <i>设 甲 是不同数据中点</i>				
		$\beta=-2.5$	$\beta=-1$	$\beta=0.5$	$\beta=1.5$	$\beta=3$
-3	-1	-1	-1	-1	-1	-1
-2	+1	+1	-1	-1	-1	-1
0	+1	+1	+1	-1	-1	-1
1	+1	+1	+1	+1	-1	-1
2	-1	+1	+1	+1	+1	-1
4	-1	+1	+1	+1	+1	+1
# mistakes		2	3	4	5	4

Example: Decision stumps

- Lets use the classifier:

$$f(x) = \text{sgn}(\beta - x) = \begin{cases} +1, & \text{if } x < \beta \\ -1, & \text{if } x \geq \beta \end{cases}$$

- Count the number of mistakes for all thresholds β

x	y	$f(x)$				
		$\beta=-2.5$	$\beta=-1$	$\beta=0.5$	$\beta=1.5$	$\beta=3$
-3	-1	+1	+1	+1	+1	+1
-2	+1	-1	+1	+1	+1	+1
0	+1	-1	-1	+1	+1	+1
1	+1	-1	-1	-1	+1	+1
2	-1	-1	-1	-1	-1	+1
4	-1	-1	-1	-1	-1	-1
# mistakes		4	3	2	1	2

Example: Decision stumps

- Thus our best decision stump classifier is

$$f(x) = \text{sgn}(1.5 - x) = \begin{cases} +1, & \text{if } x < 1.5 \\ -1, & \text{if } x \geq 1.5 \end{cases}$$

- We consider all classifiers of the form (for all $\beta \in \mathbb{R}$)

$$f(x) = \text{sgn}(x - \beta) = \begin{cases} +1, & \text{if } x > \beta \\ -1, & \text{if } x \leq \beta \end{cases}$$

$$f(x) = \text{sgn}(\beta - x) = \begin{cases} +1, & \text{if } x < \beta \\ -1, & \text{if } x \geq \beta \end{cases}$$

- Although these are simple classifiers, the set of decision stump classifiers \mathcal{F} is uncountable (there are as “many” as real values)

△???

β 的取值有无数个

Example: Growth function of Decision stumps

- Consider all possible ways we can classify data

$$\textcircled{1} f(x) = \text{sgn}(x - \beta) = \begin{cases} +1, & \text{if } x > \beta \\ -1, & \text{if } x \leq \beta \end{cases}$$

$$\textcircled{2} f(x) = \text{sgn}(\beta - x) = \begin{cases} +1, & \text{if } x < \beta \\ -1, & \text{if } x \geq \beta \end{cases}$$

①的-1和②的+1结果一样，去掉；②的-1和①的+1一样，去掉。

x	$f(x)$					
	$\beta=-2.5$	$\beta=-1$	$\beta=0.5$	$\beta=1.5$	$\beta=3$	$\beta=\infty$
-3	-1	-1	-1	-1	-1	-1
-2	+1	-1	-1	-1	-1	-1
0	+1	+1	-1	-1	-1	-1
1	+1	+1	+1	-1	-1	-1
2	+1	+1	+1	+1	-1	-1
4	+1	+1	+1	+1	+1	-1

x	$f(x)$					
	$\beta=-2.5$	$\beta=-1$	$\beta=0.5$	$\beta=1.5$	$\beta=3$	$\beta=\infty$
-3	+1	+1	+1	+1	+1	+1
-2	-1	+1	+1	+1	+1	+1
0	-1	-1	+1	+1	+1	+1
1	-1	-1	-1	+1	+1	+1
2	-1	-1	-1	-1	+1	+1
4	-1	-1	-1	-1	-1	+1

- A **dichotomy** (in blue) is one way of classifying the 6 samples
- We have **12 unique dichotomies**

$2 \times 7 - 2 = 12$ unique
个数的

Dichotomies

- Given dataset $\mathcal{X} = \{x_1, \dots, x_n\}$ of size $|\mathcal{X}| = n$ and a classifier $f \in \mathcal{F}$, a **dichotomy** is the pattern of labels (n -dimensional vector of labels) produced by f on \mathcal{X}

$$(f(x_1), \dots, f(x_n)) \in \{-1, +1\}^n.$$

每个数据点进入模型要+1要-1 产生不同的组合, 一组就是一个 dichotomy.

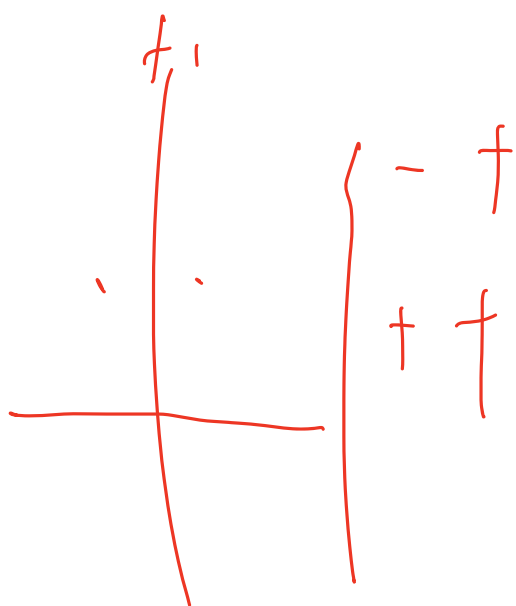
- Unique dichotomies:** unique patterns of labels possible with all classifiers in the model class \mathcal{F}

$$\mathcal{F}(\mathcal{X}) = \{(f(x_1), \dots, f(x_n)) : f \in \mathcal{F}\}$$

- * Even when \mathcal{F} infinite, $|\mathcal{F}(\mathcal{X})| \leq 2^n$ (why?) $(1+1)^n$
- * For \mathcal{F} countably finite, $|\mathcal{F}(\mathcal{X})| \leq |\mathcal{F}|$ (why?)

以上模型产生的所有可能的标签模式
每个样本两种结果
这是一系列所有可能的结果

最终: 每个 dichotomy 都不同



Growth Function

- The **growth function**

$$S_{\mathcal{F}}(n) = \sup_{|\mathcal{X}|=n} |\mathcal{F}(\mathcal{X})|$$

unique dichotomy 的个数
 对 n 样本, \mathcal{F} 可生成的最多的 unique dichotomy 数量
 (unique dichotomy 的上界; 与样本点数量有关的函数)
- is the maximum number of label patterns achievable by classifiers in the model class \mathcal{F} for any set of n samples.
 - * Even when \mathcal{F} infinite, $S_{\mathcal{F}}(n) \leq 2^n$ (why?)
 - * For \mathcal{F} countably finite, $S_{\mathcal{F}}(n) \leq |\mathcal{F}|$ (why?)

$S_{\mathcal{F}}(n)$ 衡量了 \mathcal{F} 的表达能力 (容量): 告诉我们 \mathcal{F} 内的分类器可以对多少不同尺寸的样本 n 进行分类

Example: Growth function of Decision stumps

- In general, the set of decision stump classifiers lead to $2n$ unique dichotomies for n samples
 - We classify the n samples as -1's followed by +1's (先 +1 -1) $(n+1)$
 - We also classify the n samples as +1's followed by -1's (再反过来) $(n+1)$
- Thus, $S_{\mathcal{F}}(n) = 2n$

有所重复的 $(n+1)$

$\rightarrow 2 \times (n+1) - 2 = 2n$
- More complex classifiers would lead to more than $2n$ unique dichotomies for n samples

对 decision stump 来说 $S_{\mathcal{F}}(n)=2n$ ，对更复杂的分类器来说

$S_{\mathcal{F}} \subset \text{Stump}$

Growth-Function Generalisation Theorem

- For a model class \mathcal{F} with growth function $S_{\mathcal{F}}(n)$, without knowing the data distribution P , with probability $\geq 1 - \delta$ over the choice of the training set D of n i.i.d. samples

$$R[\hat{f}_D] \leq \hat{R}_D[\hat{f}_D] + \sqrt{8 \frac{\log S_{\mathcal{F}}(2n) + \log(4/\delta)}{n}}$$

(Proof outside scope of COMP90051)

* $|\mathcal{F}|$ becomes $S_{\mathcal{F}}(2n)$, and few negligible extra constants

* If $S_{\mathcal{F}}(n)$ grows exponentially in n , e.g., $S_{\mathcal{F}}(n) = 2^n$ then $\frac{\log S_{\mathcal{F}}(2n)}{n} = 2$, the bound does not decay with more samples n

这么多 data $n \uparrow \rightarrow$ 希望 bound 更紧 才 学习 好

Mini Summary

- Better to organise families by possible patterns of labels on a data set: the dichotomies of the model class
- Counting possible dichotomies gives the growth function
- Generalisation bound with growth function potentially tackles general (countably infinite and uncountable) families provided growth function is sub-exponential in data size

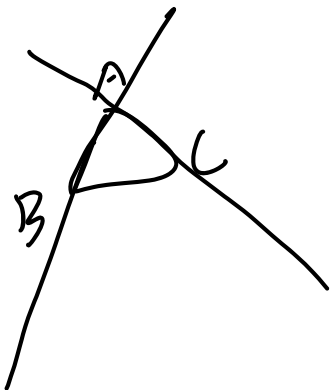
Next: VC dimension for a computable bound on growth functions, with the polynomial behaviour we need! Gives our final VC generalisation bound

The VC dimension

Computable, bounds growth function

Vapnik-Chervonenkis dimension

- The **VC dimension** $VC(\mathcal{F})$ of a model class \mathcal{F} is the largest n such that $S_{\mathcal{F}}(n) = 2^n$. (使得增长函数等于 2^n 的最大的 n)
(等于二分类数量)
- Set of samples $\mathcal{X} = \{x_1, \dots, x_n\}$ are **shattered** by \mathcal{F} if $|\mathcal{F}(\mathcal{X})| = 2^n$, that is, if \mathcal{X} can be classified in all possible ways. (VC 能够打乱出 2^n 种可能组合)
(下中的元素把 \mathcal{X} 的 $\{1, \dots, n\}$ 可能理都表达出来 \Rightarrow shatter 粉碎)
- $VC(\mathcal{F})$ is the size of the largest set of samples shattered by \mathcal{F} .



Shatter: (能够分类的最大样本数量)

VC: 衡量模型容量的指标 $VC \uparrow \Rightarrow$ model 越复杂
 \Rightarrow 能够打乱更多的数据

Example: VC Dimension of Decision Stumps

VC dimension: the largest number n such that $S_{\mathcal{F}}(n) = 2^n$

- Recall that for decision stump classifiers $S_{\mathcal{F}}(n) = 2n$
- Find the maximum n for which $2n = 2^n$
- The VC dimension is $VC(\mathcal{F}) = 2$

n	$2n$	2^n
1	2	2
2	4	4
3	6	8

- For more intuition, see the $2n$ ways of classifying n samples

$n=1$

+1	-1
----	----

$n=2$

+1	+1	-1	-1
+1	-1	+1	-1

$n=3$

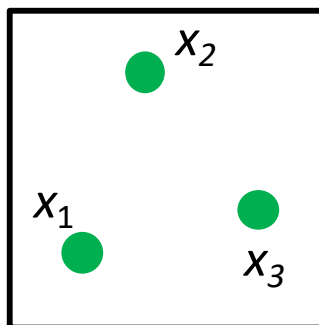
+1	+1	+1	+1	-1	-1	-1	-1
+1	+1	-1	-1	+1	+1	-1	-1
+1	-1	+1	-1	+1	-1	+1	-1

2 ways ($2^3 - 2 \cdot 3 = 2$) of classifying (in red) are not -1's followed by +1's, neither +1's followed by -1's

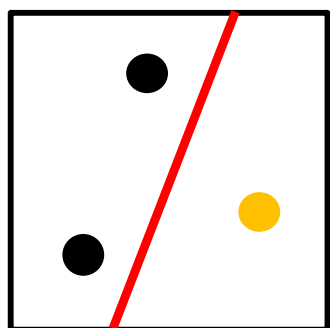
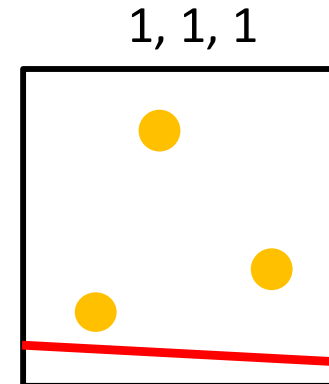
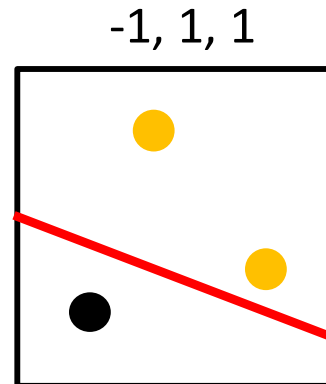
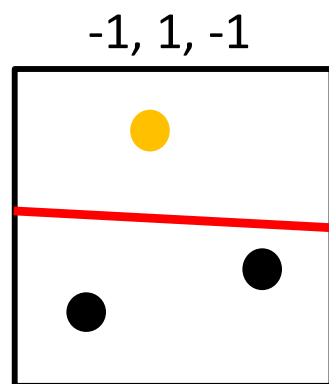
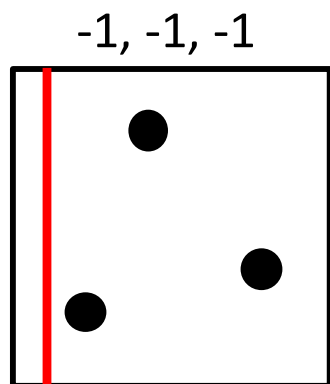
Example 2: Growth function for linear classifiers in 2D

● Black means $f(x)=-1$

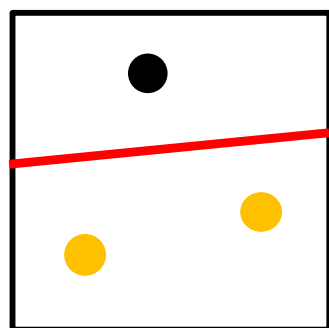
● Yellow means $f(x)=1$



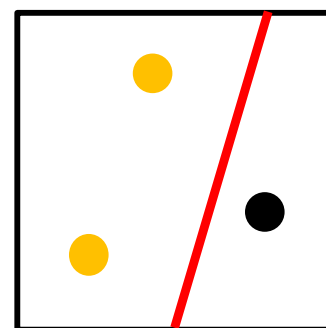
pD
 $S_{\mathcal{F}}(3) = 8$
 最大的 unique dichotomy 数量.



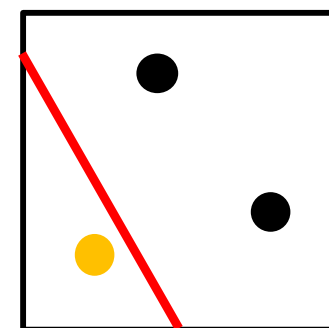
-1, -1, 1



1, -1, 1

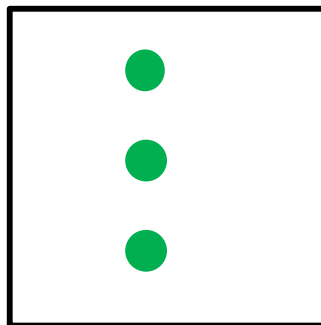


1, 1, -1

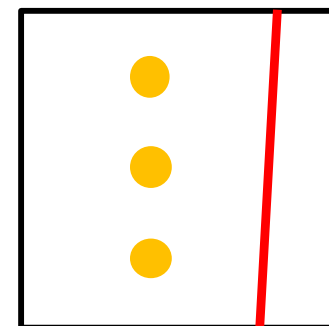
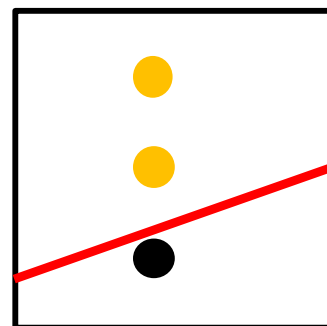
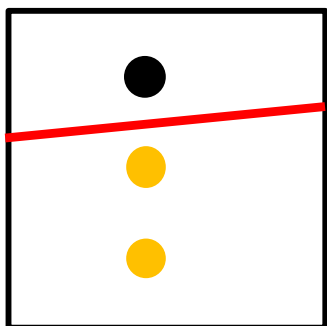
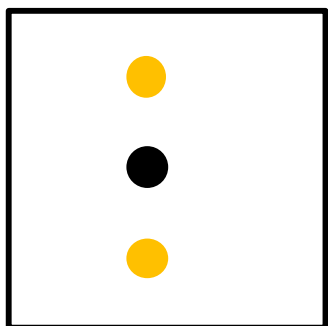
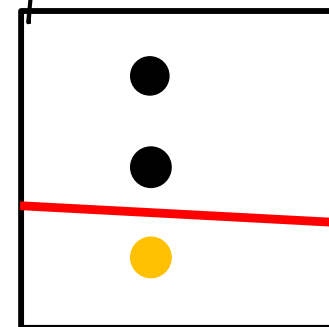
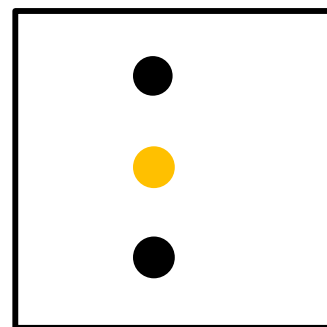
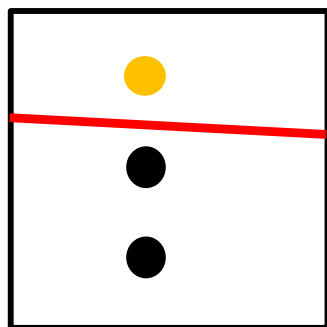
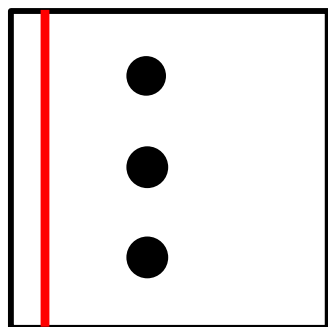


1, -1, -1

Example 2: Growth function for linear classifiers in 2D



The possible patterns should be

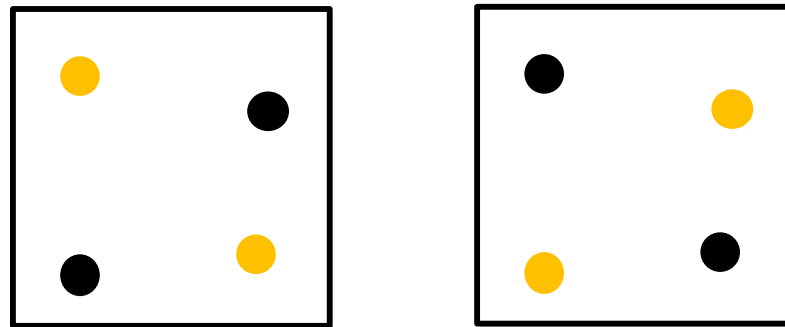


$\sup_{|X|=n} |f(x)| = 8$
 $|F(\mathbf{x})| = 6$
 but still have
 $S_{\mathcal{F}}(3) = 8$

growth function is always superermining

Example 2: Growth function for linear classifiers in 2D

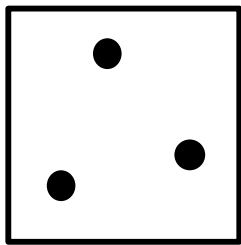
- What about $n = 4$ points?
- Can never produce the criss-cross (XOR) dichotomy



- In fact $S_{\mathcal{F}}(4) = 14 < 2^4$
 grow 小于 2^4 意味着不是所有 possible rays.

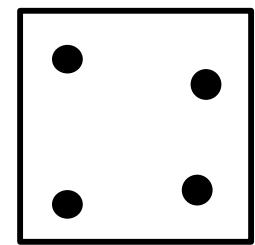
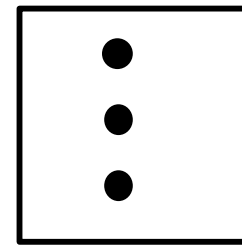
Example 2: VC dimension for linear classifiers in 2D

- Example: linear classifiers in \mathbb{R}^2 , $VC(\mathcal{F}) = 3$



Shattered

可入 shattered 的样本的数量
Not shattered

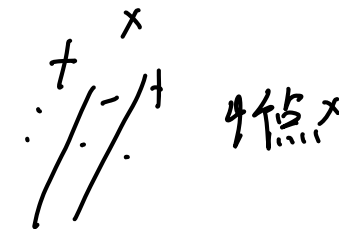


- Guess: VC dimension of linear classifiers in \mathbb{R}^d ?

增函数: 取最大的可入 shattered 的样本数量 n

\mathbb{R}^d 维下的 LR 的 VC 维 : $\Rightarrow d+1$

$d+1 \leq VC \leq d+2$ for hyperplane



Example 3: VC dimension from dichotomies on whole domain?

① 每个输入为 unique

② $\frac{n}{2} VC=4$ 且 $|X|=4 \Rightarrow$ 有 2^4 VC dimension is the points that I can classify in all possible way

x_1	x_2	x_3	x_4
0	0	0	0
0	1	1	0
1	0	0	1
1	1	0	1
0	1	0	0
1	0	1	0
1	1	1	1
0	0	1	1
0	1	0	1
1	1	1	0

③ 但现在不是 16 个 $\Rightarrow VC \neq 4$

- Each row is a dichotomy on entire input domain
- Obtain dichotomies on a subset of samples $\mathcal{X}' \subseteq \{x_1, \dots, x_4\}$ by: drop columns, drop dupe rows
- \mathcal{F} shatters \mathcal{X}' if number of rows is $2^{|\mathcal{X}'|}$

x_1	x_2	x_4
0	0	0
0	1	0
1	0	1
1	1	1
0	1	0
1	0	0
1	1	1
0	0	1
0	1	1
1	1	0

This example:

- ① Dropping column 3 leaves 8 rows behind: \mathcal{F} shatters $\{x_1, x_2, x_4\}$
- Original table has $< 2^4$ rows: \mathcal{F} doesn't shatter more than 3
- $VC(\mathcal{F}) = 3$

Note we're using labels $\{0,1\}$ instead of $\{-1,+1\}$. Why OK?

Sauer-Shelah Lemma

- Consider any model class \mathcal{F} with finite $VC(\mathcal{F})$, and any sample size n . Then

$$S_{\mathcal{F}}(n) \leq \sum_{i=0}^{VC(\mathcal{F})} \binom{n}{i}$$

(Proof outside scope of COMP90051)

原非 bound 中 growth function 不随 n 的增加而 delay.

- Since $\sum_{i=0}^k \binom{n}{i} \leq (n+1)^k$, the above implies
growth function $\leq VC$

$$\log S_{\mathcal{F}}(n) \leq VC(\mathcal{F}) \log(n+1)$$

VC Generalisation Theorem

- For a model class \mathcal{F} with VC dimension $VC(\mathcal{F})$, without knowing the data distribution P , with probability $\geq 1 - \delta$ over the choice of the training set D of n i.i.d. samples

$$R[\hat{f}_D] \leq \hat{R}_D[\hat{f}_D] + \sqrt{8 \frac{VC(\mathcal{F}) \log(2n + 1) + \log(4/\delta)}{n}}$$

有 VC 数

- * Proof-sketch: From the growth-function generalization theorem and since

$$\log S_{\mathcal{F}}(2n) \leq \overbrace{VC(\mathcal{F}) \log(2n + 1)}$$

Structural Risk Minimisation

- Choose the model class \mathcal{F} with best guarantee of generalisation:

$$\underbrace{\hat{R}_D[\hat{f}_D]}_{\text{Large for simple classifiers, small for complex classifiers}} + \underbrace{\sqrt{8 \frac{VC(\mathcal{F}) \log(2n + 1) + \log(4/\delta)}{n}}}_{\text{Small for simple classifiers (small } VC(\mathcal{F}) \text{), large for complex classifiers (large } VC(\mathcal{F}) \text{), Large for small } n \text{ (few samples), small for large } n \text{ (many samples)}}$$

Large for simple classifiers,
small for complex classifiers

Small for simple classifiers (small $VC(\mathcal{F})$),
large for complex classifiers (large $VC(\mathcal{F})$)

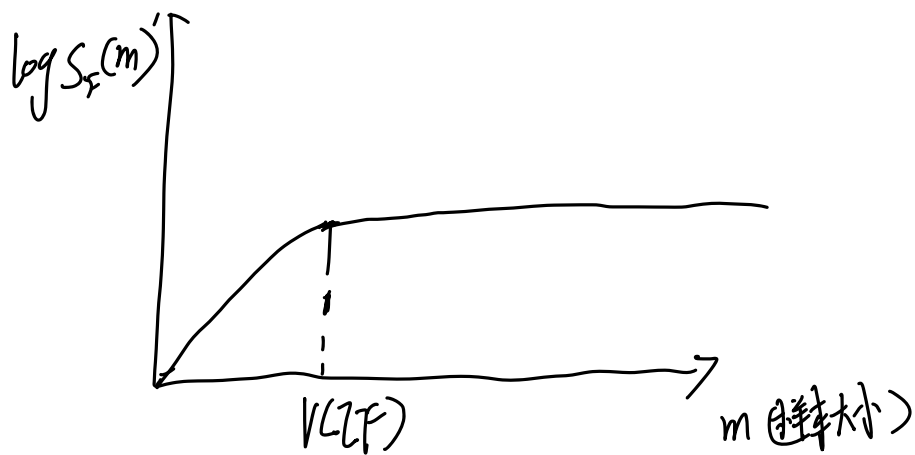
Large for small n (few samples),
small for large n (many samples)

Mini Summary

- VC dimension is the size of the largest set of samples shattered by a model class
 - * It is $d + 1$ for linear classifiers in \mathbb{R}^d
- Sauer-Shelah: The growth function grows only polynomially in the set size beyond the VC dimension
- As a result, VC generalisation bounds true risk and empirical risk deviation by $O(\sqrt{(\text{VC}(\mathcal{F}) \log n)/n})$

Much more...

- Finite VC dimension equivalent to Provably approximately correct (PAC) learning
- VC dimension is not the only tool in learning theory
 - * Some problems might have infinite VC dimension
 - * Other problems beyond classification
- The generalization of some methods require different complexity measures or analysis frameworks, such as:
 - * Fat shattering dimension
 - * Provably approximately correct (PAC) Bayes bounds
 - * Rademacher complexity



随着数据 \uparrow \rightarrow 模型表达能力 \rightarrow 后期收敛
 $VC \uparrow$ $VC \downarrow$

若 unique table 有 9 行, 则 $2^4 = 16 < 9 \Rightarrow$ 当前 model class 不能粉碎 4 \Rightarrow 达不到 16 种 \Rightarrow 所以只能粉碎 3 次

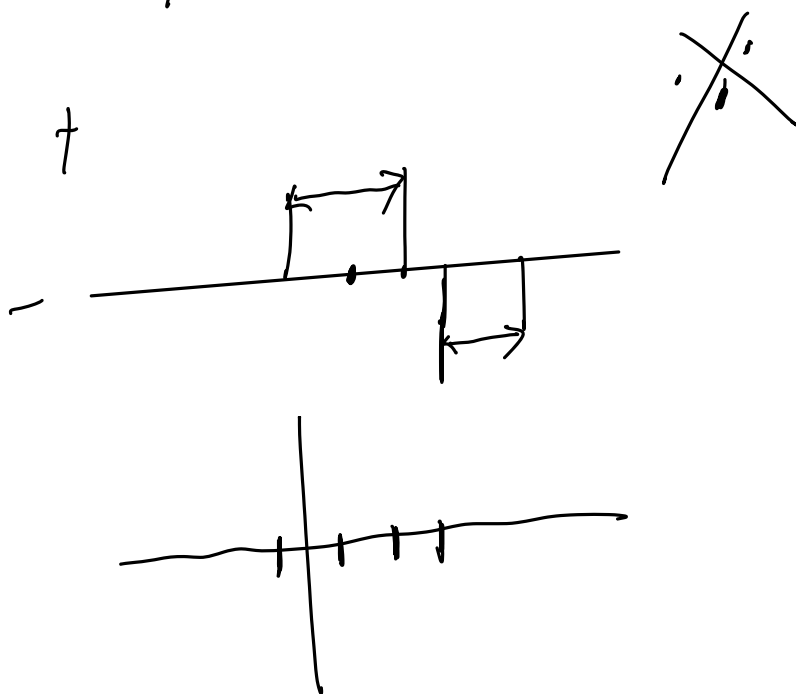
V_C dimension:

① a set of d can be shattered by H @ distance bound

② $d+1$ is upper bound

$H = \{ \text{interval on the real line} \}$

$$X = \mathbb{R}$$



$$H = \text{line}$$

$$X = \mathbb{R}^2$$