

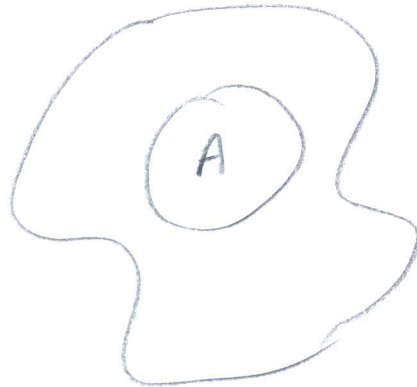
YouTube 12/09/2016

LW: HW question (?)

LW: continue talking about uniform bounds.

Suppose sample space:

- draw i.i.d.



$$P(A) = P(X \in A) \text{ (unknown prob.)}$$

$$P_n(A) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(X_i \in A) \text{ (empirical prob.)}$$

Q: How far is $P_n(A)$ from $P(A)$?

A: Hoeffding

$$P(|P_n(A) - P(A)| > \epsilon) \leq 2e^{-2n\epsilon^2}$$

LW: Guaranteeing this is small one 'thing' at a time does not mean that it holds uniformly over a class (set)

LW: Instead, we want:-

$$P\left(\sup_{A \in \mathcal{A}} |P_n(A) - P(A)| > \epsilon\right) \leq \text{small}$$

- this is a uniform bound over set \mathcal{A} ultimately, we will have a bound like Hoeffding, but with a complex term2. finite classesLet $\mathcal{A} = \{A_1, \dots, A_N\}$ - A finite set of events.

note give

$$P\left(\sup_{A \in \mathcal{A}} |P_n(A) - P(A)| > \epsilon\right) \leq c_1 K(\mathcal{A}) e^{-c_2 n \epsilon^2}$$

$$P\left(\sup_{j=1, \dots, N} |P_n(A_j) - P(A_j)| > \epsilon\right) \quad (1)$$

define B_j be the event that $B_j = \{|P_n(A_j) - P(A_j)| > \epsilon\}$ - event that probability of absolute difference of empirical and population probability of event A_j is greater than ϵ .
(true)

(*) (*)

W: An alternative to (I) is ~~from~~ observing that it's identical to saying at least one of the events B_j is true.

(*) events B_1, B_2, B_3 are NOT disjoint. Hence $P(\cup B_j) \neq \sum_j P(B_j)$

However we can use the union bound (W) (A) - Recall that? - I don't recall (Q)

$$P\left(\sup_{A \in \mathcal{A}} |P_n(A) - P(A)| > \epsilon\right) = P\left(\bigcup_{j=1}^N B_j\right) \leq \sum_{j=1}^N P(B_j)$$

Note: introduction of B_j yields $\sum_{j=1}^N P(B_j) = \sum_{j=1}^N P(|P_n(A_j) - P(A_j)| > \epsilon) = 2Ne^{-2n\epsilon^2}$

Hence apply Hoeffding:-

$$P\left(\sup_{A \in \mathcal{A}} |P_n(A) - P(A)| > \epsilon\right) \leq 2Ne^{-2n\epsilon^2}$$

Hoeffding
(with extra
term denoting
no. of events B_j
or A_j)

W: $2N$ is a 'measure' of complexity of the class/
set of events $A = \{A_1, \dots, A_N\}$ or $B = \{B_1, \dots, B_N\}$.

- yields simplest way to control probability (uniform bounds)

W: A limitation to this; what if class/set \mathcal{A} is infinite i.e. infinite cardinality, bound then goes $\rightarrow \infty$

W: need more sophisticated method than counting

Consider $\mathcal{Q} = \{\text{discs on } \mathbb{R}^2\}$ - infinite no. of discs (how to improve on counting?)

$$P\left(\sup_A |P_n(A) - P(A)| > \epsilon\right)$$

W: Invoke shattering / VC dimension

- Give me a class of sets \mathcal{Q}

(class is possibly infinite)
(set of sets)

examples of infinite classes: \mathcal{Q} or \mathcal{A} (collection of sets)

1. $A = \{[-\infty, t] : t \in \mathbb{R}\}$ - set of all half intervals on real-line (\mathcal{QF})

W: What is sample space?

- Let F be an arbitrary finite set $F = \{x_1, \dots, x_n\}$

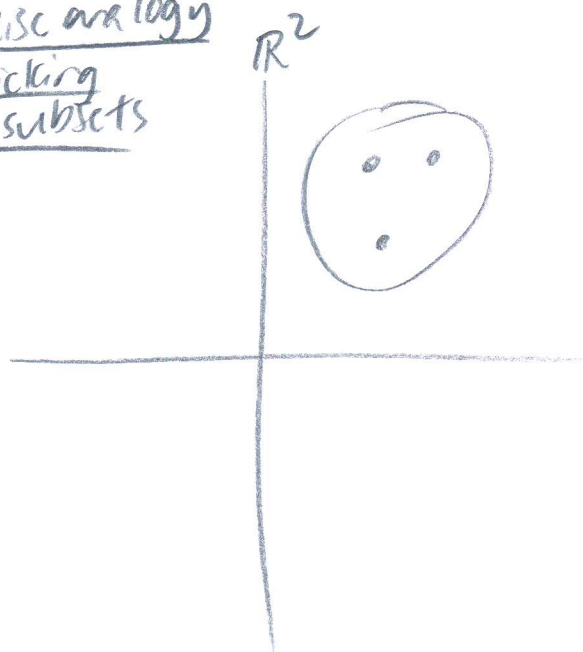
- W: So we have an arbitrary finite set F ; but also a class of sets \mathcal{Q} (*)

- W: Take any set from the \mathcal{Q} and intersect it with our finite set.

$$A \cap F = G \text{ and } G \subset F$$

- use disc analogy

for picking out subsets



(*) we say A picks out G if $A \cap F = G$ for some $A \in \mathcal{A}$
 (collection) (subset)

(*) $S(A, F)$ - no. of subsets that are picked out by A .

simple example

$$A = \{[-\infty, t] : t \in \mathbb{R}\}$$



- finite set $F = [x_1, x_2]$

- Q How many possible sets can be extracted from F ?

- possible subsets of F (i.e. the G s) - $\emptyset, x_1, x_2, [x_1, x_2]$

- can we get these by considering infinite coll of sets A ?

- $S_n(A)$ is measure of complexity of class of sets A

W: This class of sets A is not as rich as we might think as $S_n(A) \neq 2^n$ for many n .

- more complex sets $\rightarrow S_n(A) \leq 2^n$ for many n .

- Vapnik-Chervonenkis (not proved) Theorem 5

$$P\left(\sup_{A \in A} |P_n(A) - P(A)| > \epsilon\right) \leq 8 S_n(A) e^{-\frac{n\epsilon^2}{32}}$$

- "Extension of Hoeffding:

Probability that maximum of absolute diff. between empirical prob. and the probability over entire class of sets A is greater than ϵ is bounded by a factor made up of an exponential term in n, ϵ and the shattering coefficient $S_n(A)$ "

- $S_n(A)$ - shattering coefficient

W: Bound is very tight over a complex class of sets A (?)

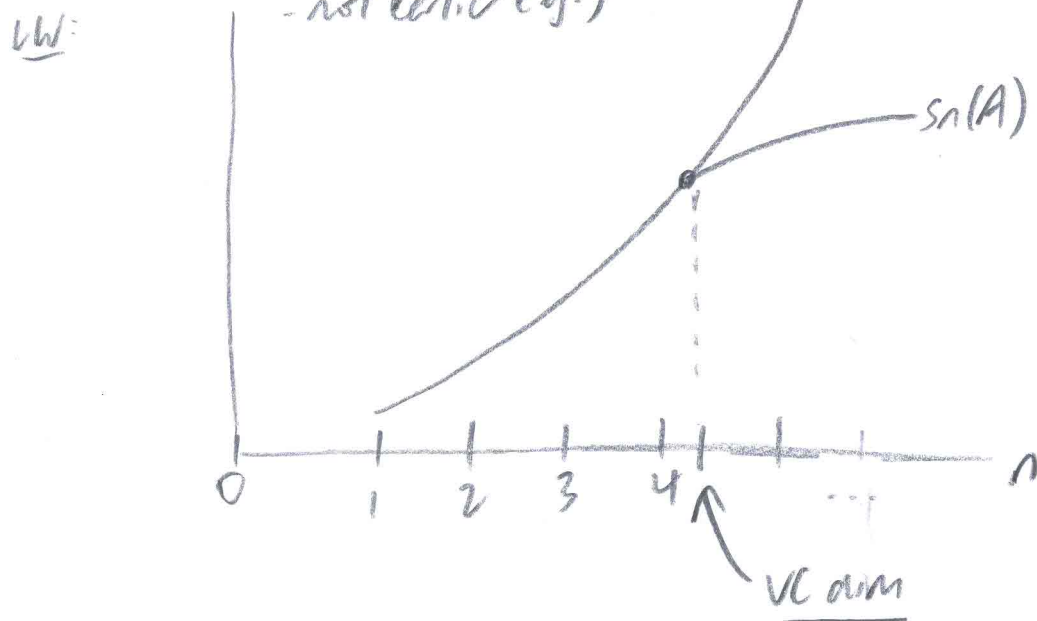
- $S_n(A)$ is potentially

- complex set class of sets $S_n(A) \rightarrow$ poor bound (?)

- Consider rates of increase of $S_n(A)$ and $e^{-\frac{n\epsilon^2}{32}}$ separately in n ;

(Q): How do they offset each other?

(General case
- not eerie e.g.)



- Assume we can plot/have algorithm for $S_n(A)$ as fn of n .

- Plot $f=(2^n)$ $g=(S_n(A))$

- VC dimension

- VC dimension:-

def $d = d(A) = \text{largest } n: s_n(A) = 2^n$

- d is the size of the largest set that can be shattered (ie. $s(A, \mathcal{F}) = 2^d$)

W: Sauer's lemma (*) As $s_n(A)$ increases ⁿ before $n = \text{VC dimension} = d$;
rate of increase is exponential

As $s_n(A)$ increases in n after $n = \text{VC dim} = d$; rate of increase is polynomial.

Theorem 7: Suppose A has a finite VC dimension d . Then for all $n \geq d$:-

$$s(A, n) \leq (n+1)^d$$

(*) At this particular
polynomial rate

- After VC dimension d :- ie. $n \geq d$; $n \rightarrow \infty$

$$8 s_n(A) e^{-\frac{ne^2}{32}}$$

decreases exponentially

increases
at polynomial rate

- i.e. an extension of Hoeffding; with some add. info

W: How to compute VC dimension? (*) VC Theory & Sauer proved
m 36-702 / 10-702

- Sometimes easy / trivial

Assume we have a VC dimension dictionary (for course purposes)

WAS - Try for \mathbb{R} - VC dim of 2. (see Table 1) (?)

- W(*):- Know VC dim \rightarrow know shattering $\rightarrow s_n(A)$...

(*) - simple probability bounds \rightarrow Hoeffding

collections
of sets \rightarrow use VC Theory/inequality
(extension of Hoeffding).

YouTube 12/09/2016

Convergence theory

• LW: answers question of what happens when we get more data (asymptotics)

- 2 key results - LLN, CLT

- deal with sequence of r.v.s. $X_1, \dots, X_n \sim F$

- concerned with statistics $T_n = g(X_1, \dots, X_n)$ - a function of data of sequence of r.v.s.

- e.g. sample mean $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ - (*) sequence of sample means (as r.v.s. are NOT iid) ✓

(*) A sequence of statistics may not be iid, whilst underlying r.v.s. are iid

more abstractly

original sequence Y_1, Y_2, Y_3 (likely iid)

$$X_n = g_n(Y_1, \dots, Y_n)$$

- compute sth from first n data points

- new (~~arbitrary~~) sequence of statistics; can be viewed as arbitrary seq of r.v.s. :-

(use/have $X_i = \bar{X}_n$ for concrete example).

$X_1, X_2, X_3, X_4, \dots$ where $X_i = g(Y_1, \dots, Y_i)$

Q: what does it mean for sequence X_1, X_2, \dots to converge??

- more complex than calculus for r.v.s.

- find many different types of convergence.

LW: will be confusing → will take a few lectures

At
Almost sure convergence (to a constant)

$X_n \xrightarrow{\text{a.s.}} c$

$P\left(\lim_{n \rightarrow \infty} X_n = c\right) = 1$
 conv. calc. definite

probabilistic conv. defn.

Convergence in probab. (to an r.v. / constant) (*) essential

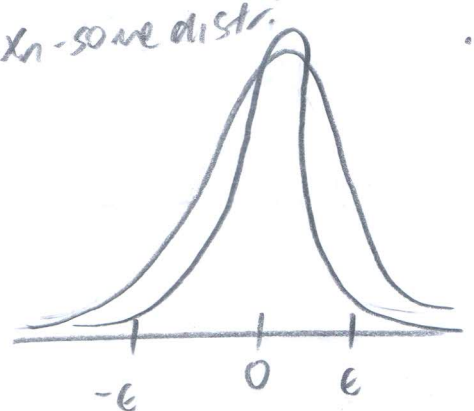
$X_n \xrightarrow{P} X$

$X_n \xrightarrow{P} c \quad (*) \text{essential}$

$\forall \epsilon > 0 \quad P(|X_n - X| > \epsilon) \rightarrow 0 \quad \forall \epsilon > 0 \quad P(|X_n - c| > \epsilon) \rightarrow 0$

As $n \rightarrow \infty$

W: X_n - some distr.



As n gets large, distr gets squished towards 0

Already seen:-

$X_n - c = o_p(1)$

W: convergence in probability
 not almost sure convergence not the same

Convergence in quadratic mean \rightarrow E^2 or $E[\]^2$? $\checkmark E[\]^2$

$X_n \xrightarrow{q.m.} X \quad E[(X_n - X)^2] \rightarrow 0$

- convergence in L_2

$X_n \xrightarrow{q.m.} c \quad E[(X_n - c)^2] \rightarrow 0$

Convergence in distri (*) - essential

$X_n \rightsquigarrow X$ also $X_n \xrightarrow{d} X$

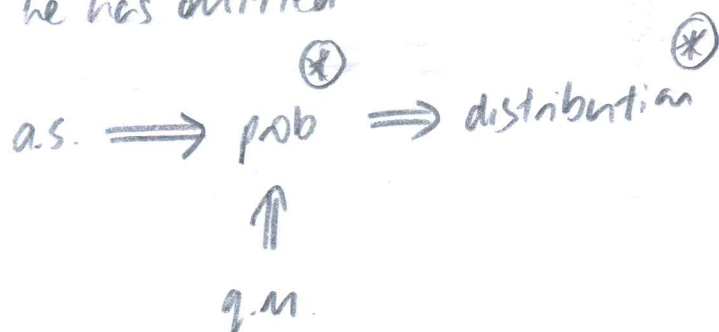
W: CDF of $X_n \rightarrow$ CDF of X

$\lim_{n \rightarrow \infty} F_n(t) = F(t)$

$F_n(t) \rightarrow F(t)$ as $n \rightarrow \infty$; at all t for which F is continuous

(*) - show by example.

(W2): Review and polish/check other definitions he has omitted



} A taxonomy of convergence relationships

- W: The \circledast are what we are concerned with
- a.s. \rightarrow subtle (not used for our purposes)
 - q.m. \rightarrow stepping stone to proving prob.

W: make clear through proofs, examples

- All above are general (no ass. about discreteness/continuous etc.)

W: Theorem 3: for curiosity

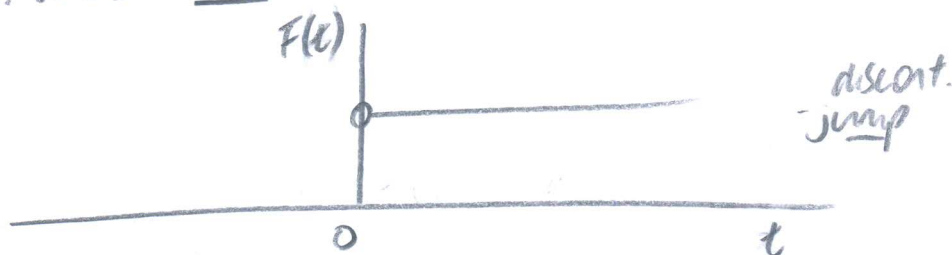
Example 4

- Consider following seq. of r.v.s. with $X_n \sim N(0, \frac{1}{n})$

X_1, X_2, \dots

- Intuitively; concentrates around 0 as $n \rightarrow \infty$

- Consider r.v. whose value is always 0.



CDF $F(t) = \begin{cases} 0 & t < 0 \\ 1 & t \geq 0 \end{cases}$

• Note $\sqrt{n}X_n \sim N(0, 1)$

- Showing $X_n \rightarrow 0$ \Rightarrow show $F_n(t) \rightarrow F(t)$ as $n \rightarrow \infty$; note

$F_n(t) = P(X_n \leq t) = P(\sqrt{n}X_n \leq \sqrt{n}t) = P(Z \leq \sqrt{n}t)$

consider two cases

$$t > 0 \quad P(Z \leq \sqrt{n}t) \rightarrow 1 = F(t) \\ \text{as } n \rightarrow \infty$$

$$t < 0 \quad P(Z \leq \sqrt{n}t) \rightarrow 0 = F(t) \\ \text{as } n \rightarrow \infty$$

Q: How about at 0?
i.e. $t=0$

(?) (W) (A4)

$$F_n(0) = P(X_n \leq 0) \quad \checkmark \cdot F$$

$$= \frac{1}{2}$$

$$F_n(0) \neq F\left(\frac{1}{2}\right) = 1$$

• This does not converge

• convergence in dist only requires convergence
at continuity points

• exclude discontinuity

Convergence fails

• We have proved :-

$$X_n \xrightarrow{d} X$$

where X is always equal to 0 i.e. $P(X=0)=1$

$$X_n \xrightarrow{d} 0$$

in next lecture, more examples

formal proofs to develop intuition