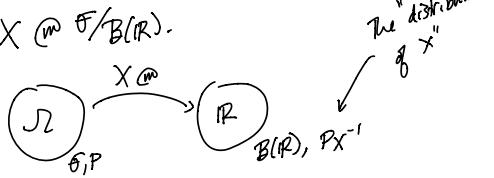


Lecture 11: Random variable densities and expected values

Recall that random variable is a map $X: \Omega \rightarrow \mathbb{R}$ along with a probability measure (Ω, \mathcal{F}, P) s.t. $X \in \mathcal{F}/B(\mathbb{R})$.



Def: The expected value of X , denoted $E(X)$, is defined as

$$E(X) = \int_{\Omega} X dP = \int_{\Omega} X^+ dP - \int_{\Omega} X^- dP$$

when it is defined, i.e. when $X \in \mathcal{F}/B(\Omega, \mathcal{F}, P)$.

You should think of X as a placeholder for a random number $X(w)$ obtained by choosing $w \in \Omega$ at random according to P . Then $E(X)$ is essentially what you "expect" X to be.

e.g. For $\theta \in [0, 1]$ the r.v. X has a Bernoulli θ distribution, denoted

$X \sim \text{Ber}(\theta)$, if

$$P(X=1) = \theta$$

$$P(X=0) = 1 - \theta.$$

Let $A = \{w : X(w) = 1\}$ so that

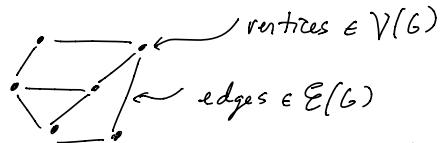
$$X(w) = I_A(w) \quad P\text{-a.e.}$$

$$\begin{aligned} \therefore E(X) &= E(I_A(w)) \quad \text{by "a.e. useful"} \\ &= P(A) \quad \text{by def.} \\ &= \theta \quad \text{since } A = \{X=1\}. \end{aligned}$$

Note: It is always the case that (2)

$$E(I_A) = P(A) \quad \text{whenever } A \in \mathcal{F}.$$

e.g. This example uses expected value & probability to prove a "non-probabilistic" statement about graphs G :



These types of proofs were made famous by Erdős.

Claim: Every graph G has a bipartite subgraph H for which $\#E(H) \geq \frac{1}{2}\#E(G)$.

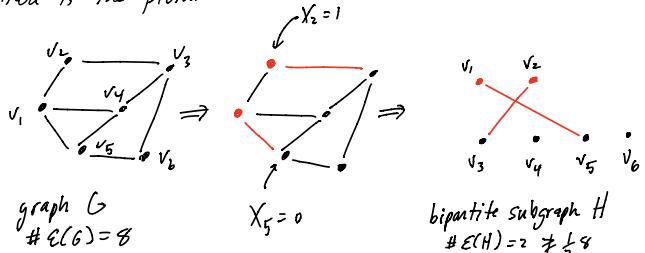
Proof:

Suppose G has n vertices labeled v_1, v_2, \dots, v_n . Let X_1, X_2, \dots, X_n be n independent $\text{Ber}(\frac{1}{2})$ r.v.s defined on some probability space (Ω, \mathcal{F}, P) . Define the subgraph H as follows

$$V(H) := V(G)$$

$$E(H) := \left\{ v_i v_j \in E(G) : (X_i, X_j) = (1, 0) \text{ or } (X_i, X_j) = (0, 1) \right\}$$

Here is the picture



Notice that $\#E(H)$ is a r.v. in $\mathcal{Y}_S(\Omega, \mathcal{F})$.

Let $\mathcal{Z} = \{(i, j) : i > j \text{ & } v_i v_j \in E(G)\}$ index all edges $E(G)$ so that

Now

$$\begin{aligned}
 E(\#\mathcal{E}(H)) &= E\left(\sum_{(i,j) \in \mathcal{X}} I_{\{(X_i, X_j) = (1,0)\}} + I_{\{(X_i, X_j) = (0,1)\}}\right) \quad (3) \\
 &\stackrel{\text{B.g. 3}}{=} \sum_{(i,j) \in \mathcal{X}} P(X_i=1, X_j=0) + P(X_i=0, X_j=1) \\
 &= \sum_{(i,j) \in \mathcal{X}} \left(\frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2}\right) \\
 &= \frac{1}{2} \underbrace{\sum_{(i,j) \in \mathcal{X}} 1}_{\#\mathcal{E}(G)} \quad (*)
 \end{aligned}$$

Now if $\#\mathcal{E}(\tilde{H}) < \frac{1}{2}\#\mathcal{E}(G)$ for all bipartite subgraphs \tilde{H} of G then

$$E(\#\mathcal{E}(H)) < \frac{1}{2}\#\mathcal{E}(G)$$

\curvearrowleft easy to see since
 $\#\mathcal{E}(H)$ is a simpler r.v.

This contradicts (*) so that we must have
 \exists a bipartite subgraph \tilde{H} s.t.

$$\#\mathcal{E}(\tilde{H}) \geq \frac{1}{2}\#\mathcal{E}(G). \quad \underline{\text{QED}}$$

Let's look at a couple fundamental properties of expected value.

Theorem (Jensen's inequality)

If $\varphi: \mathbb{R} \rightarrow \mathbb{R}$ is a convex function and $X \in L_1(\Omega, \mathcal{F}, P)$ then $\varphi(X)$ is quasi-integrable

and $\varphi(E(X)) \leq E(\varphi(X)).$

Proof:

Let $x \mapsto ax+b$ be a supporting line of φ passing through the point $(E(X), \varphi(E(X)))$

\curvearrowleft finite.

(3)

In particular let $a, b \in \mathbb{R}$ satisfy

$$\begin{cases} ax+b \leq \varphi(x) \text{ for } x \in \mathbb{R} \\ aE(X)+b = \varphi(E(X)) \end{cases}$$

$$\therefore aX+b \leq \varphi(X) \quad (*)$$

Notice that $aX+b$ is integrable since X is integrable.

Also φ is convex & mapping into \mathbb{R} $\Rightarrow \varphi$ is continuous

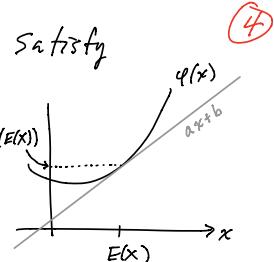
$$\Rightarrow \varphi \text{ is C}^1$$

$\Rightarrow \varphi(X)$ is a r.v.
and in L^1 by (*).

$$\therefore \underbrace{aE(X)+b}_{\parallel} \leq E(\varphi(X)) \quad \text{by B.g. 3}$$

$$\varphi(E(X))$$

QED



e.g. Let's use Jensen's inequality to get an understanding of why type of maximal fluctuation we might expect to see if we stare at a large number of random variables

Theorem (What to expect of the max)

Let X_1, X_2, \dots, X_n be r.v.s in $L_1(\Omega, \mathcal{F}, P)$.

Suppose $\exists \sigma > 0$ s.t.

$$E(e^{tx_i}) \leq \exp\left(\frac{t^2\sigma^2}{2}\right), \quad \forall t > 0, \forall i \in n$$

$$\text{then } E\left(\max_{1 \leq i \leq n} X_i\right) \leq \sigma \sqrt{2 \log n}.$$

Proof:

$$\begin{aligned}
 \exp(t E(\max_{i \leq n} X_i)) &\stackrel{\text{Jensen}}{\leq} E(\exp(t \max_{i \leq n} X_i)) \\
 &= E\left(\max_{i \leq n} \exp(t X_i)\right) \\
 &\quad \text{since } e^{tx} \uparrow \text{ in } x \\
 &\leq E\left(\sum_{i \leq n} \exp(t X_i)\right) \\
 &\quad \text{since these are } \geq 0 \\
 &\stackrel{\text{Big 3}}{=} \sum_{i \leq n} E(\exp(t X_i)) \\
 &\leq n \exp\left(\frac{t^2 \sigma^2}{2}\right) \text{ by assumption.}
 \end{aligned}$$

Taking log gives

$$E(\max_{i \leq n} X_i) \leq \underbrace{\frac{\log n}{t} + \frac{t \sigma^2}{2}}_{\text{Now choose a good } t} \quad \forall t > 0.$$

Notice

$$\begin{aligned}
 \frac{d}{dt} \left(\frac{\log n}{t} + \frac{t \sigma^2}{2} \right) &= -\frac{\log n}{t^2} + \frac{\sigma^2}{2} = 0 \\
 \Updownarrow \\
 t &= \frac{\sqrt{2 \log n}}{\sigma}
 \end{aligned}$$

Plugging this into our inequality gives

$$\begin{aligned}
 E(\max_{i \leq n} X_i) &\leq \frac{\sigma}{\sqrt{2 \log n}} \log n + \frac{\sqrt{2 \log n}}{\sigma} \frac{\sigma^2}{2} \\
 &= \frac{\sigma}{\sqrt{2}} \sqrt{\log n} + \frac{\sigma}{\sqrt{2}} \sqrt{\log n} \\
 &= \sigma \sqrt{2 \log n} \quad \underline{\text{QED}}
 \end{aligned}$$

(5)

Theorem (expected value factors on indep r.v.s)

Suppose X and Y are (possibly extended) independent r.v.s on (Ω, \mathcal{F}, P) . If $X \geq 0$ & $Y \geq 0$ or $X, Y \in L_1(\Omega, \mathcal{F}, P)$ then $XY \in Q(\Omega, \mathcal{F}, P)$ and $E(XY) = E(X)E(Y)$.

(6)

Proof:

Case 1: Suppose $X, Y \in \mathcal{H}_S(\Omega, \mathcal{F}, P)$ so that

$$X = \sum_{i=1}^n a_i I_{A_i} \quad \& \quad Y = \sum_{j=1}^m b_j I_{B_j}$$

where a_1, \dots, a_n are distinct and A_1, \dots, A_n are a disjoint measurable partition of Ω (and same for b_j 's & B_j 's).

Note that A_i is indep of B_j since

$$\begin{aligned}
 A_i &= \{X = a_i\} \in \sigma\langle X \rangle \quad \xrightarrow{\text{indep } \sigma\text{-fields}} \\
 B_j &= \{Y = b_j\} \in \sigma\langle Y \rangle
 \end{aligned}$$

$$\begin{aligned}
 \therefore E(XY) &= E\left(\sum_{i,j} a_i b_j I_{A_i \cap B_j}\right) \\
 &= \sum_{i,j} a_i b_j \underbrace{P(A_i \cap B_j)}_{P(A_i)P(B_j)} \\
 &\quad \text{Note that } a_i \text{ or } b_j \text{ could be } \infty \text{ but } \mathcal{H}_S \text{ allows } I_{A_i \cap B_j} \text{ to apply} \\
 &= \left(\sum_i a_i P(A_i)\right) \left(\sum_j b_j P(B_j)\right) \\
 &= E(X)E(Y).
 \end{aligned}$$

Case 2: Suppose $X, Y \in \mathcal{H}(\Omega, \mathcal{F})$.

Notice that we also have that

$$X \in \mathcal{H}(\Omega, \sigma\langle X \rangle) \quad \& \quad Y \in \mathcal{H}(\Omega, \sigma\langle Y \rangle).$$

Therefore the Structure Theorem
implies there exists $X_n \in \mathcal{N}(\mathcal{A}, \sigma\langle X \rangle)$
and $Y_n \in \mathcal{N}(\mathcal{A}, \sigma\langle Y \rangle)$ such that

$$X_n \uparrow X \text{ and } Y_n \uparrow Y.$$

$$\therefore E(X_n Y_n) = E(X_n) E(Y_n) \text{ by case 1.}$$

$$\begin{aligned}\therefore E(XY) &= E(\lim_n X_n Y_n) \\ &= \lim_n E(X_n Y_n) \text{ by little 3} \\ &= \lim_n E(X_n) E(Y_n) \\ &= E(X) E(Y).\end{aligned}$$

Case 3: Suppose $X, Y \in L_1(\mathcal{A}, \mathcal{F}, P)$.

Notice that

$$\begin{aligned}(XY)^+ &= X^+ Y^+ + X^- Y^- \\ (XY)^- &= X^+ Y^- + X^- Y^+\end{aligned}$$

Therefore

$$E(XY)^+ = E(X^+) E(Y^+) + E(X^-) E(Y^-) < \infty$$

$$E(XY)^- = E(X^+) E(Y^-) + E(X^-) E(Y^+) < \infty$$

by little 3, case 2 and the fact that
 $\sigma\langle X^+ \rangle \subset \sigma\langle X \rangle, \dots \text{ & } \sigma\langle Y^- \rangle \subset \sigma\langle Y \rangle$.

(Notice I'm implicitly using "X @ w.r.t $\sigma\langle Y \rangle$ thm" & "subclasses".)

$\therefore XY \in L_1(\mathcal{A}, \mathcal{F}, P)$ and

$$\begin{aligned}E(XY) &= E(XY)^+ - E(XY)^- \\ &= E(X^+) E(Y^+) + E(X^-) E(Y^-) \\ &\quad - E(X^+) E(Y^-) - E(X^-) E(Y^+) \\ &= E(X^+) E(Y) - E(X^-) E(Y) \\ &= E(XY).\end{aligned}$$

QED

⑦

Notice this fully generalizes to situations like this: Suppose X_1, X_2, \dots are independent $L_1(\mathcal{A}, \mathcal{F}, P)$ r.v.s then

$$E(X_1 X_2 \cdots X_n) = E(X_1) E(X_2) \cdots E(X_n)$$

and

$$\begin{aligned}E(g(X_1, X_2, \dots) h(X_2, X_4, \dots)) \\ = E g(X_1, X_3, \dots) E h(X_2, X_4, \dots)\end{aligned}$$

if $g, h \in \mathcal{N}(\mathbb{R}^\infty, \mathcal{B}(\mathbb{R}^\infty))$. To show these just use the previous Thm, ANOVA & "X @ w.r.t $\sigma\langle Y \rangle$ Thm".

Def: If X and Y are two r.v.s on $(\mathcal{A}, \mathcal{F}, P)$

then set

$$\text{var}(X) = \text{"the variance of } X\text{"} := E(X - E(X))^2$$

$$\text{sd}(X) = \text{"the standard deviation of } X\text{"} := \sqrt{\text{var}(X)}$$

$$\text{cov}(X, Y) = \text{"the covariance b/w } X \text{ & } Y\text{"}$$

$$:= E[(X - E(X))(Y - E(Y))]$$

when they are defined.

Note: $\text{var}(X)$, $\text{sd}(X)$ & $\text{cov}(X, Y)$ may not be defined if the expectations which define them do not exist.

Remark: $\text{sd}(X)$ measures how spread out X is on \mathbb{R} & $\text{cov}(X, Y)$ measure how X & Y co-vary together.

(9)

Later we will see that $sd(X)$ and $cov(X, Y)$ are essentially the functional analysis equivalent to L_2 norm & L_2 inner product. Here is a hint as to why.

Theorem (Hölder)

Let X and Y be two R.V.s on (Ω, \mathcal{F}, P) . If $p, q > 0$ s.t. $\frac{1}{p} + \frac{1}{q} = 1$ then

$$E|XY| \leq (E|X|^p)^{\frac{1}{p}} (E|Y|^q)^{\frac{1}{q}} \quad (*)$$

even if X and Y are not quasi-integrable.

If $XY \in L_1(\Omega, \mathcal{F}, P)$ then

$$|E(XY)| \leq (E|X|^p)^{\frac{1}{p}} (E|Y|^q)^{\frac{1}{q}}. \quad (**)$$

Proof:

$(*)$ is trivially true if any one of the following is true:

$E|X|^p = \infty$, $E|X|^p = 0$, $E|X|^q = \infty$, $E|X|^q = 0$.

Since "a.e. useful facts"

So suppose $E|X|^p, E|Y|^q \in (0, \infty)$. implies $|X|^p \stackrel{\text{a.e.}}{=} 0$ so that $|XY| \stackrel{\text{a.e.}}{=} 0$.

Define $Z := \frac{X}{(E|X|^p)^{\frac{1}{p}}} \& W := \frac{Y}{(E|Y|^q)^{\frac{1}{q}}}$

Now we simply show

$$E|ZW| \leq 1.$$

Use Young's inequality

$$a^{w_1} b^{w_2} \leq w_1 a + w_2 b \quad (***)$$

when $a, b \geq 0$ & $w_1, w_2 > 0$ s.t. $w_1 + w_2 = 1$

Young's inequality follows since \log is concave so that

$$w_1 \log a + w_2 \log b \leq \log(w_1 a + w_2 b).$$

(10)

Now $E|ZW| = E\left(\underbrace{(|Z|^p)^{\frac{1}{p}}}_{\text{has the form } (***)} (\underbrace{|W|^q)^{\frac{1}{q}}}_{\text{when }}\right)$

$a = (|Z|^p)^{\frac{1}{p}}$, $b = (|W|^q)^{\frac{1}{q}}$,

$w_1 = \frac{1}{p}$ & $w_2 = \frac{1}{q}$.

$\stackrel{(***)}{\leq} \underbrace{\frac{1}{p} E|Z|^p}_{=1} + \underbrace{\frac{1}{q} E|W|^q}_{=1} \text{ using Big 3 (2) [1]}$

$= 1.$

If $XY \in L_1(\Omega, \mathcal{F}, P)$ then

$$|E(XY)| \leq E|XY|$$

by corollary to Big 3.

QED.

Corollary:

If X and Y are two r.v.s on (Ω, \mathcal{F}, P) s.t. $E(X^2) < \infty$ & $E(Y^2) < \infty$ then $cov(X, Y)$, $sd(X)$ & $sd(Y)$ are well defined and

$$|cov(X, Y)| \leq sd(X) sd(Y)$$

Proof: By Hölder $E|X| \leq \sqrt{E|X|^2} < \infty$ so that $X, Y \in L_1(\Omega, \mathcal{F}, P)$ and $E(X) < \infty$, $E(Y) < \infty$. Let $\tilde{X} = X - E(X)$, $\tilde{Y} = Y - E(Y)$. Also by Hölder we have $E|\tilde{X}\tilde{Y}| \leq \sqrt{E(\tilde{X}^2)} \sqrt{E(\tilde{Y}^2)} < \infty$

$$\therefore \tilde{X}\tilde{Y} \in L_1 \& |cov(X, Y)| \leq \sqrt{E(\tilde{X}^2)} \sqrt{E(\tilde{Y}^2)}$$

$$= sd(X) sd(Y) \quad \underline{QED}$$

Corollary:

If X and Y are two independent r.v.s s.t. $E(X^2) < \infty$ & $E(Y^2) < \infty$ then $\text{cov}(X, Y)$, $E(X)$ and $E(Y)$ are well defined, finite and

$$\begin{aligned}\text{cov}(X, Y) &= E[(X - EX)(Y - EY)] \\ &= \underbrace{E(X - EX)}_{=0} \underbrace{E(Y - EY)}_{=0}\end{aligned}$$

Once we prove the following theorem we will discuss how it gives an interesting relation b/w Lebesgue integration and Riemann integration.

Theorem:

Suppose X is a r.v. on (Ω, \mathcal{F}, P) s.t.

$X \geq 0$ p.a.e. Then

$$\begin{aligned}E(X) &= \int_0^\infty P(X > t) dt \quad \leftarrow \\ &\stackrel{(**)}{=} \int_0^\infty P(X \geq t) dt \quad \leftarrow \text{Lebesgue integral}\end{aligned}$$

Proof:

First notice that $t \mapsto P(X \geq t)$ & $t \mapsto P(X > t)$ are both monotonic and therefore $\mathcal{B}(\mathbb{R})/\mathcal{B}(\mathbb{R})$.

Clearly these two functions are in $\mathcal{C}_b^-([0, \infty], \mathcal{B}([0, \infty]), \mathbb{Z}')$ so the integrals in

(*) and (**) are well defined.

To show (*) use the "1-2 argument".

(1)

Step 1: Show (*) holds for $X \in \mathcal{N}_s(\Omega, \mathcal{F})$ (12)

$\therefore X$ can be written in the form

$$X = \sum_{i=1}^n c_i I_{A_i} \quad \leftarrow \text{measurable partition}$$

$$\text{so that } E(X) = \sum_{i=1}^n c_i P(A_i).$$

Now

$$\int_0^\infty P(X > t) dt = \int_0^\infty \sum_{i=1}^n P(\{X > t\} \cap A_i) dt$$

since A_i 's partition Δ

$$\begin{aligned}&\stackrel{\text{Big 3}}{=} \sum_{i=1}^n \int_0^\infty P(\{X > t\} \cap A_i) dt \\ &\quad \leftarrow \text{on } A_i, X = c_i \\ &= \sum_{i=1}^n \int_0^\infty \underbrace{P(\{c_i > t\} \cap A_i)}_{\begin{cases} 0 & \text{if } c_i \leq t \\ P(A_i) & \text{o.w.} \end{cases}} dt \\ &= \sum_{i=1}^n c_i P(A_i) \\ &= E(X)\end{aligned}$$

Step 2:

If $X \in \mathcal{N}(\Omega, \mathcal{F})$ then $\exists X_n \in \mathcal{N}_s(\Omega, \mathcal{F})$ s.t.

$$X_n \uparrow X.$$

$$\begin{aligned}\therefore E(X) &= E(\lim_n \uparrow X_n) \\ &= \lim_n \uparrow E(X_n), \text{ little 3} \\ &= \lim_n \uparrow \int_0^\infty P(X_n > t) dt \\ &= \int_0^\infty \lim_n \uparrow P(X_n > t) dt, \text{ little 3} \\ &= \int_0^\infty P(X > t) dt \\ &\quad \text{by CFB since } \{X_n > t\} \uparrow \{X > t\} \\ &\quad (\text{but not } \{X_n \geq t\} \uparrow \{X \geq t\})\end{aligned}$$

This gives (*).

To show (**) simply notice that (13)

$$P(X \geq t) = P(X > t) + \underbrace{P(X=t)}$$

This can only be non-zero
for at most countably
many t 's.

$$= P(X > t) \quad \mathcal{L}^1\text{-a.e.}$$

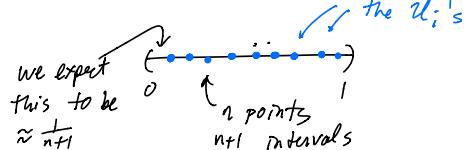
QED

The previous theorem can be useful for computing the expected value of the minimum of a sequence of r.v.

e.g. Let U_1, U_2, \dots, U_n be iid $\text{Unif}(0,1)$ random variables. Then

$$\begin{aligned} E(\min_{1 \leq i \leq n} U_i) &= \int_0^\infty P\left(\min_{1 \leq i \leq n} U_i \geq t\right) dt \\ &\quad \text{This event holds iff each } U_i \geq t. \\ &= \int_0^\infty P(U_1 \geq t, \dots, U_n \geq t) dt \\ &= \int_0^\infty P(U_1 \geq t)^n dt \\ &\quad \text{since } U_i \text{'s are indep and identically distributed} \\ &= \int_0^1 (1-t)^n dt \\ &= -\frac{(1-t)^{n+1}}{n+1} \Big|_0^1 = \frac{1}{n+1} \end{aligned}$$

In some sense we get the following picture



(14) Using improper Riemann integration to define $E(X)$.

To relate the formula $E(X) = \int_0^\infty P(X > t) dt$ to Riemann integration first notice

Thm (CDF's have countably many jumps)

If X is a.r.v. with cdf $F(t) := P(X \leq t)$ then $\{t \in \mathbb{R} : F \text{ is discontinuous at } t\}$ is countable

Proof:

$$\begin{aligned} \text{Since } F(t) \text{ is right continuous} \\ F \text{ is discontinuous at } t &\iff F(t) \neq F(t-) \\ &\iff F(t) - F(t-) > 0 \\ &\iff P(X=t) > 0 \end{aligned}$$

But $\{\{X=t\} : t \in \mathbb{R}\}$ is a countable collection of disjoint events $\implies P(X=t) > 0$ for at most countably many $t \in \mathbb{R}$ (by a Thm in Lecture 5).

QED

This means that $t \mapsto P(X > t) = P(X \leq t)$ is Riemann integrable on any bdd interval $t \in [a, b]$ (since it is bdd and has countably many discontinuities)

Now when $X \geq 0$

$$E(X) = \int_0^\infty P(X > t) dt$$

abstract integral
 $\int \chi_{(t, \infty)} dP(t)$

This can be interpreted as a
improper Riemann integral on $[0, \infty]$.
... this can be used to define $E(X)$
with just improper Riemann integration!

Probability inequalities involving expected value (15)

It is often the case that computing bounds for expected value is easy. These inequalities then yield probability bounds

Theorem (Markov's inequality)

If X is a r.v. on (Ω, \mathcal{F}, P) s.t. $X \geq 0$ P-a.e. then $P(X \geq \alpha) \leq \frac{E(X)}{\alpha}$ for all $\alpha > 0$.

Proof: $X \geq 0$ P-a.e. implies that $X \in \Omega^c(\Omega, \mathcal{F}, P)$ and $X I_{\{X \geq \alpha\}} \geq \alpha I_{\{X \geq \alpha\}}$ "Indicate what you want trick"

$$\begin{aligned} \therefore E(X) &\geq E(X I_{\{X \geq \alpha\}}) \quad \text{by Prop 3(i)} \text{ and } X \geq 0 \text{ P-a.e.} \\ &\geq E(\alpha I_{\{X \geq \alpha\}}) \quad \text{by Prop 3(i)} \\ &= \alpha P(X \geq \alpha) \quad \text{definition of } \int_X dP \end{aligned}$$

QED.

Corollary:

- For any r.v. X on (Ω, \mathcal{F}, P) and any $\alpha > 0$:
- $P(|X| \geq \alpha) \leq \frac{E(|X|^k)}{\alpha^k}$ Chernoff's method
 - $P(X \geq \alpha) \leq \inf_{t > 0} \frac{E(e^{tX})}{e^{t\alpha}}$
 - If $E(X^2) < \infty$ then $E(X) < \infty$ and Chebyshev's inequality
- $$P(|X - EX| \geq \alpha) \leq \frac{\text{var}(X)}{\alpha^2}$$

Note: i) says that if $E|X|^k < \infty$ for a large $k > 0$ then $P(|X| \geq \alpha)$ decays quickly as $\alpha \rightarrow \infty$. We used ii) in Lecture 1 when illustrating the "modern" way to prove Borel's Normal number theorem. iii) shows why $\text{sd}(X)$ controls the "spread" of X .

Let's use Chebyshev to show a famous theorem in analysis.

Lagrange

Theorem (Weierstrass Approximation)

If $f: [0,1] \rightarrow \mathbb{R}$ is continuous then for any $\epsilon > 0$ there exists a polynomial $p(x)$ s.t.

$$\sup_{x \in [0,1]} |f(x) - p(x)| < \epsilon.$$

Proof: Let U_1, U_2, \dots be iid r.v.s uniformly distributed on $[0,1]$. For each $x \in [0,1]$ let $F_x(\cdot)$ be the C.D.F. of a $\text{Ber}(x)$ r.v. and set

$$S_n^x := F_x^{-1}(U_1) + \dots + F_x^{-1}(U_n)$$

$\nwarrow \uparrow$
independent $\text{Ber}(x)$ r.v.s

Note that S_n^x is a collection of coupled r.v.s induced by x

A simple counting argument shows

$$S_n^x \sim \text{Bin}(n, x)$$

$$P(S_n^x = m) = \binom{n}{m} x^m (1-x)^{n-m}$$

for $m = 0, 1, \dots, n$. Moreover,

$$E(S_n^x) = x \quad \text{and} \quad \text{var}(S_n^x) = \frac{x(1-x)}{n}.$$

Also notice that

$$P_n(x) := E f\left(\frac{S_n^x}{n}\right) = \underbrace{\sum_{m=0}^n f\left(\frac{m}{n}\right)}_{\text{simple r.v.}} \underbrace{P(S_n^x = m)}_{\text{Polynomial in } x \text{ of degree } n.}$$

Since f is uniformly continuous on $[0,1]$ let

$$M := \sup_{x \in [0,1]} |f(x)| < \infty$$

$$S(\epsilon) := \sup \{|f(x) - f(y)| : |x - y| < \epsilon\}$$

and for a given $\epsilon > 0$ and $x \in [0,1]$ let

$$A_{\epsilon,x} := \left\{ w \in \mathbb{N} : \left| \frac{S_n^x(w)}{n} - x \right| < \epsilon \right\}$$

The definition of $A_{\varepsilon,x}$ now implies (17)

$$w \in A_{\varepsilon,x} \Rightarrow |f(\frac{S_n^x(n)}{n}) - f(x)| \leq \delta(\varepsilon)$$

$$w \in A_{\varepsilon,x}^c \Rightarrow |f(\frac{S_n^x(n)}{n}) - f(x)| \leq 2M$$

so that

$$\begin{aligned} |f(\frac{S_n^x(n)}{n}) - f(x)| &= |f(\frac{S_n^x(n)}{n}) - f(x)| (I_{A_{\varepsilon,x}(n)} + I_{A_{\varepsilon,x}^c(n)}) \\ &\leq \delta(\varepsilon) I_{A_{\varepsilon,x}(n)} + 2M I_{A_{\varepsilon,x}^c(n)} \end{aligned}$$

To finish

$$\begin{aligned} |P_n(x) - f(x)| &= |E f(\frac{S_n^x}{n}) - f(x)| \\ &\leq E |f(\frac{S_n^x}{n}) - f(x)| \text{ by Jensen} \\ &\leq \delta(\varepsilon) P(|\frac{S_n^x}{n} - x| < \varepsilon) + 2M P(|\frac{S_n^x}{n} - x| \geq \varepsilon) \\ &\leq \delta(\varepsilon) + 2M \frac{\text{Var}(\frac{S_n^x}{n})}{\varepsilon^2} \text{ by Chebyshev} \\ &= \delta(\varepsilon) + 2M \frac{x(1-x)}{n\varepsilon^2} \end{aligned}$$

Since $x(1-x) \leq \frac{1}{4}$ $\forall x \in [0,1]$ we have

$$\sup_{x \in [0,1]} |P_n(x) - f(x)| \leq \delta(\varepsilon) + \frac{M}{2n\varepsilon^2} \quad (*)$$

By replacing ε with $\frac{1}{n^{1/3}}$ (for example) and choosing n large enough I can make $(*)$ as small as I want (since $\delta(\varepsilon) \rightarrow 0$ as $\varepsilon \rightarrow 0$).

QED

Chernoff's method is especially useful since (18)
 $t \mapsto E(e^{tX})$ is an important function called the moment generating function. More on that later.

Lets look at a special case which will give us the SLLN for bdd r.v.s.

Theorem (Hoeffding's inequality)

Let X_1, X_2, \dots, X_n be iid r.v.s on (Ω, \mathcal{F}, P) . If there exists finite real numbers $a \leq b$ s.t.

$a \leq X_i \leq b$
 p-a.e. $\forall i = 1, \dots, n$, then

$$P\left(\left|\frac{S_n}{n} - \mu\right| \geq \varepsilon\right) \leq 2 \exp\left(-\frac{2n\varepsilon^2}{(b-a)^2}\right)$$

$\forall \varepsilon > 0$ where $S_n = X_1 + \dots + X_n$ and $\mu := E(X_1)$.

Note: The iid assumption in Hoeffding can be relaxed and extended to martingales. We will cover this next quarter when we study dependence in random variables.

Note: The Hoeffding bound gives the exact same bound we derived by hand for our rademacher coin flip r.v.s R_1, R_2, \dots from lecture I:

$$P\left(\left|\frac{S_n}{n}\right| \geq \varepsilon\right) \leq 2e^{-n\varepsilon^2/2}$$

where $\frac{S_n}{n} = \frac{R_1 + \dots + R_n}{n}$ and $-1 \leq R_i \leq 1$.

Lemma: If X is a r.v. on (Ω, \mathcal{F}, P)
s.t. \exists finite $a < b$ s.t. $a \leq X \leq b$ & $E(X) = 0$
Then $E(e^{tX}) \leq e^{t^2(b-a)^2/8} \quad \forall t \geq 0$.

Proof:

Let $w = \frac{x-a}{b-a}$ so that $0 \leq w \leq 1$ and
 $X = wb + (1-w)a$. By convexity we have

$$e^{tX} \leq we^{tb} + (1-w)e^{ta}$$

$$\begin{aligned} \therefore E(e^{tX}) &\leq E(w) e^{tb} + E(1-w) e^{ta} \\ &= \frac{E(X)-a}{b-a} e^{tb} + \frac{b-E(X)}{b-a} e^{ta} \\ &= -\frac{a}{b-a} e^{tb} + \frac{b}{b-a} e^{ta} \\ &= (1-\theta)e^{-u\theta} + \theta e^{u(1-\theta)} \end{aligned}$$

$$\text{where } u = -(b-a)t \quad \& \quad \theta = \frac{b}{b-a}.$$

Notice the assumption $E(X)=0$ implies
 $\theta \in [a, b]$ so that $0 \leq \theta \leq 1$.

If we can show

$$(1-\theta)e^{-u\theta} + \theta e^{u(1-\theta)} \leq e^{u^2/8} \quad (*)$$

$\forall u \in \mathbb{R}$ and $\forall \theta \in [0, 1]$ we are done.

Taking log it is sufficient to show

$$\log((1-\theta)e^{-u\theta} + \theta e^{u(1-\theta)}) \leq u^2/8$$

!!

$$\log(e^{-u\theta}[1-\theta + \theta e^u])$$

!!

$$-u\theta + \log(1-\theta + \theta e^u)$$

!!

$$K(u)$$

$$\text{Now } K'(u) = -\theta + \frac{\theta e^u}{1-\theta + \theta e^u} = -\theta + \frac{\theta}{\theta + (1-\theta)e^{-u}}$$

$$K''(u) = \frac{\theta(1-\theta)e^{-u}}{(\theta + (1-\theta)e^{-u})^2}$$

(19)

Now Taylor's thm gives

$$\begin{aligned} K(u) &= K(0) + uK'(0) + \frac{u^2}{2} K''(u^*) \quad u^* \in [0, u] \\ &= 0 + 0 + \frac{u^2}{2} \left(\underbrace{\frac{\theta}{\theta + (1-\theta)e^{-u^*}}}_{\in [0, 1]} \right) \left(1 - \underbrace{\frac{\theta}{\theta + (1-\theta)e^{-u^*}}}_{\in [0, 1]} \right) \\ \therefore K(u) &\leq \frac{u^2}{8} \quad \underbrace{\leq \frac{1}{4}}_{QED.} \end{aligned}$$

(20)

Proof of Hoeffding's inequality:

We can suppose w.l.g that $E(X_i) = \mu = 0$.

Now

$$\begin{aligned} P\left(\left|\frac{S_n}{n}\right| \geq \varepsilon\right) &= P\left(\left\{\frac{S_n}{n} \geq \varepsilon\right\} \cup \left\{-\frac{S_n}{n} \geq \varepsilon\right\}\right) \\ &= P\left(\frac{S_n}{n} \geq \varepsilon\right) + P\left(-\frac{S_n}{n} \geq \varepsilon\right) \end{aligned}$$

using Chernoff's method for any $t > 0$

$$\begin{aligned} P\left(\frac{S_n}{n} \geq \varepsilon\right) &\leq P\left(e^{tS_n} \geq e^{t\varepsilon n}\right) \\ &\leq \frac{E(e^{tS_n})}{e^{t\varepsilon n}} \quad \text{by Markov's reg.} \\ &= e^{-tn\varepsilon} \prod_{i=1}^n E(e^{tX_i}) \quad \text{by indep.} \\ &\leq e^{t^2(b-a)^2/8} \quad \text{by lemma} \\ &\leq e^{-tn\varepsilon} \underbrace{e^{nt^2(b-a)^2/8}}_{\text{minimized at } t = \frac{4\varepsilon}{(b-a)^2}} \end{aligned}$$

$$\begin{aligned} \therefore P\left(\frac{S_n}{n} \geq \varepsilon\right) &\leq e^{-4n\varepsilon^2/(b-a)^2} e^{n4^2\varepsilon^2/8(b-a)^2} \\ &= e^{-2n\varepsilon^2/(b-a)^2} \end{aligned}$$

Similar arguments give the exact same upper bound for $P\left(-\frac{S_n}{n} \geq \varepsilon\right)$.

QED

Here is an example of the utility of Hoeffding.

(21)

e.g.

Theorem: (SLLN for bounded r.v.s)

Let X_1, X_2, \dots be iid r.v.s on (Ω, \mathcal{F}, P) s.t. $|X_i| \leq c$ for some finite c . Then

$$\frac{S_n}{n} \xrightarrow{n \rightarrow \infty} E(X_i) \text{ p-a.e.}$$

where $S_n = X_1 + \dots + X_n$.

Proof:

By Hoeffding's inequality

$$P\left(\left|\frac{S_n}{n} - E(X_i)\right| \geq \varepsilon\right) \leq 2 \exp\left(-\frac{2n\varepsilon^2}{(2c)^2}\right) \quad \forall \varepsilon > 0.$$

(underbrace)
finite sum over n
for each $\varepsilon > 0$.

\therefore The First Borel-Cantelli lemma applies so

$$\text{that } P\left(\bigcap_{\varepsilon \in \mathbb{Q}^+} \left\{\left|\frac{S_n}{n} - E(X_i)\right| < \varepsilon \text{ a.a.m}\right\} = 1$$

$$\therefore P\left(\frac{S_n}{n} \rightarrow E(X_i)\right) = 1. \quad \text{QED.}$$

Notice that this tells us we have the right definition for $E(X) := \int x dP$, i.e. $E(X)$ tells us the long run average of independent samples of X (with probability 1).

After we cover densities, later in this lecture, we will discuss the Glivenko-Cantelli Theorem which is an important application of the SLNN for bdd r.v.s.

In the Law of the iterated log we needed lower bounds to get tight control on the behavior of $P(S_n \geq \alpha)$. Let's look at one example of a lower bound.

(22)

Theorem: (Paley-Zygmund Ineq.)

If X is a non-negative r.v. s.t. $E(X^2) < \infty$

$$(1-\alpha)^2 \frac{(EX)^2}{E(X^2)} \leq P(X \geq \alpha EX)$$

for all $\alpha \in (0, 1)$.

Proof: First notice that

$$X = X I_{\{X \leq \alpha EX\}} + X I_{\{X \geq \alpha EX\}}.$$

Taking expected values gives

$$\begin{aligned} E(X) &\leq \alpha E(X) + E(X I_{\{X \geq \alpha EX\}}) \\ &\leq \alpha E(X) + \sqrt{E(X^2)} \sqrt{E(I_{\{X \geq \alpha EX\}}^2)} \\ &\quad \text{by Hölder} \\ &= \alpha E(X) + \sqrt{E(X^2)} \sqrt{P(X \geq \alpha EX)} \end{aligned}$$

Now since $E(X), E(X^2)$ are both finite ($E(X) < \infty$ by Hölder) we can shuffle terms around to get the result.

QED.

Note: Jensen's inequality shows $(EX)^2 \leq E(X^2)$

$$\therefore \text{the lower bound } (1-\alpha)^2 \frac{(EX)^2}{E(X^2)} \leq 1$$

Maximal Inequalities

(23)

We used two maximal inequalities in lecture 7 which were custom to our coinflip model. In this section we state general versions and use them later (after we get the CLT) to establish Kolmogorov's 3 series theorem.

Theorem (Kolmogorov's Maximal Inequality)

Let X_1, \dots, X_n be independent r.v.s s.t. $E(X_k^2) < \infty$ and $E(X_k) = 0$ a.s. If $\alpha > 0$ then

$$P\left(\max_{1 \leq k \leq n} |S_k| \geq \alpha\right) \leq \frac{1}{\alpha^2} \underbrace{\text{var}(S_n)}_{= E(S_n^2)}$$

where $S_n := X_1 + \dots + X_n$.

Proof: The proof is exactly similar to the one we did for coin flips in lecture 7. Let's do it again using our theory of integration.

Define $F_k := \{|S_1| < \alpha, \dots, |S_{k-1}| < \alpha, |S_k| \geq \alpha\}$.

$$\begin{aligned} E(S_n^2) &= \int_S S_n^2 dP \\ &\geq \int_S S_n^2 \sum_{k=1}^n I_{F_k} dP \quad \text{since } F_k \text{'s are disjoint so } \sum I_{F_k} \leq 1. \\ &= \sum_{k=1}^n \int_S S_n^2 I_{F_k} dP \quad S_n^2 = (S_n \pm S_k)^2 \\ &= \sum_{k=1}^n \int_S [S_k^2 + 2S_k(S_n - S_k) + (S_n - S_k)^2] I_{F_k} dP \\ &\geq \sum_{k=1}^n \int_S S_k^2 I_{F_k} dP + 2 \int_S S_k I_{F_k} (S_n - S_k) dP \end{aligned}$$

Notice $S_k \in \sigma(X_1, \dots, X_k)$ since S_k is a measurable function of X_1, \dots, X_k . (24)

Also $I_{F_k} \in \sigma(X_1, \dots, X_k)$ since $F_k \in \sigma(X_1, \dots, X_k)$.

Therefore $S_k I_{F_k} \in \sigma(X_1, \dots, X_k)$ so that

$$\sigma(S_k I_{F_k}) \subset \sigma(X_1, \dots, X_k)$$

and similarly

$$\sigma(S_n - S_k) \subset \sigma(X_{k+1}, \dots, X_n)$$

Since the X_k 's are indep we have

$$\begin{aligned} E(S_k I_{F_k} (S_n - S_k)) &= E(S_k I_{F_k}) E(S_n - S_k) \\ &\stackrel{indep}{=} 0 \\ \therefore E(S_n^2) &\geq \int_S \sum_{k=1}^n S_k^2 I_{F_k} dP \quad \text{on } |S_k| \geq \alpha \text{ on } F_k \\ &\geq \alpha^2 \int_S \sum_{k=1}^n I_{F_k} dP \quad \xrightarrow{\text{max}} I_{\{\max_{1 \leq k \leq n} |S_k| \geq \alpha\}} \\ &\geq \alpha^2 P\left(\max_{1 \leq k \leq n} |S_k| \geq \alpha\right) \end{aligned}$$

Q.E.D.

Here is Etemadi's neg which can be used without the moment assumptions on X_k in Kolmogorov's max neg.

Theorem (Etemadi's maximal inequality)

Suppose X_1, X_2, \dots, X_n are independent r.v.s and $\alpha > 0$. Then

$$P\left(\max_{1 \leq k \leq n} |S_k| \geq 3\alpha\right) \leq 3 \max_{1 \leq k \leq n} P(|S_k| \geq \alpha)$$

Proof: The proof is similar to what we did for the coin flips.

$$\text{Let } F_k := \{|S_1| < 3\alpha, \dots, |S_{k-1}| < 3\alpha, |S_k| \geq 3\alpha\}$$

Notice that the F_k 's are disjoint, (25)

$$\bigcup_{k=1}^n F_k = \left\{ \max_{1 \leq k \leq n} |S_k| \geq 3\alpha \right\} \text{ and}$$

$$\begin{aligned} w \in F_k \cap \{|S_n| < \alpha\} &\Rightarrow |S_k(w)| \geq 3\alpha \text{ and } |S_n(w)| < \alpha \\ &\stackrel{(*)}{\Rightarrow} |S_n(w) - S_k(w)| > 2\alpha \\ &\text{and } w \in F_k \end{aligned}$$

Now

$$\begin{aligned} P\left(\max_{1 \leq k \leq n} |S_k| \geq 3\alpha\right) &= P\left(\max_{1 \leq k \leq n} |S_k| \geq 3\alpha, |S_n| < \alpha\right) \\ &\quad + P\left(\max_{1 \leq k \leq n} |S_k| \geq 3\alpha, |S_n| \geq \alpha\right) \\ &\leq P(|S_n| \geq \alpha) + \sum_{k=1}^{n-1} P(F_k \cap \{|S_n| < \alpha\}) \\ &\quad \text{drop } k=n \text{ term} \\ &\quad \text{since } F_n \cap \{|S_n| < \alpha\} = \emptyset \\ &\leq P(|S_n| \geq \alpha) + \sum_{k=1}^{n-1} P(F_k \cap \{|S_n - S_k| > 2\alpha\}) \\ &\quad \text{by } (*) \\ &= P(|S_n| \geq \alpha) + \sum_{k=1}^{n-1} P(F_k) P(|S_n - S_k| > 2\alpha) \\ &\quad \text{since } F_k \subset \sigma(x_1, \dots, x_k) \text{ and} \\ &\quad \{|S_n - S_k| > 2\alpha\} \subset \sigma(x_{k+1}, \dots, x_n) \\ &\leq P(|S_n| \geq \alpha) + \max_{1 \leq k \leq n} P(|S_n - S_k| > 2\alpha) \\ &\quad \text{by } (*) \\ &\leq 3 \max_{1 \leq k \leq n} P(|S_k| \geq \alpha). \end{aligned}$$

QED.

(26)

Characterizing and Constructing probability measures with densities

In undergraduate probability we are taught that probability densities characterize "continuous r.v.s" & probability mass functions characterize "discrete r.v.s".

e.g.

$$\text{if } P(X \in B) = \int_B e^{-x} I_{(0, \infty)}(x) dx$$

then this is the density

but if

$$P(X \in B) = \sum_{k \in B} \binom{n}{k} \left(\frac{1}{3}\right)^k \left(\frac{2}{3}\right)^{n-k}$$

then this is the probability mass function

It was annoying to me that you had to figure something out about X , i.e. continuous or discrete, before trying to compute probabilities. Our general integration theory unifies both cases and gives an accessible method for characterizing and constructing Prob measures.

For the rest of this section let $(\Omega, \mathcal{F}, \mu)$ denote a measure space (unless stated otherwise).

Def: If $f \in \mathcal{Q}(\Omega, \mathcal{F}, \mu)$ then the set function $\int f d\mu : \mathcal{F} \rightarrow \overline{\mathbb{R}}$ mapping

$$A \in \mathcal{F} \mapsto \int_A f d\mu := \int_{\Omega} f I_A d\mu$$

is called the indefinite integral of f with respect to μ .

Theorem ($\int f d\mu$ is σ -additive)

If $f \in \mathcal{Q}(\Omega, \mathcal{F}, \mu)$ then $\int f d\mu$ is countably additive over disjoint \mathcal{F} -sets.

Proof: Let F_1, F_2, \dots be disjoint \mathcal{F} -sets. Use the "2-3" argument.

Step 2: Suppose $f \in \mathcal{H}(\Omega, \mathcal{F})$.

$$\int_{\cup F_k} f d\mu = \int_{\Omega} \sum_{k=1}^{\infty} f I_{F_k} d\mu \quad \text{since } F_k \text{'s disjoint}$$

$$= \int_{\Omega} \limsup_n \sum_{k=1}^n f I_{F_k} d\mu \quad \text{since there are } \geq 0$$

$$\stackrel{\text{By 3(3)}}{=} \limsup_n \int_{\Omega} \sum_{k=1}^n f I_{F_k} d\mu$$

$$\stackrel{\text{By 3(3)}}{=} \limsup_n \sum_{k=1}^n \int_{\Omega} f I_{F_k} d\mu$$

$$= \sum_{k=1}^{\infty} \int_{F_k} f d\mu$$

\therefore the theorem holds over $\mathcal{H}(\Omega, \mathcal{F})$.

Step 3: Suppose $f \in \mathcal{Q}(\Omega, \mathcal{F}, \mu)$.

Certainly $f I_{U_k F_k} \in \mathcal{Q}(\Omega, \mathcal{F}, \mu)$ since

$$(f I_{U_k F_k})^{\pm} = f^{\pm} I_{U_k F_k} \leq f^{\pm}$$

(27)

$$\begin{aligned} \therefore \int_{\cup F_k} f d\mu &\stackrel{\text{def}}{=} \int_{\Omega} f^+ d\mu - \int_{\Omega} f^- d\mu \\ &= \sum_{k=1}^{\infty} \int_{F_k} f^+ d\mu - \sum_{k=1}^{\infty} \int_{F_k} f^- d\mu \\ &= \sum_{k=1}^{\infty} \left[\int_{F_k} f^+ d\mu - \int_{F_k} f^- d\mu \right] \\ &\quad \text{by Big 3 for counting measure.} \\ &= \sum_{k=1}^{\infty} \int_{F_k} f d\mu. \end{aligned}$$

QED.

(28)

Corollary:

i) If $f \in \mathcal{H}(\Omega, \mathcal{F})$ then

$\int f d\mu$ is a measure on (Ω, \mathcal{F})

ii) If $f \in \mathcal{H}(\Omega, \mathcal{F})$ & $\int f d\mu = 1$ then

$\int f d\mu$ is a probability measure on (Ω, \mathcal{F}) .

Definition:

For any two measures μ, ν on (Ω, \mathcal{F}) if $\exists s \in \mathcal{H}(\Omega, \mathcal{F})$ s.t.

$$\nu(A) = \int_A s d\mu, \quad \forall A \in \mathcal{F}$$

Then s is a density of ν with respect to μ .

Here is an important question:

(Are densities unique?)

the next example shows the answer must be

Not without assumptions

e.g. Let $\Omega = \mathbb{R}$ (2a)
 $\mathcal{F} = \{\emptyset, \Omega, (-\infty, a), [0, \infty)\}$
 $\therefore \mathcal{L}'(\Omega) = \int_A 1 d\mathcal{L}' = \int_A 2 d\mathcal{L}' \quad \forall A \in \mathcal{F}$
 these are both densities w.r.t \mathcal{L}' resulting in the same measure.

The following theorem gives sufficient conditions for uniqueness (as a corollary)

Thm: Suppose $f, g \in L^1(\Omega, \mathcal{F}, \mu)$. If $f \in L_1(\Omega, \mathcal{F}, \mu)$ or $g \in L_1(\Omega, \mathcal{F}, \mu)$ or μ is σ -finite then

$$\int f d\mu \leq \int g d\mu \text{ on } \mathcal{F} \iff f \leq g \text{ } \mu\text{-a.e.}$$

Proof:

\Leftarrow : Follows directly from Big 3 [1]

\Rightarrow :

Case 1: Suppose $f \in L_1(\Omega, \mathcal{F}, \mu)$ or $g \in L_1(\Omega, \mathcal{F}, \mu)$. We will show $\mu(f > g) = 0$ by the "indicate what you want trick".

$$\begin{aligned} f I_{\{f > g\}} &\geq g I_{\{f > g\}} \\ \Rightarrow \int_{f > g} f d\mu &\geq \int_{f > g} g d\mu \text{ by Big 3 [1].} \\ \Rightarrow \int_{f > g} f d\mu &= \int_{f > g} g d\mu \text{ since } \int f d\mu = \int g d\mu \text{ on } \mathcal{F}. \\ \Rightarrow \underbrace{\int_{f > g} (f-g) d\mu}_{\geq 0} &= 0 \text{ by Big 3 [2] since } f \in L_1 \text{ or } g \in L_1. \\ \Rightarrow (f-g) I_{\{f > g\}} &= 0 \text{ } \mu\text{-a.e. by a.e. useful facts} \\ \Rightarrow \underbrace{\mu(f > g)}_{\text{since } f(w) > g(w) \Rightarrow (f(w)-g(w)) I_{\{f > g\}(w)} = 1} &= 0 \end{aligned}$$

Case 2: Suppose μ is finite. (30)

First notice that

$$f \leq g \text{ on } \{f=\infty\} \cap \{g=\infty\} \quad (1)$$

$$f \leq g \text{ on } \{f=-\infty\} \cap \{g=\infty\} \quad (2)$$

$$f \leq g \text{ on } \{f=-\infty\} \cap \{g=-\infty\} \quad (3)$$

We also have

$$\mu(\{f=\infty\} \cap \{g=-\infty\}) = 0 \quad (4)$$

otherwise it would contradict the assumption

$$\int f d\mu \leq \int g d\mu. \text{ Now we just show}$$

$$f \leq g \text{ } \mu\text{-a.e. on } \{|f| < \infty\} \cup \{|g| < \infty\}. \quad (5)$$

Let $A_n := \{|f| < n\}$. Since μ is a finite measure

$$f I_{A_n} \in L_1 \text{ & } g I_{A_n} \in L_1.$$

Also

$$\int f I_{A_n} d\mu = \int_{A_n} f d\mu \leq \int_{A_n} g d\mu = \int g I_{A_n} d\mu$$

on \mathcal{F} . Since $f I_{A_n} \in L_1$, Case 1 implies

$$f I_{A_n} \leq g I_{A_n} \text{ } \mu\text{-a.e.}$$

i.e. $f \leq g \text{ } \mu\text{-a.e. on } \{|f| < n\}$

$$\therefore f \leq g \text{ } \mu\text{-a.e. on } \{|f| < \infty\} = \bigcup_{n=1}^{\infty} \{|f| < n\}.$$

A similar argument shows

$$f \leq g \text{ } \mu\text{-a.e. on } \{|g| < \infty\}.$$

Now the union of (1)-(5) gives

$$f \leq g \text{ } \mu\text{-a.e.}$$

Case 3: Suppose μ is σ -finite.

Let $F_k \in \mathcal{F}$ s.t. $\mu(F_k) < \infty$ and $\bigcup_{k=1}^{\infty} F_k = \Omega$.

$$\begin{aligned} \therefore \mu(f > g) &= \sum_{k=1}^{\infty} \underbrace{\mu(\{f > g\} \cap F_k)}_{= \mu_k(f > g) \text{ where } \mu_k(\cdot) := \mu(\cdot \cap F_k)} \\ &= \mu_k(f > g) \end{aligned}$$

Now case 2 applies to the finite measure μ_p since (31)

$$\begin{aligned} \int f d\mu_p &= \int f I_{F_p} d\mu \quad \text{by a "1-2-3" argument} \\ &= \int f d\mu \\ &\quad \bullet_{\text{if } F_p} \\ &\leq \int g d\mu = \int g d\mu_p \end{aligned}$$

\therefore case 2 implies $\mu_p(f > g) = 0$.

$\therefore \mu(f > g) = 0$. QED

Corollary (uniqueness of densities)

Let $f, g \in Q(\Omega, \mathcal{F}, \mu)$. If f or g is integrable or μ is σ -finite then

$$\int f d\mu = \int g d\mu \text{ on } \mathcal{F} \Leftrightarrow f = g \text{ } \mu\text{-a.e.}$$

Note: If P is a probability measure and μ is a measure (over (Ω, \mathcal{F})) then

$$\begin{aligned} i) P(\cdot) &= \int \delta d\mu \Rightarrow \text{both } \int \delta^+ d\mu < \infty \text{ and} \\ &\quad \int \delta^- d\mu < \infty \\ &\Rightarrow \delta \in L_1(\Omega, \mathcal{F}, \mu) \\ &\Rightarrow \delta \text{ is unique } \mu\text{-a.e.} \end{aligned}$$

$$ii) \mu(\cdot) = \int \delta dP \Rightarrow \delta \text{ is unique } P\text{-a.e.}$$

since P is σ -finite.

The next theorem shows how to compute $\int f d\nu$ when ν has density s w.r.t. μ . (32)

Theorem (Slap on the density: $d\nu = s d\mu$)

Let ν and μ be densities on (Ω, \mathcal{F}) where ν has density s w.r.t. μ .

Then

$$f \in Q^\pm(\Omega, \mathcal{F}, \nu) \Leftrightarrow fs \in Q^\pm(\Omega, \mathcal{F}, \mu)$$

and either one implies

$$\int f d\nu = \int f s d\mu \quad \begin{matrix} \text{i.e. } d\nu = s d\mu \\ \text{i.e. } s = \frac{d\nu}{d\mu} \end{matrix}$$

Proof: Again use "1-2-3 argument".

Step 1: Suppose $f \in \mathcal{N}_s(\Omega, \mathcal{F})$ so that

$$f = \sum_{k=1}^n c_k I_{A_k} \text{ for } A_k \in \mathcal{F} \text{ & } c_k \geq 0.$$

Since $s \in \mathcal{N}(\Omega, \mathcal{F})$ clearly

$$f \in Q^-(\Omega, \mathcal{F}, \nu) \Leftrightarrow fs \in Q^-(\Omega, \mathcal{F}, \mu).$$

$$\begin{aligned} \therefore \int f d\nu &= \sum_{k=1}^n c_k \nu(A_k) \\ &= \sum_{k=1}^n c_k \int_A s d\mu \quad \begin{matrix} \text{since } s \text{ is a} \\ \text{density for } \nu \\ \text{w.r.t. } \mu \end{matrix} \\ &= \int \left(\sum_{k=1}^n c_k I_{A_k} \right) s d\mu \quad \text{by little 3} \\ &\quad \underbrace{\hspace{1cm}}_f \end{aligned}$$

and this implies

$$f \in Q^+(\Omega, \mathcal{F}, \nu) \Leftrightarrow fs \in Q^+(\Omega, \mathcal{F}, \mu).$$

Step 2: Suppose $f \in \mathcal{N}(\Omega, \mathcal{F})$. Then the result follows similarly by little 3.

Step 3: From step 2

$$\int f^\pm d\nu = \int f^\pm s d\mu = \int (fs)^\pm d\mu$$

$$\therefore f \in Q^\pm(\Omega, \mathcal{F}, \nu) \Leftrightarrow fs \in Q^\pm(\Omega, \mathcal{F}, \mu) \quad \underline{\text{QED.}}$$

Notation:

(33)

Suppose ν and μ are measures on (Ω, \mathcal{F}) . If ν has a density w.r.t μ I will denote it:

$$\frac{d\nu}{d\mu}$$

↳ Non-negative (in) function mapping $\Omega \rightarrow \mathbb{R}$ s.t.

$$\nu(A) = \int_A \frac{d\nu}{d\mu} d\mu$$

Moreover when I say $\frac{d\nu}{d\mu}$ exists I mean there exists some density δ^ν w.r.t. μ (it will be unique μ -a.e. if $\frac{d\nu}{d\mu} \in L_1(\Omega, \mathcal{F}, \mu)$ or μ is σ -finite)

Theorem (Chain rule)

Suppose ν, ρ, μ are measures on (Ω, \mathcal{F}) with μ σ -finite. If $\frac{d\rho}{d\nu}$ and $\frac{d\nu}{d\mu}$ exists then

$$\frac{d\rho}{d\mu} = \frac{d\rho}{d\nu} \frac{d\nu}{d\mu} \quad \mu\text{-a.e.}$$

↳ even if μ isn't σ -finite
this serves as a density of ρ w.r.t. μ

Proof:

$$\begin{aligned} \int \frac{d\rho}{d\nu} \frac{d\nu}{d\mu} d\mu &= \int \frac{d\rho}{d\nu} d\nu \quad \text{by step in the density} \\ \uparrow &= \int \rho d\nu \quad \text{by step in the density} \\ \in \mathcal{N}(\Omega, \mathcal{F}) &= \rho(\omega) \quad \text{def} \end{aligned}$$

over \mathcal{F} .

Uniqueness follows since μ is σ -finite.

GED

Our final result is important for two reasons. First it will be used to significantly simplify the proof of the Radon-Nikodym theorem. Second it shows that measure theory (for σ -finite measures) can be considered a sub-theory of probability!

(34)

Theorem (Probabilist's world view)

i.e. $\exists A \text{ s.t. } \mu(A) > 0$

If μ is a non-trivial σ -finite measure on a measure space (Ω, \mathcal{F}) , then there exists a density $\delta: \Omega \xrightarrow{\text{onto}} (0, \infty)$ and a probability measure P on (Ω, \mathcal{F}) s.t.

$$d\mu = \delta dP$$

$$\text{i.e. } \mu(A) = \int_A \delta dP \quad \forall A \in \mathcal{F}.$$

Proof:

Let $\Omega = \bigcup_{k=1}^{\infty} A_k$ where $0 < \mu(A_k) < \infty$
 ↓ disjoint
 $A_k \in \mathcal{F}$ to get this just
 $\mu(A_j) = 0$ into an A_k with positive measure.

Choose any sequence w_1, w_2, \dots s.t. $w_k > 0$ and

$$\sum_{k=1}^{\infty} w_k = 1. \quad \text{Now define}$$

$$\delta^* := \sum_{k=1}^{\infty} \frac{w_k}{\mu(A_k)} I_{A_k} \in \mathcal{N}(\Omega, \mathcal{F})$$

$$P(\cdot) := \int_{\Omega} \delta^* d\mu$$

$$\text{Notice } P(\Omega) = \int_{\Omega} \delta^* d\mu$$

$$= \int_{\Omega} \limsup_{n \rightarrow \infty} \sum_{k=1}^n \frac{w_k}{\mu(A_k)} I_{A_k} d\mu$$

$$\stackrel{\text{B3'}}{=} \sum_{k=1}^{\infty} \int_{\Omega} \frac{w_k}{\mu(A_k)} I_{A_k} d\mu$$

$$= 1$$

By σ -additivity of $\int \delta^* d\mu$ and $\delta^* \in \mathcal{N}(\mathbb{R}, \mathcal{F})$
 $P(\cdot)$ is a probability measure s.t. (35)

$$dP = \delta^* d\mu$$

To finish notice that $0 < \delta^*(w) < \infty$ $\forall w \in \mathbb{R}$
and define $\delta(w) := \frac{1}{\delta^*(w)}$ which maps into
 $(0, \infty)$ and is \mathcal{F} by closure.

Now $\forall A \in \mathcal{F}$

$$\mu(A) = \int_A \delta \delta^* d\mu = \int_A \delta dP$$

by step in the density
since $dP = \delta^* d\mu$.

QED

Note: Suppose $\delta \in \mathcal{N}(\mathbb{R}, \mathcal{F})$ and ν, μ are two measures s.t.

$$d\nu = \delta d\mu$$

In an exercise you will show

- i) ν is finite $\Leftrightarrow \delta \in L_1(\mathbb{R}, \mathcal{F}, \mu)$
- ii) ν is σ -finite $\Rightarrow \delta < \infty$ μ -a.e.
- iii) $\delta < \infty$ μ -a.e. $\begin{cases} \text{and} \\ \mu \text{ is } \sigma\text{-finite} \end{cases} \Rightarrow \nu \text{ is } \sigma\text{-finite.}$

e.g. Let X be a r.v. which has a density w.r.t. Lebesgue measure. By this

I mean

$$dPX^{-1} = \delta d\lambda$$

for some $\delta \in \mathcal{N}(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. In particular

$$P(X \in B) = \int_B \delta(x) dx \quad \forall B \in \mathcal{B}(\mathbb{R}).$$

unique \mathcal{F} -a.e. since P is a prob

BTW some people write $P(X \in dx) = \delta(x)dx$ (36)
as synonymous for $dPX^{-1} = \delta dx$. It sorta makes sense since $P(X \in dx) = PX^{-1}(dx)$ and

$$\begin{aligned} P(X \in A) &= \int_A dPX^{-1}(x) && \text{different notation for the same thing} \\ &= \int_A \underbrace{PX^{-1}(dx)}_{\text{think of this as } P(X \in dx)} \\ &= \int_A \delta(x) dx \quad dPX^{-1} = \delta dx \end{aligned}$$

Now suppose $g: \mathbb{R} \xrightarrow{\text{def}} \mathbb{R}$. If $g(X)$ is quasi-integrable we have

$$\begin{aligned} \therefore E(g(X)) &= \int g(x) dP \\ &= \int g(x) dPX^{-1}(x) && \text{by change of variables} \\ &= \int_{\mathbb{R}} g(x) \delta(x) dx && \text{by step in the density.} \\ &&& \text{This is what everybody sees in undergrad probability.} \end{aligned}$$

L-9.

Let X_1, X_2, \dots, X_n be indep. r.v. s.t.

$$X_i = \begin{cases} 1 & \text{with prob } \theta \\ 0 & \text{with prob } 1-\theta \end{cases}$$

where $0 \leq \theta \leq 1$. Set $S_n = X_1 + \dots + X_n$.

S_n has a binomial distribution denoted

$$S_n \sim \text{Bin}(n, \theta).$$

$$\text{Let } \delta(k) = \begin{cases} \binom{n}{k} \theta^k (1-\theta)^{n-k} & \text{if } k=0, \dots, n \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned} \text{Now } P(S_n \in B) &= \int_B \delta(k) d\lambda(k) && \text{Counting measure} \\ &\quad \downarrow \\ \therefore dPS_n^{-1} &= \delta d\lambda \end{aligned}$$

e.g. Probably the most important r.v., besides the coin flip $\text{Ber}(e)$, is the

Gaussian r.v.. A r.v. X is said to be Gaussian, denoted $X \sim N(\mu, \sigma^2)$, if

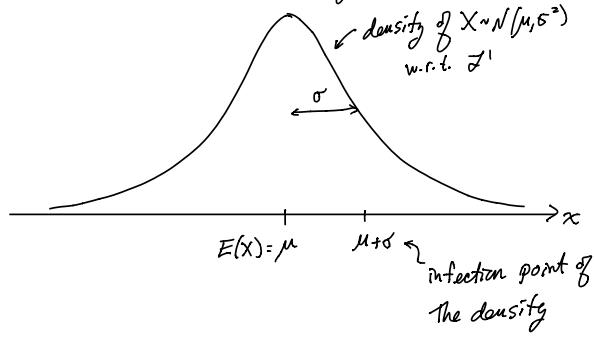
$$dP_{X^{-1}} = \frac{1}{\sqrt{\sigma^2 2\pi}} \exp\left(-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right) dx.$$

Notice that $X \sim N(\mu, \sigma^2)$ implies

$$E(X) = \mu$$

$$\text{var}(X) = \sigma^2$$

so one has the following picture



If X_1, \dots, X_n are r.v.s defined on (Ω, \mathcal{F}, P) then the random vector $X = (X_1, \dots, X_d)^T$ is said to be jointly Gaussian, denoted $X \sim N_d(\mu, \Sigma)$, if

$$dP_{X^{-1}} = \frac{1}{\sqrt{\det(\Sigma \cdot \pi)}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right) d\pi^d(x)$$

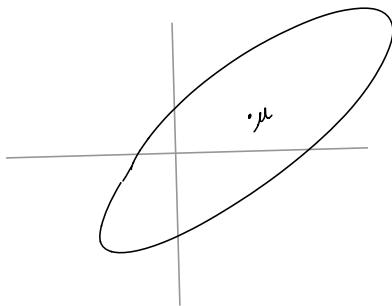
↑
 induced measure on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$

↓
 d dimensional vector
 ↓
 d x d positive definite matrix

In this case

$$\mu = \begin{pmatrix} E(X_1) \\ \vdots \\ E(X_d) \end{pmatrix} \text{ and } \text{cov}(X_i, X_j) = \Sigma_{ij}$$

The contours of $X \sim N_2(\mu, \Sigma)$ look like



Warning! It is possible that both X_1 and X_2 are Gaussian r.v.s but (X_1, X_2) is not jointly Gaussian.

As an exercise check that if

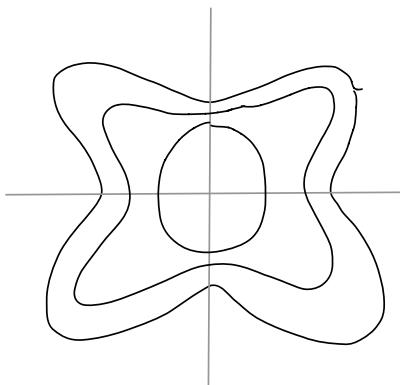
$$S_1(x_1, x_2) \text{ is the density of } N_2(0, \begin{pmatrix} 1 & \frac{1}{2} \\ \frac{1}{2} & 1 \end{pmatrix})$$

$$S_2(x_1, x_2) \text{ is the density of } N_2(0, \begin{pmatrix} 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{pmatrix})$$

Then any random vector (X_1, X_2) with

$$\text{density } S(x_1, x_2) := \frac{S_1(x_1, x_2)}{2} + \frac{S_2(x_1, x_2)}{2}$$

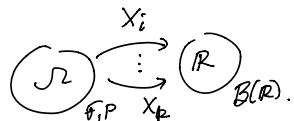
has Gaussian marginals, i.e. $X_1 \sim N_1(0, 1)$ and $X_2 \sim N_1(0, 1)$, but (X_1, X_2) are not jointly Gaussian since the contours of $S(x_1, x_2)$ look like



Glivenko-Cantelli

(39)

In statistics one often observes r.v.s X_1, X_2, \dots, X_n which are iid



Suppose each X_k has density $\delta(x)$ w.r.t. λ , i.e.

$$dP_{X_k} = \delta(x) dx.$$

If $\delta(x)$ is unknown it is natural to consider estimating it with the empirical measure based on X_1, X_2, \dots, X_n defined by

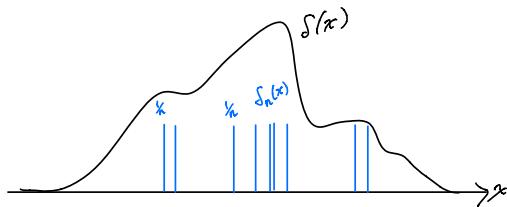
$$dP_n = \left[\frac{1}{n} \sum_{k=1}^n I_{\{X_k\}}(x) \right] d\lambda$$

↑
empirical measure on $(R, B(R))$

:= $\delta_n(x)$

↑ counting measure

Here is the picture:



For some events $B \in B(R)$, such as intervals, one would expect

$$\frac{\#\{X_k \in B\}}{n} = \int_B \delta_n(x) d\lambda(x) \approx \int_B \delta(x) dx = P(X_k \in B)$$

But for other $B \in B(R)$, such as $\{x : X_k = x\}$ one has

$$\frac{1}{n} \int_B \delta_n(x) d\lambda(x) \neq \int_B \delta(x) dx = 0$$

So δ_n doesn't uniformly estimate δ over Borel events. (40)

if we instead compare the c.d.f.'s of the two measures $\delta(x)dx$ and $\delta_n(x)d\lambda(x)$ we get uniform approximation. This is the Glivenko-Cantelli theorem.

As setup for the theorem let

$$F(t) := \int_{(-\infty, t]} \delta(x) dx = P(X_k \leq t)$$

$$F_n(t) := \int_{(-\infty, t]} \delta_n(x) d\lambda(x) = \frac{1}{n} \sum_{k=1}^n I_{(-\infty, t]}(x_k)$$

of k's b/w 1 & n
s.t. $X_k \leq t$

F_n is called the empirical c.d.f

Theorem (Glivenko-Cantelli)

If X_1, X_2, \dots are iid r.v.s on (Ω, \mathcal{F}, P) .

Then

$$\sup_{t \in \mathbb{R}} |F(t) - F_n(t)| \xrightarrow{n \rightarrow \infty} 0 \quad P\text{-a.e.}$$

Proof:

Notice that $I_{(-\infty, t]}(X_1), \dots, I_{(-\infty, t]}(X_n)$ are iid bdd r.v.s. Therefore the SLLN for bdd r.v.s applies and gives

$$\underbrace{\frac{1}{n} \sum_{k=1}^n I_{(-\infty, t]}(X_k)}_{= F_n(t)} \xrightarrow{n \rightarrow \infty} \underbrace{E I_{(-\infty, t]}(X_1)}_{= P(X_1 \leq t) = F(t)} \quad P\text{-a.e.}$$

for each fixed $t \in \mathbb{R}$. Similarly

$$\underbrace{\frac{1}{n} \sum_{k=1}^n I_{(-\infty, t)}(X_k)}_{= F_n(t-)} \xrightarrow{n \rightarrow \infty} \underbrace{E I_{(-\infty, t)}(X_1)}_{= P(X_1 < t) = F(t-)} \quad P\text{-a.e.}$$

(41)

For each $t \in \mathbb{R}$ let $\mathcal{N}_t^\circ \subset \mathcal{N}$ be the measurable event s.t.

- $P(\mathcal{N}_t^\circ) = 1$
- $F_n(t) \rightarrow F(t)$ & $F_n(t^-) \rightarrow F(t^-)$ everywhere on \mathcal{N}_t°

Now fix $\varepsilon > 0$.

Choose m points (depending on ε) s.t.

$$0 < u_1 < u_2 < \dots < u_m < 1$$

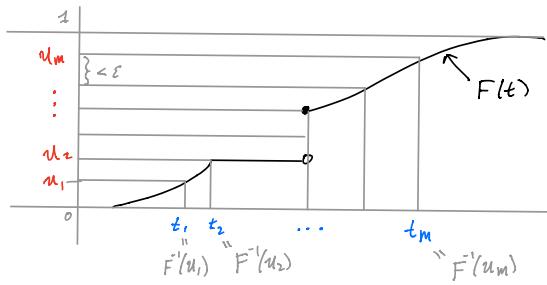
↑ ↑ ... ↑ ↑

with spacing less than ε

and define

$$t_k = F^{-1}(u_k) := \inf\{t : F(t) \geq u_k\}.$$

Here is the picture



Recall the "cdf sandwich Lemma"

$$F(F^{-1}(u)^-) \leq u \leq F(F^{-1}(u)) \quad \forall u \in (0,1).$$

Therefore

$$(*) \quad k=1, \dots, m \Rightarrow F(t_k^-) \leq u_k$$

$$(**) \quad k=2, \dots, m+1 \Rightarrow u_{k-1} \leq F(t_{k-1})$$

(42)

Let $\mathcal{N}_\varepsilon := \bigcap_{k=1}^m \mathcal{N}_{t_k}$ so that

$$w \in \mathcal{N}_\varepsilon \Rightarrow \begin{cases} F_n(t_k) \rightarrow F(t_k) \\ F_n(t_k^-) \rightarrow F(t_k^-) \\ \forall k = 1, \dots, m \end{cases}$$

$\Rightarrow \exists$ a finite $N \equiv N(n, \varepsilon)$ s.t.

$$(i) \quad \sup_{1 \leq k \leq m} |F_n(t_k) - F(t_k)| \leq \varepsilon$$

$$(ii) \quad \sup_{1 \leq k \leq m} |F_n(t_k^-) - F(t_k^-)| \leq \varepsilon$$

$\forall n \geq N.$



e.g. take N to be the max over N_k 's and M_k 's s.t.

$$n \geq N_k \Rightarrow |F_n(t_k) - F(t_k)| \leq \varepsilon$$

$$n \geq M_k \Rightarrow |F_n(t_k^-) - F(t_k^-)| \leq \varepsilon$$

To finish we show

$$w \in \mathcal{N}_\varepsilon \Rightarrow \sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \leq 2\varepsilon, \quad \forall n \geq N(n, \varepsilon)$$

This is sufficient since $P(\mathcal{N}_\varepsilon) = 1$ hence

$$P\left(\bigcap_{\varepsilon \in \mathbb{Q}^+} \mathcal{N}_\varepsilon\right) = 1$$

$$\sup_{\substack{t \in \mathbb{R} \\ t \in \mathcal{N}}} |F_n(t) - F(t)| \rightarrow 0 \quad \text{for these } w's$$

Case 1: $t \in (t_{k-1}, t_k)$ & $n \geq N$.

(43)

$$\begin{aligned} \therefore F_n(t) &\leq F_n(t_{k-1}) \text{ since } P_n([-v, t]) \leq P_n([-v, t_k]) \\ &\leq F(t_{k-1}) + \varepsilon \text{ by (ii)} \\ &\leq u_k + \varepsilon \text{ by (*)} \\ &= u_{k-1} + \underbrace{(u_k - u_{k-1})}_{\leq \varepsilon} + \varepsilon \\ &\leq F(t_{k-1}) + 2\varepsilon \text{ by (**)} \\ &\leq F(t) + 2\varepsilon \text{ by monotonicity} \end{aligned}$$

Also $F(t) \leq F(t_{k-1})$

$$\begin{aligned} &\vdots \quad \leftarrow \text{same as above without extra } \varepsilon \\ &\leq F(t) + \varepsilon \\ &\leq F_n(t) + 2\varepsilon \text{ by (i)} \end{aligned}$$

$$\therefore \forall t \in \mathbb{R}_0 \Rightarrow \sup_{t \in (t_{k-1}, t_k)} |F_n(t) - F(t)| \leq 2\varepsilon, \forall n \geq N$$

Case 2: $t < t_1$ & $n \geq N$.

Similar to case 1 using the fact that $0 \leq F(t_1^-) \leq F(t_1) \leq \varepsilon$ by (*).

$$\therefore \forall t \in \mathbb{R}_0 \Rightarrow \sup_{t < t_1} |F_n(t) - F(t)| \leq 2\varepsilon, \forall n \geq N$$

Case 3: $t > t_m$ & $n \geq N$.

Similar to case 2 using $1 - F(t_m) \leq \varepsilon$ by (**).

$$\therefore \forall t \in \mathbb{R}_0 \Rightarrow \sup_{t > t_m} |F_n(t) - F(t)| \leq 2\varepsilon, \forall n \geq N$$

Finally by cases 1, 2, 3, (i), (ii) we have

$$\forall t \in \mathbb{R}_0 \Rightarrow \sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \leq 2\varepsilon, \forall n \geq N$$

As was to be shown.

QED

Scheffé's Theorem

(44)

The last theorem we will cover in this lecture shows that pointwise convergence of probability densities implies uniform convergence of probabilities over the σ -field.

Theorem (Scheffé)

Let P_n & P be two probability measures on (Ω, \mathcal{F}) which have densities f_n & f w.r.t some measure μ on (Ω, \mathcal{F}) .

If $f_n \rightarrow f$ μ -a.e. as $n \rightarrow \infty$ then

$$\|P_n - P\|_{TV} := \sup_{A \in \mathcal{F}} |P_n(A) - P(A)| \stackrel{(i)}{\leq} \int |\delta f_n - \delta f| d\mu \stackrel{(ii)}{\rightarrow} 0$$

called the total variation norm as $n \rightarrow \infty$.

Proof:

To show (i) notice that

$$\begin{aligned} |P_n(A) - P(A)| &= \left| \int_A \delta f_n d\mu - \int_A \delta f d\mu \right| \\ &= \left| \int_A (\delta f_n - \delta f) d\mu \right| \quad \text{by Big 3 since } \delta f_n, \delta f \in L^1(\Omega, \mathcal{F}, \mu) \\ &\leq \int_A |\delta f_n - \delta f| d\mu \\ &\leq \int_{\Omega} |\delta f_n - \delta f| d\mu \end{aligned}$$

To prove (ii) set $\Delta_n := f - f_n$. Now

$$\begin{aligned} \int |\Delta_n| d\mu &= \int_{\Delta_n > 0} \Delta_n d\mu - \int_{\Delta_n < 0} \Delta_n d\mu \\ &= 2 \int_{\Delta_n > 0} \Delta_n d\mu \quad \text{since } \int_{\Delta_n < 0} \Delta_n d\mu = \int_{\Omega} \delta f_n - \delta f d\mu = 0 \\ &= 2 \int_{\Delta_n > 0} \Delta_n^+ d\mu + \int_{\Delta_n < 0} \Delta_n^- d\mu \\ &\xrightarrow{n \rightarrow \infty} 0 \quad \text{by DCT since } \lim_n \Delta_n^+ = (\lim_n \Delta_n)^+ = 0 \text{ } \mu\text{-a.e.} \\ &\quad \text{and } \Delta_n^+ = (\delta f - \delta f_n)^+ \leq \delta^+ \in L^1(\Omega, \mathcal{F}, \mu) \\ &\quad \text{monotonicity of } L^1 \end{aligned}$$

QED

We will use this theorem in the (45)
lecture on convergence in distribution for a nice
proof of a result by Schoenberg proving
that functions of the form

$$K(h, g) := \int_0^\infty \exp(-t^2 \|h-g\|_{\mathcal{H}}^2) dP(t)$$

are positive definite over all $h, g \in \mathcal{H}$ where
 \mathcal{H} is a Hilbert space (Ref: Steerneman &
van Perlo-ten Kleij).