

Suppose you know that on an average 500 customers visit a supermarket every day. Also suppose it is known that every customer spends <sup>Rs. 200</sup> ~~Rs. 200~~ on an average in each of his visit. Now can we compute the average money received by the supermarket in a day (i.e. average money spent by all customers in a day)?

A any school kid will say the answer is  $\text{Rs. } 200 \times 500 = 100000$

Now let us try to answer this question using the theory (n.v.) we learnt till now:

Suppose  $X_i$  is the amount of money spent by customer  $i$  in a visit. We are given that  $E[X_i] = 200 \quad \forall i$ . ~~Now consider~~

~~$$E[X_1 + X_2 + \dots + X_N]$$~~

Let  $N$  be the no. of customers visiting in a day. we are given that  $E[N] = 500$ . Now consider a new n.v.:

$$X = X_1 + X_2 + \dots + X_N \quad (\text{sum of random numbers of random variables!})$$

In words,  $X$  is the total money received by supermarket in a day. We want to compute  $E[X]$ .

i.e.  $E\left[\sum_{i=1}^N X_i\right]$

This indeed seems to be a very difficult problem  
How do you find exp. of sum of random number of  
random variables? That too we are not given the  
distributions of  $X_i$  &  $N$  !!

But somehow a school kid knows:

$$E\{X\} = E\left[\sum_{i=1}^N X_i\right] \stackrel{??}{\stackrel{\text{How?}}{=}} E\{N\} E\{X_i\} = 500 \times 200 = 100000$$

(Though we are <sup>reasonably</sup> not as clever as the kid, let us try to prove  
his statement).

With a little thought we would realize that if  
it were <sup>exp</sup> sum of fixed number of r.v.'s then we know the  
answer is sum of expectations of r.v. So can we exploit  
this? In particular can we (say) condition on  $N$  (i.e.  
assume  $N = n$  say)?

i.e. Consider the random variable  $X/N = n$ .

we know  $X/N = n$  is nothing but  $\sum_{i=1}^n X_i$

$$\text{So } E[X/N = n] = E\left[\sum_{i=1}^n X_i\right] = n E\{X_i\} = 200n$$

Now  $\downarrow$  expectation seems to be something dependent on the  
value taken by the conditioned r.v.  $N$  !!



In other words, we can view  $E\{X/N=n\}$  as a function of the random variable  $N$ !

How? →

Consider  $g(n) = 200n$

and  $g(N) = 200N$

This is ~~exactly the random variable~~ <sup>quantity</sup>  $E\{X/N=n\}$

→ We denote this random variable as  $E\{X/N\}$

In summary,

$E\{X/N\}$  is a random variable which takes on values  $E\{X/N=n\} \forall n$  ( $n$  is the values taken by  $N$ ).

~~Now  $E\{X/N\}$~~   $E\{X/N\}$  is called the conditional expectation of  $X$  given  $N$

→ Now  $E\{X/N\}$  is a random variable, in fact a function of r.v.  $N$ ! So we know how to compute its expectation!

$$E\{E\{X/N\}\} = E\{g(N)\} = \sum_{n} g(n) P\{N=n\} = 200 \times 500 = 100000$$

~~We claim that~~

$$\underline{\underline{E\{E\{X/N\}\} = E\{X\}}}$$

Finally we got the answer the kid got but we still need to prove the kid's claim that

$$\underline{\underline{E\{E\{X/N\}\} = E\{X\} !}}$$

We will use more general results in the following:

- Consider two RV's  $X, Y$  (joint dist. events) and three cases we were always considering:
- (a)  $X, Y$  both discrete
  - (b)  $X, Y$  jointly conts
  - (c)  $X/Y=y$  conts. &  $Y$  is discrete
  - (d)  $Y/X=x$  discrete &  $X$  is conts
- $\left. \begin{array}{l} \text{---} \\ \text{---} \end{array} \right\} \begin{array}{l} X \text{ is conts.} \\ Y \text{ is discrete} \end{array}$
- Recall that ~~remembered~~ in each of these <sup>4</sup> cases we placed total prob. rule & Baye's rule!

~~We can summarize these cases as two cases:~~

- (i)  ~~$X$  and  $X/Y=y$  are conts. i.e.  $f_x, f_{X/Y}$  are p.d.f's~~
- (ii)  ~~$X$  and  $X/Y=y$  are discrete i.e.  $f_x, f_{X/Y}$  are p.m.f's~~

~~(iii)  $X$  and  $X/Y=y$  are discrete i.e.  $f_x, f_{X/Y}$  are p.m.f's~~



Lets prove  $E[E[g(x)/y]] = E[g(x)]$  for ~~all~~

~~the (i) cases (and the others)~~ all the <sup>fair</sup> ~~the~~ general cases.

Proof  ~~$E[g(x)/y]$~~  Now suppose ~~the~~  $X, X/Y=y$  are conts.

$$E[g(x)/y=y] = \int_{-\infty}^{\infty} g(x) f_{X/Y}(x/y) dx$$

This is a function of  $y$ ! Let it be  $h(y)$

$$h(y) = \int_{-\infty}^{\infty} g(x) f_{X/Y}(x/y) dx$$

$E[g(x)/y]$  is the random variable  $h(y)$ .

(BTW you can show that  $h(y)$  is a valid n.v. if  $g(x)$  is a n.v. &  $E[g(x)] < \infty$  (not in this class))

Now,  $E[g E[g(x)/y]] = E[h(y)]$

~~Suppose  $Y$  is conts~~  $= \int_{-\infty}^{\infty} h(y) f_Y(y) dy$  if  $Y$  is conts

$= \sum_{y} h(y) f_Y(y)$  if  $Y$  is discrete

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x) f_{X/Y}(x/y) dx f_Y(y) dy = \int_{-\infty}^{\infty} g(x) \left[ \int_{-\infty}^{\infty} f_{X/Y}(x/y) f_Y(y) dy \right] dx$$

(if  $Y$  is conts)

$$= \int_{-\infty}^{\infty} g(x) f_X(x) dx = E[g(x)]$$
  

$$= \sum_{y} \int_{-\infty}^{\infty} g(x) f_{X/Y}(x/y) dx f_Y(y) = \int_{-\infty}^{\infty} g(x) \left[ \sum_{y} f_{X/Y}(x/y) f_Y(y) \right] dx = \int_{-\infty}^{\infty} g(x) f_X(x) dx$$

(if  $Y$  is discrete)

$\downarrow$   
 $f_X(x)$

So proof is done for two cases.

Next case is say  $X, X/Y=y$  are discrete r.v's.

Then,

$$E[g(X)/Y=y] = \sum_{\neq x} g(x) f_{X/Y}(x/y)$$

This is a function of  $y$ ! Let it be  $h(y)$

$$h(y) = \sum_{\neq x} g(x) f_{X/Y}(x/y)$$

$\Rightarrow E[g(X)/Y]$  is the random variable  $h(Y)$

(Again if  $g(x)$  is a r.v. &  $E[g(x)] < \infty$ , then  $h(y)$  is a r.v.!)  
valid

$$E[E[g(X)/Y]] = E[h(Y)]$$

$$= \begin{cases} \sum_{\neq y} h(y) f_Y(y) = \sum_{\neq y} \sum_{\neq x} g(x) f_{X/Y}(x/y) f_Y(y) & \text{if } Y \text{ is discrete} \\ \int_{-\infty}^{\infty} h(y) f_Y(y) dy = \int_{-\infty}^{\infty} \sum_{\neq x} g(x) f_{X/Y}(x/y) f_Y(y) dy & \text{if } Y \text{ is cont.} \end{cases}$$

$$= \begin{cases} \sum_{\neq x} g(x) \left[ \sum_{\neq y} f_{X/Y}(x/y) f_Y(y) \right] = \sum_{\neq x} g(x) f_X(x) = E[X] & \text{if } Y \text{ is discrete} \\ \sum_{\neq x} g(x) \left[ \int_{-\infty}^{\infty} f_{X/Y}(x/y) f_Y(y) dy \right] = \sum_{\neq x} g(x) f_X(x) = E[X] & \text{if } Y \text{ is cont.} \end{cases}$$

So all cases are proved. Hence now we are as smart as a school kid!! :D



Now we can easily verify the following properties of conditional expectation:

(i) linearity:  $E[a_1 h_1(x) + a_2 h_2(x) / y] = a_1 E[h_1(x) / y] + a_2 E[h_2(x) / y]$   
Cond. exp.

(ii) If  $h(x) \geq 0$ , then  $E[h(x) / y] \geq 0$

(with prob. 1)  $\Downarrow$  as a result (almost surely) (almost surely)

If  $X_1 \geq X_2$ , then  $E[X_1 / y] \geq E[X_2 / y]$

(almost surely)

(almost surely)

(iii)  $E[c / y] = c$  (c is const.)

(c is const.)

(iv)  $E[h(y) / y] = h(y)$

also if  $X, Y$  are independent then

$$E[g(x) / y] = E[g(x)] \text{ etc.}$$

also if  $X, Y$  are dependent then

$$E[g(x)h(y) / y] = h(y) E[g(x)]$$

\* (Realize, where we used this in supermarket problem)

Consider series of coin flips with usual assumptions. Suppose we are interested in average no. flips for seeing the first 'HT'. Here's a way to do it using the concept of conditional expectation:

Let  $X_1$ : r.v. representing no. flips for seeing first 'H' (geometric r.v.)  
 $X_2$ : " " " " " " 'HT'

Let's look at the following r.v. which reads as "no. flips for seeing first HT given first head was at  $x_1$ 'th toss":

$$X_2 / X_1 = x_1 = \begin{cases} x_1 + 1 & 1-p \\ x_1 + 2 & p(1-p) \\ x_1 + 3 & p^2(1-p) \\ \vdots & \vdots \end{cases}$$

→ note that this is indeed a valid pmf and indeed  $X_2 / X_1 = x_1$  is a valid r.v.

$$\begin{aligned} \text{Now } E[X_2 / X_1 = x_1] &= \sum_{i=1}^{\infty} (x_1 + i) p^{i-1} (1-p) \\ &= x_1 \sum_{i=1}^{\infty} p^{i-1} (1-p) + \underbrace{\sum_{i=1}^{\infty} i p^{i-1} (1-p)}_{\text{exp. of geometric r.v. with parameter } 1-p!} \\ &= x_1 + \frac{1}{1-p} \end{aligned}$$

$$\Rightarrow E[X_2 / X_1] = x_1 + \frac{1}{1-p} \quad \left. \vphantom{E[X_2 / X_1]} \right\} \text{the cond. exp. r.v.}$$

$$\Rightarrow E[X_2] = E\{E[X_2 / X_1]\} = E\left\{x_1 + \frac{1}{1-p}\right\} = E[X_1] + \frac{1}{1-p} = \frac{1}{p} + \frac{1}{1-p}$$

(\* Put  $p = \frac{1}{2}$  and realize what is happening)



Now suppose we want "avg. no. flips for seeing the first 'HH'."

Again let  $X_1$ : ~~no.~~ no. flips for seeing first head (geometric r.v.)  
 $X_2$ : " " " " 'HH'

Again look at:

$$Y = X_2 / X_1 = x_1 = \begin{cases} x_1 + 1 & p \\ x_1 + 1 + X_2 & 1-p \end{cases} \rightarrow \text{again valid pmf} \quad \textcircled{\text{II}}$$

But constraint this defn. and  $\textcircled{\text{I}}$ . In fact we have never seen such a defn! The values taken by  $Y$  must be numbers. Here it looks like <sup>fraction</sup> another r.v.!

With a little thought you can convince yourself that this is indeed a valid defn.  $\rightarrow$  just means is:

To see this let  $X_2$  r.v. take values in set  $E$  &  $P\{X_2 = \delta_i\} = p_i$   
 (of course  $\sum_{i=1}^{\infty} p_i = 1, p_i \geq 0$ ) we can re-write  $\textcircled{\text{II}}$  as:

$$X_2 / X_1 = x_1 = \begin{cases} x_1 + 1 & p \\ x_1 + 1 + \delta_i & (1-p)p_i \end{cases} \rightarrow \text{again valid pmf and indeed a "usual way of def'ing a r.v.!"}$$

~~any case~~ Now,

$$E\{X_2 / X_1 = x_1\} = (x_1 + 1)p + \sum_{\delta_i} (x_1 + 1 + \delta_i)(1-p)p_i = x_1 + 1 + (1-p) \sum_{\delta_i} \delta_i p_i$$

$$= x_1 + 1 + (1-p)E\{X_2\}$$

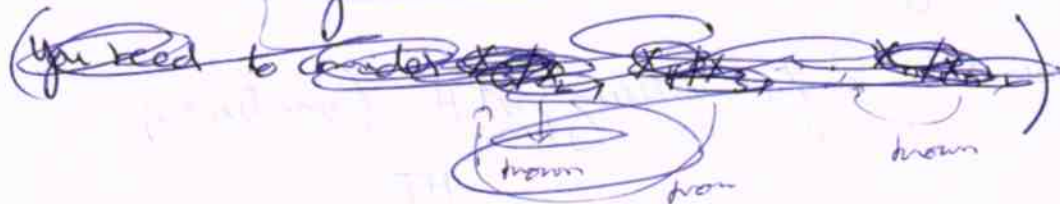
$$\Rightarrow E\{X_2 / X_1\} = x_1 + 1 + (1-p)E\{X_2\}$$

$$\Rightarrow E\{X_2\} = E\{E\{X_2 / X_1\}\} = E\{x_1 + 1 + (1-p)E\{X_2\}\}$$

$$\Rightarrow pE\{X_2\} = 1 + E\{X_1\} \Rightarrow E\{X_2\} = \frac{1 + \frac{1}{p}}{p} = \frac{p+1}{p^2}$$

(\* put  $p = \frac{1}{2}$  compare with prev. result. Explain intuitively!)

~~How can you think about recursive heads here?~~



~~easy to see this if  $E\{X_2\} = E\{X_1\} + 1$  this recursive relation of~~  
~~we know boundaries  $E\{X_2\}, E\{X_1\}$~~

→ Another eg. where cond. exp. is used:

Consider a quicksort algorithm where pivot is always chosen as the first number. Lets do amortized analysis of computational cost. In other words assume the algorithm gets all possible orderings of the same 'n' numbers with equal probability and we want to find the average <sup>no.</sup> comparisons<sup>†</sup> need to be done to sort the 'n' numbers.

Let  $C_n$  be the r.v. representing no. comparisons for sorting 'n' numbers. It is a random variable because it depends on the order in which ~~no~~ numbers are given!

Let  $Y_n$  be the <sup>collect</sup> position of pivot (obtained after doing  $n-1$  comparisons)

Consider the r.v.  $C_n / Y_n = y$



(By the nature of quicksort alg. we have)

$$C_n / X_n=y = C_{y-1} + C_{n-y} + n-1$$

$$\Rightarrow E\{C_n / X_n=y\} = E\{C_{y-1}\} + E\{C_{n-y}\} + n-1 \\ = \bar{C}_{y-1} + \bar{C}_{n-y} + n-1$$

(let us denote  $E\{C_n\} = \bar{C}_n$ ).

Now we can consider the random variable  $E\{C_n / X_n\}$  which takes values  $E\{C_n / X_n=y\}$  for  $y=1$  to  $n$  and of course the ~~prob~~ dist. of  $X_n$  is uniform dist. (no ~~prob~~  $P\{X_n=y\} = \frac{1}{n}$  for  $y=1$  to  $n$ ).

$$\Rightarrow \bar{C}_n = E\{C_n\} = \cancel{E\{C_n\}} E\{E\{C_n / X_n\}\} = \sum_{y=1}^n (\bar{C}_{y-1} + \bar{C}_{n-y} + n-1) \frac{1}{n} \\ = \cancel{(n-1)} + \cancel{\frac{2}{n}(\bar{C}_0 + \bar{C}_n)}$$

$$\Rightarrow n\bar{C}_n = n(n-1) + 2 \sum_{y=1}^{n-1} \bar{C}_y$$

solving recursion we get,  $\bar{C}_n \sim 2(n+1) \log(n+1)$   
which is a well-known result.

→ another eg. for use of cond. exp.:

Suppose  $X, Y$  are jointly conts. then,

$$E\{g(x, y)\} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f_{XY}(x, y) dx dy.$$

$$\text{If } X, Y \text{ are discrete then } E\{g(x, y)\} = \sum_{x=1}^{\infty} \sum_{y=1}^{\infty} g(x, y) f_{XY}(x, y).$$

Now can we still compute  $E\{g(x,y)\}$  in the ~~case~~ <sup>case</sup>:

\*  $X/y$  is conts.,  $Y$  is discrete  $\Leftrightarrow X$  is conts. &  $Y/x$  is discrete  
(~~in other words, if  $X$  is conts.,  $Y$  is discrete~~)

In this case we noted many times that  $F_{xy}$  obviously exists but there is no notion of " $f_{xy}$ ". ( $f_{xy}$  does not make sense of for this case)

Using notion of conditional exp. we can still compute  $E\{g(x,y)\}$ :

$$\text{Consider } E\{g(x,y)/Y=y\} = E\{g(x,y)/Y=y\}$$
$$= \int_{-\infty}^{\infty} g(x,y) f_{X/Y}(x/y) dx$$

exp. of function of  $x$   
wrt.  $f_{X/Y}$

Now again we can consider the random variable  $Y$  which takes values (at  $Y=y$ ). Now:

$$E\{g(x,y)\} = E\{E\{g(x,y)/Y\}\}$$

exp. of function of  $Y$  wrt.  $f_Y$

$$= \sum_{+y} E\{g(x,y)/Y=y\} f_Y(y)$$

$$= \sum_{+y} \int_{-\infty}^{\infty} g(x,y) f_{X/Y}(x/y) dx f_Y(y)$$



$$\Rightarrow E[g(x, y)] = \int_{-\infty}^{\infty} \left( \sum_{y} g(x, y) f_{X/Y}(x/y) f_Y(y) \right) dx$$

(this is a useful expression.)

also by Baye's rule we have:

$$\begin{aligned} \Rightarrow E[g(x, y)] &= \int_{-\infty}^{\infty} \sum_{y} g(x, y) f_{Y/X}(y/x) f_X(x) dx \\ &= \sum_{y} \left( \int_{-\infty}^{\infty} g(x, y) f_{Y/X}(y/x) f_X(x) dx \right) \end{aligned}$$

(this is also a useful expression.)

In lecture we saw how this can be used as starting step for Bayesian Decision Theory.

### Conditional Variance

Similar to notion of conditional exp. we can define conditional variance:

$$\text{var}(X/Y) = E[X^2/Y] - (E[X/Y])^2 \quad - \textcircled{I}$$

↓  
(new defn.)  
conditional variance

again a p.v. and  
in fact a function of Y.

Now we also know through notion of  $\text{var}(x)$  and notion of conditional exp. ~~we have~~ the following:

$$\text{var}(x) = E\{x^2\} - (E\{x\})^2 \quad (\text{defn. of } \text{var}(x))$$

$$= E\{E\{x^2/y\}\} - (E\{E\{x/y\}\})^2 \quad (\text{result we proved in case of cond. exp.})$$

Ⓓ

using Ⓓ & Ⓓ we can prove (see assign.) the following <sup>useful</sup> result:

$$\text{var}(x) = \text{var}(E\{x/y\}) + E\{\text{var}(x/y)\} \quad \text{Ⓓ}$$

Let's look at an application of this:

Let us compute the variance of  $X = \sum_{i=1}^N X_i$

(This is supermarket eg. assume again  $X_i$  &  $N$  are independent)

$$\begin{aligned} \text{var}(X/N=n) &= \text{var}\left(\sum_{i=1}^N X_i / N=n\right) = \text{var}\left(\sum_{i=1}^n X_i / N=n\right) \\ &= \text{var}\left(\sum_{i=1}^n X_i\right) \\ &= n \text{var}(X_i) \quad \text{if } X_i \text{ are iid.} \end{aligned}$$

$$\Rightarrow \text{var}(X/N) = N \text{var}(X_i)$$

(note that this is indeed a fraction of  $N$ )

Using Ⓓ:

$$\begin{aligned} \text{var}(X) &= \text{var}(E\{X/N\}) + E\{N \text{var}(X_i)\} \\ &= \text{var}(NE\{X_i\}) + E\{N\} \text{var}(X_i) = (E\{X_i\})^2 \text{var}(N) + E\{N\} \text{var}(X_i) \end{aligned}$$

(we will later use these results in C.L.T.)  
Central limit theorem.



mgf of a mrv

In case of a r.v.  $X$  we have defined mgf as:  $M_X(\delta) = E[e^{\delta X}]$ .

Now for a m.r.v.  $X$  we define mgf analogously:  $M_X(\underline{\delta}) = E[e^{\delta^T X}]$

If  $X$  is <sup>the</sup> a collection of the r.v.s  $X_1, X_2, \dots, X_n$ , then ~~it is also known~~ is also known as the joint mgf of  $X_1, X_2, \dots, X_n$  and is represented as follows:

$$M_{X_1, X_2, \dots, X_n}(\underbrace{\delta_1, \delta_2, \dots, \delta_n}_{\underline{\delta}}) = M_X(\underline{\delta}) = E[e^{\delta^T X}] = E[e^{\delta_1 X_1 + \delta_2 X_2 + \dots + \delta_n X_n}]$$

Now if  $X_1, X_2, \dots, X_n$  are independent, then we already know that expectations factorize here we have:

$$\begin{aligned} M_{X_1, \dots, X_n}(\delta_1, \dots, \delta_n) &= E[e^{\delta_1 X_1 + \dots + \delta_n X_n}] = E[e^{\delta_1 X_1} e^{\delta_2 X_2} \dots e^{\delta_n X_n}] \\ &= E[e^{\delta_1 X_1}] E[e^{\delta_2 X_2}] \dots E[e^{\delta_n X_n}] \quad \left. \begin{array}{l} X_1, X_2, \dots, X_n \\ \text{are} \\ \text{independent} \end{array} \right\} \\ &= M_{X_1}(\delta_1) M_{X_2}(\delta_2) \dots M_{X_n}(\delta_n). \end{aligned}$$

Hence if  $X_1, \dots, X_n$  are independent then the joint mgf factorizes!

Also it is easy to see that:

$$\left. \frac{\partial^{m_1 + m_2 + \dots + m_n} M_X(\underline{\delta})}{\partial \delta_1^{m_1} \partial \delta_2^{m_2} \dots \partial \delta_n^{m_n}} \right|_{\underline{\delta} = 0} = E[X_1^{m_1} X_2^{m_2} \dots X_n^{m_n}]$$

↙ This shows that all moments can be "generated" from the mgf of  $X$ !

One can also show that (not in this course) ~~that~~ given the joint mgf, the joint dist. function of  $X_1, \dots, X_n$  is determined and vice-versa.

Here <sup>joint</sup> mgf is indeed a useful function for characterizing the collection of r.v.s.

→ Another look at mean vector & covariance matrix of m.s.v.

Given a m.s.v.  $X = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}$ , we have defined:

$$E[X] \equiv \begin{bmatrix} E[X_1] \\ \vdots \\ E[X_n] \end{bmatrix}$$

(mean vector of  $X$ ) (vector of means)

&

$$\text{cov}(X) \equiv \text{matrix with } (i,j)^{\text{th}} \text{ element as } \text{cov}(X_i, X_j).$$

(covariance matrix of  $X$ )

Now let's look at a linearly related m.s.v.:  $Y = \underline{A}X + \underline{b}$ . Note that  $\underline{A}$  is a  $n \times n$  matrix &  $\underline{b}$  is a  $n \times 1$  vector.  $Y$  is another m.s.v. ~~linearly related~~ which is a linear function of  $X$ .

(affine)

It is an exercise to show that: (at least take  $n=2, 3$  and convince yourself)

$$E[Y] = \underline{A} E[X] + \underline{b}$$

$$\text{cov}(Y) = \underline{A} \text{cov}(X) \underline{A}^T \quad \textcircled{I}$$

(in fact there are nearly alternate ways of stating the linearity prop. of  $E$ )



Comparing the diagonal terms in the matrix equality (I) we have:

$$\text{var}(Y_i) = \underline{a}_i^T \text{Cov}(X) \underline{a}_i \quad (\text{why?})$$

(Here  $\underline{a}_i$  is the  $i^{\text{th}}$  column in the matrix  $\underline{A}$ )

Now  $\text{var}(Y_i) \geq 0 \Rightarrow \underline{a}_i^T \text{Cov}(X) \underline{a}_i \geq 0$  ~~for all~~ <sup>irrespective of what is</sup>  $\underline{a}_i$

Also by the very construction of  $\text{Cov}(X)$ , the  $(i, j)^{\text{th}}$  element =  $\text{Cov}(X_i, X_j)$   
 $(j, i)^{\text{th}}$  element =  $\text{Cov}(X_j, X_i)$   
=  $\text{Cov}(X_i, X_j)$

Here  $\text{Cov}(X)$  is a matrix which satisfies two properties:

(i) it is symmetric

(ii)  $\forall \underline{a} \in \mathbb{R}^n, \underline{a}^T \text{Cov}(X) \underline{a} \geq 0$

Matrices satisfying the above two properties are well-studied and are known as positive semi-definite (psd) matrices!

In other words, covariance matrix of any m.s.v. is psd.

On a similar note ~~we~~ <sup>there</sup> are matrices which are positive definite (pd). They satisfy the ~~two~~ below two prop:

(i) Symmetric

(ii)  $\forall \underline{a} \in \mathbb{R}^n, \underline{a} \neq \underline{0}$  <sub>(non-zero vector)</sub>,  $\underline{a}^T \text{Cov}(X) \underline{a} > 0$

Of course all pd matrices are psd but converse is not true.

One interesting property of pd matrices is worth noting:

$$M \text{ is pd} \iff M = LDL^T$$

where  $L$  is an orthogonal matrix

$D$  is a diagonal matrix with pos entries.

(we will later see that such a decomposition of matrix is known as the eigenvalue decomposition of  $M$ !)

Rather than attempting to prove this (we will do it later) we will not understand the ~~the~~ relation.

An orthogonal matrix is a matrix whose column vectors are  $(l_i)$

(i) Unit length  $\|l_i\| = 1$  ( $\& l_i^T l_i = 1$ )

(ii) orthogonal  $\Rightarrow l_i^T l_j = 0$  ( $i \neq j$ ).

With this it is easy to see that  $L^{-1} = L^T$ . ( $\because L^T L = I, L L^T = I$ ).

Also  $|\det(L)|$  is nothing but the ( $n$ -dim) area of the hyper-rectangle formed by  $l_1, \dots, l_n$  vectors; ~~for already~~ <sup>which</sup> which is 1 since each of  $l_i$  is a unit vector.

This says that  $\det M = \det L \det D \det L^T = \det D = \prod_{i=1}^n d_i$   
~~product of~~  
 (assuming  $D = \begin{bmatrix} d_1 & 0 \\ 0 & d_n \end{bmatrix}$ )

Also,  $M^2 = LDL^T LDL^T = LD^2L^T$  ||| by  $M^p = LD^pL^T$ !  
 (at least for rational  $p \geq 0$ )

~~Consider~~ ~~the~~ ~~matrix~~

Consider  $N = LD^{-1}L^T$ . It is again easy to see that  $MN = I, NM = I$ .

$\Rightarrow M^{-1} = LD^{-1}L^T$  ||| by  $M^n = LD^nL^T$ !  
 (at least for all rational  $n < 0$ )

(+)



In other words computing <sup>only</sup> ~~all~~ rational powers of  $M$  is trivial & given the eigen-value decomposition of  $M$ !

In the following we will use another important feature of this decomposition.

## → Multivariate Normal r.v.

In one of prev. lectures we defined a multivariate std. Normal r.v.  
 Now we define  $X$  to be multivariate Normal & multivariate Gaussian r.v. iff  $\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$

$$f_X(\underline{x}) = \frac{1}{(2\pi)^{n/2} |\underline{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\underline{x}-\underline{\mu})^T \underline{\Sigma}^{-1}(\underline{x}-\underline{\mu})} \quad \forall \underline{x} \in \mathbb{R}^n$$

(where  $\underline{\Sigma}$  is a pd matrix).

Also, if this holds,  $x_1, x_2, \dots, x_n$  are said to be jointly Normal or jointly Gaussian.

Let us first prove that  $f_X$  is a valid pdf. Indeed it is  $\geq 0$ .  
 Need to check if area under it is 1:

i.e. TST  $\iint \dots \int e^{-\frac{1}{2}(\underline{x}-\underline{\mu})^T \underline{\Sigma}^{-1}(\underline{x}-\underline{\mu})} d\underline{x} = (2\pi)^{n/2} |\underline{\Sigma}|^{1/2}$

i.e. TST  $\iint \dots \int e^{-\frac{1}{2} \underline{y}^T \underline{\Sigma}^{-1} \underline{y}} d\underline{y} = (2\pi)^{n/2} |\underline{\Sigma}|^{1/2}$   
 (Put  $\underline{y} = \underline{x} - \underline{\mu}$ )

$$\text{i.e. } \underline{\underline{TST}} \int \dots \int e^{-\frac{1}{2} \underline{y}^T \underline{\Sigma}^{-1} \underline{L}^T \underline{y}} \underline{dy} = (2\pi)^{n/2} \prod_{i=1}^n d_i^{1/2}$$

$$\left( \text{we } \underline{\Sigma} = \underline{L} \underline{D} \underline{L}^T \Rightarrow \underline{\Sigma}^{-1} = \underline{L} \underline{D}^{-1} \underline{L}^T \text{ and } \det \underline{\Sigma} = \prod_{i=1}^n d_i \right)$$

(events because  $\underline{\Sigma}$  is pd)  $\rightarrow \begin{bmatrix} d_1 & 0 \\ 0 & d_n \end{bmatrix}$

$$\text{i.e. } \underline{\underline{TST}} \int \dots \int e^{-\frac{1}{2} \underline{z}^T \underline{D}^{-1} \underline{z}} \underline{dz} = (2\pi)^{n/2} \prod_{i=1}^n d_i^{1/2}$$

$$\left( \text{Put } \underline{z} = \underline{L}^T \underline{y} \Rightarrow \underline{y} = \underline{L} \underline{z} \Rightarrow |\underline{J}| = |\det \underline{L}| = 1 \right)$$

(magnification factor Jacobian)

Now LHS =

$$\text{i.e. } \underline{\underline{TST}} \int \dots \int e^{-\frac{1}{2} \sum_{i=1}^n z_i^2 / d_i} dz_1 \dots dz_n$$

$$= \prod_{i=1}^n d_i^{1/2} \int \dots \int e^{-\frac{1}{2} \sum_{i=1}^n z_i^2} dt_1 \dots dt_n \quad \left( \text{Put } \frac{z_i}{\sqrt{d_i}} = t_i \right)$$

$$= \prod_{i=1}^n d_i^{1/2} \prod_{i=1}^n \int e^{-\frac{1}{2} t_i^2} dt_i = (2\pi)^{n/2} \prod_{i=1}^n d_i^{1/2} = \text{RHS. Here proved.}$$

Also in the process we have showed that the transformation

$$\underline{x} = \underline{y} + \underline{\mu} = \underline{L} \underline{z} + \underline{\mu} = \underline{L} \underline{D}^{-1/2} \underline{t} + \underline{\mu}$$

takes the initial multiple integral and completely factorizes it.

$\Rightarrow \underline{X} = \underline{L} \underline{D}^{-1/2} \underline{T} + \underline{\mu}$  is a transformation which takes the ~~multivariate~~ multivariate Normal r.v.  $\underline{X}$  to m.r.v.  $\underline{T}$  which is nothing but a multivariate ~~std~~ Normal r.v.!



Continuing the discussion in the prev. lecture, we saw that the following series of change of variables has done the job of factoring the multiple integral into product of 1-d integrals:

$$\underline{x} = \underline{y} + \underline{\mu} = \underline{Lz} + \underline{\mu} = \underline{LD}^{\frac{1}{2}}\underline{t} + \underline{\mu}$$

Now look at the corresponding ~~transformed~~ transformed r.v.s:

$$\underline{X} = \underline{Y} + \underline{\mu} = \underline{LZ} + \underline{\mu} = \underline{LD}^{\frac{1}{2}}\underline{T} + \underline{\mu}$$

Let us introduce our notation: " $\underline{X} \sim N(\underline{\mu}, \underline{\Sigma})$ " means that  $\underline{X}$  is a multivariate Normal r.v. with parameters  $\underline{\mu}, \underline{\Sigma}$ .

It is easy to see the following from the derivation in the prev. lecture:

$$\underline{Y} \sim N(\underline{0}, \underline{\Sigma}), \quad \underline{Z} \sim N(\underline{0}, \underline{D}), \quad \underline{T} \sim N(\underline{0}, \underline{I})$$

Here  $\underline{X}$  is nothing but an ~~affine~~ affine transformation (which simply rotates & translates axis) of the r.v.  $\underline{T}$   $\longrightarrow$

$d$  in other words  $\underline{T}$  is a multivariate std Normal r.v.!

$d$  in other words it is collection of  $n$  independent r.v.  $T_1, T_2, \dots, T_n$  & all of which are std. Normal.

~~Here we can show that also~~

$$E[\underline{X}] = \underline{LD}^{\frac{1}{2}} E[\underline{T}] + \underline{\mu} = \underline{\mu}$$

$$\text{Cov}(\underline{X}) = \underline{LD}^{\frac{1}{2}} \text{Cov}(\underline{T}) \underline{D}^{\frac{1}{2}} \underline{L}^T = \underline{LDL}^T = \underline{\Sigma}$$

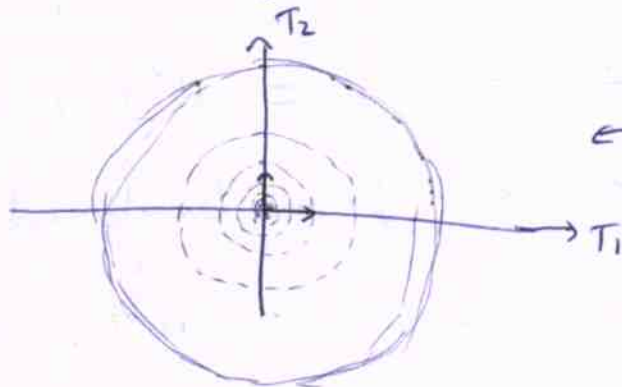
$\downarrow$   
 $\underline{I}$

~~Why we can show that~~ Hence the parameters  $\underline{\mu}, \underline{\Sigma}$  are nothing but the mean vector & cov. matrix!

In other words, given the pdf expansion one can just read out the mean vector & covariance matrix for a (multivariate) Normal r.v.

Here how the data generated from each of r.v.'s  $T, Z, Y, X$  would look like (the eigen vectors & values are also indicated in each case):  
(2-d example)

$$T \sim N(\underline{0}, \underline{I})$$



(imagine a bell curve representing the prob. density sitting on this)

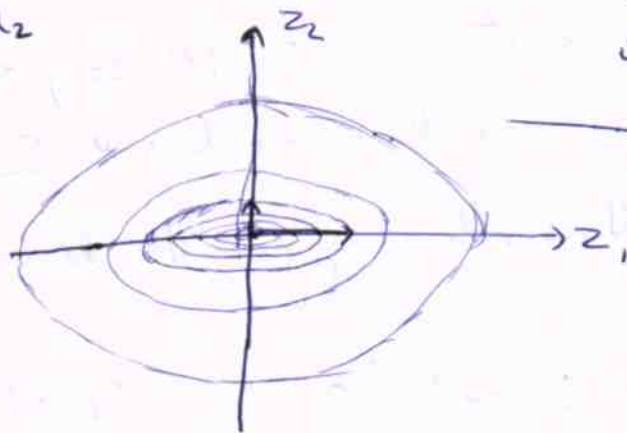
Here spacing between concentric circles is prop. to ~~prob.~~ prob. density.

$$Z \sim N(\underline{0}, \underline{0})$$

$$\begin{bmatrix} d_1 & 0 \\ 0 & d_2 \end{bmatrix}$$

let  $d_1 > d_2$

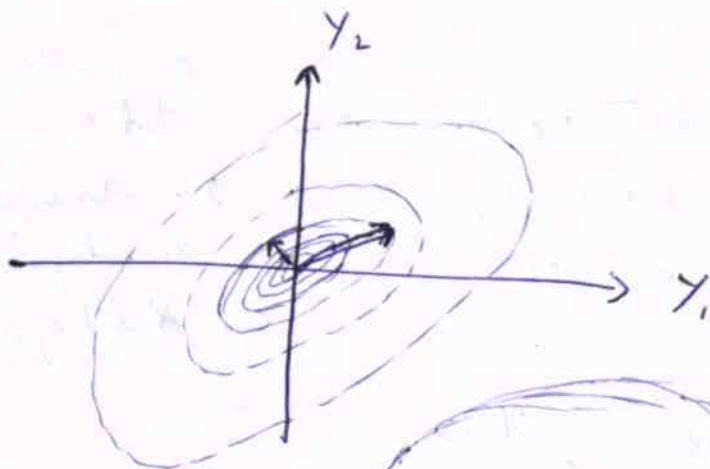
(in other words  $Z = DT$   
 $z_1 = d_1 T_1$   
 $z_2 = d_2 T_2$ )  
(transformation representing scaling of axis)



imagine a ~~skewed~~ skewed bell curve sitting on this plot representing prob. density.

$$Y \sim N(\underline{0}, \underline{\Sigma})$$

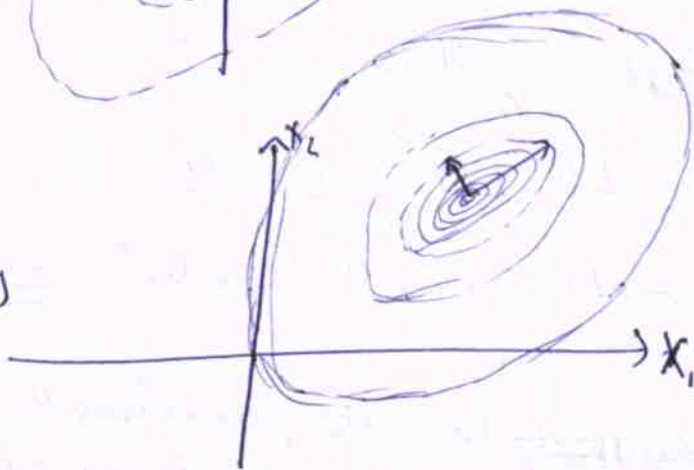
$Y = LZ$   
(transformation representing rotation of axis)



A point with its rating is the ~~needed~~ ellipse will be || to axis if cov. matrix is diagonal

$$X \sim N(\underline{\mu}, \underline{\Sigma})$$

$X = Y + \underline{\mu}$   
(transformation representing shifting of axis)





In the above figures the ellipses/circles are nothing but the cross-sections of the bell-shaped pdf at equal intervals (heights).

In words, any multivariate <sup>Normal</sup> r.v. can be produced by shifting, rotating and scaling ~~of~~ a multi-std. Normal r.v. (which is nothing but a collection of independent std. Normals).

→ What is mgf of  $X \sim N(\underline{\mu}, \underline{\Sigma})$ ?

$$\underline{\text{Ans}} \quad M_X(\underline{\Delta}) = E[e^{\underline{\Delta}^T X}] = \frac{1}{(\sqrt{2\pi})^n |\underline{\Sigma}|^{1/2}} \int \dots \int e^{\underline{\Delta}^T \underline{x}} e^{-\frac{1}{2}(\underline{x}-\underline{\mu})^T \underline{\Sigma}^{-1}(\underline{x}-\underline{\mu})} d\underline{x}$$

Again do the change of variables:  $\underline{x} = \underline{L} \underline{D}^{1/2} \underline{t} + \underline{\mu}$ , we get:

$$M_X(\underline{\Delta}) = \frac{1}{(\sqrt{2\pi})^n} \int \dots \int e^{\underline{\Delta}^T (\underline{L} \underline{D}^{1/2} \underline{t} + \underline{\mu})} e^{-\frac{1}{2} \underline{t}^T \underline{t}} d\underline{t}$$

$$= e^{\underline{\Delta}^T \underline{\mu}} \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{a_i t_i - \frac{1}{2} t_i^2} dt_i \quad \text{Put } \underline{\Delta}^T \underline{L} \underline{D}^{1/2} = \underline{a}^T$$

$$= e^{\underline{\Delta}^T \underline{\mu}} \prod_{i=1}^n e^{a_i^2 / 2} \quad \text{nothing but } M_{T_i}(a_i) \rightarrow T_i \text{ is std. Normal r.v.}$$

$$= e^{\underline{\Delta}^T \underline{\mu} + \frac{1}{2} \underline{a}^T \underline{a}} = e^{\underline{\Delta}^T \underline{\mu} + \frac{1}{2} \underline{\Delta}^T \underline{L} \underline{D} \underline{L}^T \underline{\Delta}} = e^{\underline{\Delta}^T \underline{\mu} + \frac{1}{2} \underline{\Delta}^T \underline{\Sigma} \underline{\Delta}}$$

$$\Rightarrow \boxed{M_X(\underline{\Delta}) = e^{\underline{\Delta}^T \underline{\mu} + \frac{1}{2} \underline{\Delta}^T \underline{\Sigma} \underline{\Delta}}}$$

This is also worth remembering.

Following things abt  $M_X(\Delta)$  are notable:

(i) It is again a Gaussian function! (exp. of a quadratic term)

(ii) It shows that all moments (and in fact the pdf) <sup>depend</sup> only on  $\underline{\mu}$ ,  $\underline{\Sigma}$  & in other words the entire distribution is completely specified by the first two moments!

(iii) From the joint mgf  $M_X(\Delta)$  we can easily <sup>write</sup> ~~put~~ down the mgf of  $X_i$ :

choose  $\Delta = \begin{bmatrix} 0 \\ \vdots \\ \delta_i \\ \vdots \\ 0 \end{bmatrix} \rightarrow a_i$  Using the mgf expression of  $M_X(\Delta)$  we have

$$M_X(\Delta) = E[e^{\Delta^T X}] = e^{\delta_i \mu_i + \frac{1}{2} \delta_i^2 \Sigma_{ii}}$$

$\downarrow$  also equal to  $M_{X_i}(\delta_i)$      
  $\downarrow$   $\delta_i \mu_i$  is constant of  $\underline{\mu}$      
  $\downarrow$   $\frac{1}{2} \delta_i^2 \Sigma_{ii}$  is  $i^{\text{th}}$  diagonal entry in  $\underline{\Sigma}$

This shows that  $X_i$  is ~~indeed~~ a Normal r.v. with mean  $\mu_i$  & variance  $\Sigma_{ii}$ .

In other words, if  $X_1, X_2, \dots, X_n$  are jointly Normal, then each of  $X_i$  is Normal.

→ Another property:

$X_1, X_2, \dots, X_n$  are jointly Normal  
 or  
 $X = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}$  is multivariate Normal

$\iff \underline{a}^T X$  is Normal r.v.  $\forall \underline{a} \neq 0$  (II)



Proof for  $\Rightarrow$  is given below. (Proof for  $\Leftarrow$  is an assignment)

Proof: We need to show that  $Y = \underline{a}^T X$  is Normal. The idea is to show that mgf of  $Y$  looks like that of a Normal r.v.:

$$M_Y(s) = E[e^{sY}] = E[e^{s \underline{a}^T X}] \quad (\text{put } s \underline{a}^T = \underline{c}^T)$$

$$= E[e^{\underline{c}^T X}] \rightarrow \text{mgf of } X!$$

$$= e^{\underline{c}^T \underline{\mu} + \frac{1}{2} \underline{c}^T \underline{\Sigma} \underline{c}}$$

$$= e^{\underline{a}^T \underline{\mu} + \frac{1}{2} \underline{a}^T \underline{\Sigma} \underline{a}}$$

$\downarrow$   
Indeed looks like mgf of Normal r.v. with mean as  $\underline{a}^T \underline{\mu}$  & variance as  $\underline{a}^T \underline{\Sigma} \underline{a}$ !

Here Proved.

~~One~~ One consequence of this is that if  $X_1, X_2, \dots, X_n$  are jointly Normal then  $X_1 + X_2 + \dots + X_n$  (or any partial sum) is a Normal r.v.!  
(Note: adding, convoluting Gaussian functions is again Gaussian!)

$\rightarrow$  Another property:

If  $X_{n \times 1}$  is a  $n$ -dim multivariate Normal r.v. &  $A_{m \times n}$  is a  $m \times n$  matrix such that all rows of  $A$  are linearly independent then  $Y_{m \times 1} = A X_{n \times 1}$  is also a multivariate Normal r.v.!

Let the rows of  $A$  be  $\underline{a}_i^T$  ( $i=1$  to  $m$ ) then:

$$Y_i = \underline{a}_i^T X$$

we know that each of  $\underline{a}_i \neq 0$  (since ~~the~~ <sup>rows</sup> are lin. ind.)

By  $\textcircled{\text{II}}$  we get that each of  $Y_i$  are Normal r.v.s.  
( $\Rightarrow$  part)

Now ~~we~~ we want to show that all (non-trivial) linear combinations of  $Y_i$  are Normal then by  $\textcircled{\text{II}}$  we again can say that  $Y_1, Y_2, \dots, Y_m$  are jointly Normal.  
( $\Leftarrow$  part)

lets show this:

$$\text{Consider } \underline{c}^T Y = \sum_{i=1}^m c_i Y_i = \left( \sum_{i=1}^m c_i \underline{a}_i^T \right) X = \underline{c}^T X$$

$\downarrow$  if  $\underline{c}$  is arbitrary  $\underline{c} \neq 0$  then this is arbitrary lin. comb. of  $Y_i$ 's.  
 $\downarrow$  nothing but linear comb. of rows of  $A$  which is guaranteed to be a non-zero row vector  $\underline{c}^T$ . ( $\because$  rows are lin. ind.)  
 $\downarrow$  is indeed Normal by  $\textcircled{\text{II}}$  again.  $\Rightarrow$

Here Proved.

Now lets look at the case when  $X_1, X_2, \dots, X_n$  are jointly Normal. We have already seen that if the covariance matrix of  $X = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}$  is diagonal then the pdf factorizes and hence each of  $X_1, \dots, X_n$  are independent Normal r.v.s.



However covariance matrix is diagonal  $\Leftrightarrow$  all covariance terms  $\text{cov}(X_i, X_j) = 0$  ( $i \neq j$ ). In other words for Normal rvs  $X_1, X_2, \dots, X_n$  are independent  $\Leftrightarrow X_1, X_2, \dots, X_n$  are uncorrelated.

→ To summarize:

For  $X = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}$ , a multivariate Normal r.v. with parameters  $\underline{\mu}, \underline{\Sigma}$  we have  
(or in other words  $X_1, \dots, X_n$  are jointly Normal)

- ①  $\underline{\mu}$  is mean vector of  $X$  &  $\underline{\Sigma}$  is covariance matrix of  $X$ .
- ② Each of  $X_1, \dots, X_n$  are individually Normal rvs with mean  $E[X_i] = \mu_i$ ,  $\text{var}(X_i) = \Sigma_{ii}$
- ③  $X_1, X_2, \dots, X_n$  are independent  $\Leftrightarrow X_1, X_2, \dots, X_n$  are uncorrelated (i.e.  $\underline{\Sigma}$  is diagonal)
- ④ There exists an affine transformation (i.e.  $\exists A_{n \times n}$  such that):  

$$X = AT + \underline{\mu}$$
 where  $T$  is a multivariate std. Normal r.v. (i.e.  $T$  is collection of  $n$  std. Normals)  
 In fact  $A$  is nothing but  $LD^{1/2}$  where  $\underline{\Sigma} = LD^2L^T$  ( $L$  is orthogonal,  $D$  is diagonal)  
 (In other words, data from  $X$  can be generated by rotating, scaling, and shifting data generated by  $T$ )
- ⑤  $X_1, X_2, \dots, X_n$  are jointly Normal  $\Leftrightarrow \underline{a}^T X$  is Normal  $\forall \underline{a} \neq 0$
- ⑥  ~~$X$~~   $X$  is multivariate Normal  $\Rightarrow Y = AX$  is also multivariate Normal for  $A$  full row ranked  $A$ .

## Intuitive meaning of Conditional Expectation

In one of prev. lectures we saw that  $\operatorname{argmin}_c E[(X-c)^2] = E[X]$   
 i.e.  $E[X]$  is the "best" constant value that approximates  $X$ .

Now we will show that the "best" approximation of a r.v.  $Y$  as a function of another r.v.  $X$  is nothing but  $E\{Y/X\}$

(recall that this is indeed a r.v. and a function of  $X$ )

i.e. TST  $\operatorname{argmin}_{g(x)} E[(Y-g(x))^2] = E\{Y/X\}$

Proof  $E[(Y-g(x))^2] = E[(Y-E\{Y/X\}) + (E\{Y/X\}-g(x))^2]$   
 (add and subtract  $E\{Y/X\}$ )  
 $= E[(Y-E\{Y/X\})^2] + E[(E\{Y/X\}-g(x))^2]$

$+ 2 E[(Y-E\{Y/X\})(E\{Y/X\}-g(x))]$   
 (functions of  $X, Y$  let it be  $f(x, Y)$  only a function of  $X$ , let it be  $h(x)$ )

Let's write this term down as follows

~~$E[(Y-E\{Y/X\})(E\{Y/X\}-g(x))]$~~

Now this expectation is:

$$\begin{aligned} E[f(x, Y) h(x)] &= E[E[f(x, Y) h(x) / X]] \\ &= E[h(x) E[f(x, Y) / X]] \\ &= E[h(x) \underbrace{E[(Y-E\{Y/X\}) / X]}_{\text{is zero}}] = 0 \end{aligned}$$

$\therefore E[(Y-g(x))^2] = E[(Y-E\{Y/X\})^2] + E[(E\{Y/X\}-g(x))^2]$   
 $\geq E[(Y-E\{Y/X\})^2]$



In the above inequality, an equality is achieved  
iff  $E\{Y/x\} = g(x)$ .

Hence always  $E\{(Y - g(x))^2\} = E\{Y/x\}$ . Here Proved.

Sequences of RVs

eg Let us recall an example from one of the prev. lectures.

Consider  $M_n = \max(X_1, X_2, \dots, X_n)$ , where  $X_1, X_2, \dots, X_n$  are iid r.v.

We showed that  $F_{M_n}(x) = (F(x))^n \quad \forall x$  suppose  $F_{X_i} = F, f_{X_i} = f$   
(denote their common dist. by  $F$ )

~~Previously we considered the special case  $F(x) = \dots$~~

Now let each of  $X_i$  be upper bounded by some "a". i.e.

$$F(x) = \begin{cases} \text{non-monotonic } f(x) < 1 & \text{for } x < a \\ 1 & \text{for } x \geq a \end{cases} \quad \textcircled{\text{II}}$$

(Recall that in prev. lecture we restricted  $X_i$  to be a uniform r.v.) (Here we do not.)

Now, consider this sequence of r.v.s:

$$M_1, M_2, \dots, M_n, \dots$$

Intuitively,  $\lim_{n \rightarrow \infty} M_n$  (the limiting r.v.) will be the degenerate discrete r.v. taking value "a" with prob. 1. This is because if we consider "infinite" iid r.v.s  $X_i$ , somewhere we must have hit "a".

~~Here~~ Now let's look at

$$\lim_{n \rightarrow \infty} F_{M_n}(x) \quad \forall x$$

$$\hookrightarrow = \begin{cases} 0 & x < a \\ 1 & x \geq a \end{cases} \quad \text{(by } \textcircled{\text{I}} \text{ and } \textcircled{\text{II}})$$

(At "a" the limit doesn't exist. Let's choose it to be 1.)



Note that this limiting dist. (i.e.  $\lim_{n \rightarrow \infty} F_{X_n}(x)$ ) is indeed a distribution and in fact it corresponds to that of the degenerate discrete r.v. taking value "a" with prob. 1!

This motivates for the following defn. of convergence of r.v.s:

### Convergence in Law

The sequence of r.v.s  $\{X_n\}$  is said to converge to a r.v.  $X$  "in law" or "in distributions" (denoted by  $\{X_n\} \xrightarrow{d} X$  or  $\{X_n\} \xrightarrow{D} X$ ) if

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x) \quad (\forall x \text{ where } F_X \text{ is cont.})$$

\* Note that even in prev. eg,  $x=a$  was point of discontinuity of  $F_X$  and in fact the  $\lim_{n \rightarrow \infty} F_{X_n}(a)$  doesn't exist. We can choose  $F_X(a)$  conveniently to make  $F_X$  a right cont. (i.e. a valid distribution fuc.)

\* The issue of convergence of r.v.s is posed as a problem of point-wise convergence of dist. functions! (which are very families with). Hence this form of convergence is simplest to understand and is useful in practice.

\* If we consider arbitrary sequences of r.v.s, then their dist. functions may not converge (point-wise) or even if they do, the limiting function may not be a dist. function!

eg. Consider  $X_n \sim U[0, n]$  (Intuitively  $\lim_{n \rightarrow \infty} X_n$  should not exist and this is what  $\xrightarrow{d}$  shows)

It is easy to see that  $\lim_{n \rightarrow \infty} F_{X_n}(x) = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases}$  which is not a valid dist. func!

\* In the Mn eg.  $\{M_n\} \xrightarrow{d} a$  this is a degenerate discrete r.v.  
 $\downarrow$   
 these are cont. r.v.

Hence sequences of cont. r.v. may  $\xrightarrow{d}$  to a discrete r.v. etc. So in general it may not make sense to talk abt convergence in pdf/pdf etc.

\* Though useful, the notation of  $\xrightarrow{d}$  is ~~rather loose~~ <sup>loose</sup> in the sense that values the r.v.s take cannot be determined by the dist. func. alone (~~if~~ if  $X \sim U[0,1]$  then  $1-X$  is also  $\sim U[0,1]$ !)

Here we can easily generate toy eg. where the values ~~take~~ take are diff. but dist. func. is same. So we define convergence notion that some how says  $\{X_n\}$  as  $n \rightarrow \infty$  is very close to  $X$ .

### Convergence in Probability

A sequence of r.v.s  $\{X_n\}$  is said to converge to  $X$  in probability (denoted by  $\{X_n\} \xrightarrow{p} X$ ) if  $\lim_{n \rightarrow \infty} P[|X_n - X| > \epsilon] = 0 \quad \forall \epsilon > 0$ .

\* We can equivalently write down the  $\delta$ -neighbourhood version of this defn: Given  $\delta > 0$ ,  $\exists N \in \mathbb{N} \ni \forall n \geq N \quad P[|X_n - X| > \epsilon] < \delta$ . (and this happens for  $\forall \epsilon > 0$ )

\* Taking some eg of  $M_n$ , let us see if  $\lim_{n \rightarrow \infty} P[|M_n - a| > \epsilon] = 0$ ?

$$M_n \leq a \Rightarrow P[|M_n - a| > \epsilon] = P[a - M_n > \epsilon] = P[M_n < a - \epsilon]$$

$$= \text{some constant} (g(a - \epsilon))^n$$

( $g$  is some monotonic func.  $< 1$ )

$$\Rightarrow \text{indeed } \lim_{n \rightarrow \infty} P[|M_n - a| > \epsilon] = \lim_{n \rightarrow \infty} (g(a - \epsilon))^n = 0!$$

$$\text{Here } \underline{\underline{\{M_n\} \xrightarrow{p} a}}$$

\* In fact, one can show that: (not need out of scope)

$$\{X_n\} \xrightarrow{p} X \Rightarrow \{X_n\} \xrightarrow{d} X$$



Here convergence in prob. is strictly stronger notion of convergence than <sub>in</sub> dist. convergence



Now we have seen a notion of convergence which only looks at dist. functions (i.e. probabilities) and a notion which actually looks at values r.v.'s take. Now let's look at a notion of convergence of r.v.'s which uses notion of moments (i.e. convergence in averages!)

### Convergence in $n^{\text{th}}$ Moment

A sequence of r.v.'s  $\{X_n\}$  is said to converge to  $X$  in the  $n^{\text{th}}$  moment (denoted by  $\{X_n\} \xrightarrow{n} X$ ) if  $\lim_{n \rightarrow \infty} E[|X_n - X|^n] = 0$

\* Let's again work out the  $\{M_n\}$  example:

$$\text{we have } F_{M_n}(x) = (F(x))^n \Rightarrow f_{M_n}(x) = n(F(x))^{n-1} f(x)$$

Let's take the case  $F, f$  represent uniform r.v. between  $[0, a]$ . non-zero for  $x \leq a$

$$\Rightarrow F_{M_n}(x) = \begin{cases} 0 & x < 0 \\ (x/a)^n & 0 \leq x < a \\ 1 & x \geq a \end{cases} \Rightarrow f_{M_n}(x) = \begin{cases} \frac{n x^{n-1}}{a^n} & 0 \leq x < a \\ 0 & \text{otherwise} \end{cases}$$

Also consider  $n=2$  (convergence in mean squared sense)

$$\begin{aligned} E[|X_n - a|^2] &= \int_0^a (x-a)^2 \frac{n x^{n-1}}{a^n} dx = \frac{n}{a^n} \left[ \int_0^a (x-a)^2 x^{n-1} dx \right] \\ &= \frac{n}{a^n} \left[ \frac{a^{n+2}}{n+2} - \frac{2a^{n+2}}{n+1} + \frac{a^{n+2}}{n} \right] \\ &= \frac{2a^2}{(n+1)(n+2)} \end{aligned}$$

$$\Rightarrow \{M_n\} \xrightarrow{n=2} a$$

Now what is the relation between  $\{X_n\} \xrightarrow{\delta_1} X$  and  $\{X_n\} \xrightarrow{\delta_2} X$  ??

(see assign. problem)

\* It turns out that:

$$\{X_n\} \xrightarrow{n} X \implies \{X_n\} \xrightarrow{p} X \quad \left( \begin{array}{l} \text{and hence} \\ \implies \{X_n\} \xrightarrow{d} X \end{array} \right)$$

(for some  $n$ )  
~~(however!)  
 convergence is not true~~

Proof of  $\implies$  is by what is known as Markov inequality which we will see shortly. Here is counter eg for the converse:

Consider  $P\{X_n = n\} = \frac{1}{n}$   
 $P\{X_n = 0\} = 1 - \frac{1}{n}$

It is easy to see that  $\{X_n\} \xrightarrow{p} 0$ . Here is why:

$$P\{|X_n - 0| \geq \epsilon\} = P\{X_n = n\} = \frac{1}{n} \implies \lim_{n \rightarrow \infty} P\{|X_n - 0| \geq \epsilon\} = 0$$

However  $E\{|X_n - 0|^n\} = \frac{1}{n} n^n = n^{n-1}$  (of course  $n \geq 1$ )

$$\implies \lim_{n \rightarrow \infty} E\{|X_n - 0|^n\} \neq 0 \implies \{X_n\} \not\xrightarrow{d} 0.$$

### Markov Inequality

Let  $Y = g(X)$  be a non-negative r.v. Then

$$P\{g(X) > \epsilon\} \leq \frac{E\{g(X)\}}{\epsilon} \quad \forall \epsilon > 0$$

very imp. inequality.

(neg.  $X$  in cont.)  
 $\downarrow$

Proof:  $E\{g(X)\} = \int_{-\infty}^{\infty} g(u) f_X(u) du = \int_{X: g(X) \leq \epsilon} g(u) f_X(u) du + \int_{X: g(X) > \epsilon} g(u) f_X(u) du$

$\implies \int_{X: g(X) > \epsilon} g(u) f_X(u) du$   
 $(\because g(u) \geq 0 \text{ this is OK})$



$$\Rightarrow E\{g(x)\} \geq \int_{x: g(x) > \epsilon} \epsilon f_x(x) dx$$

$$= \epsilon P[g(x) > \epsilon]$$

$$\Rightarrow P[g(x) > \epsilon] \leq \frac{E[g(x)]}{\epsilon} \quad \text{Here Proved}$$

Now put  $g(x) = |x|^n$ ,  $\epsilon = \delta^n$

~~$$P[|x| > \delta] = P[|x|^n > \delta^n]$$~~

$$P[|x| > \delta] = P[|x|^n > \delta^n] \leq \frac{E[|x|^n]}{\delta^n}$$

(inequality relating absolute moments to the corresponding prob.)

Now put  $g(x) = |X_n - X|$ ,  $\epsilon = \delta^n$  we get:

$$P[|X_n - X| > \delta] \leq \frac{E[|X_n - X|^n]}{\delta^n}$$

If  $|X_n - X| \xrightarrow{n} 0$ , then  $\text{RHS} \rightarrow 0 \Rightarrow \lim_{n \rightarrow \infty} \text{LHS} = 0 \Rightarrow |X_n| \xrightarrow{P} X$

Now put  $g(x) = (x - E[x])^2$ ,  $\epsilon = \delta^2$ , we get:

$$P[|x - E[x]| > \delta] = P[(x - E[x])^2 > \delta^2] \leq \frac{E[(x - E[x])^2]}{\delta^2} = \frac{\text{var}(x)}{\delta^2}$$

$$\Rightarrow P[|x - E[x]| > \delta] \leq \frac{\text{var}(x)}{\delta^2} \quad (\text{very imp. Chebyshev's inequality})$$

Now put  $\delta = k \sqrt{\text{var}(X)}$ , then:

$$P[|X - E\{X\}| > k \sqrt{\text{var}(X)}] \leq \frac{1}{k^2}$$

↓  
(independent of  $X$ !)

Says that prob. that any r.v. deviates from its mean by  $k$  times std. dev., it's always  $\leq \frac{1}{k^2}$ .

Though it looks like we have ~~done~~ <sup>reflected</sup> ~~some~~ some terms to arrive at the inequality, it happens that ~~we~~ "given only two moments of a r.v., one cannot beat/improve ~~the~~ what the Chebyshev's inequality gives!

In other words I can come with distributions satisfying the inequality exactly (~~with~~ <sup>i.e. with</sup> equality).



In the previous lecture we saw three modes of convergence of r.v.s:

- (i)  $\{X_n\} \xrightarrow{d} X$  iff  $\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$  ( $\forall x$  where  $F_X$  is conts.)  
(Convergence in Distribution)  
 $\uparrow \downarrow$
- (ii)  $\{X_n\} \xrightarrow{p} X$  iff  $\lim_{n \rightarrow \infty} P[|X_n - X| > \epsilon] = 0 \quad \forall \epsilon$   
(Convergence in probability)  
 $\uparrow \downarrow$
- (iii)  $\{X_n\} \xrightarrow{m} X$  iff  $\lim_{n \rightarrow \infty} E[|X_n - X|^m] = 0$  ( $m$  is given,  $m \in \mathbb{N}$ )  
(Convergence in  $m^{\text{th}}$  moment)

Convergence in 1<sup>st</sup> moment is also known as convergence in mean  
 " " 2<sup>nd</sup> " " " " " " " " mean square

Convergence in distribution, though is the weakest form, is also in some sense the most useful form of convergence. One reason is this: Usually we need to guess what "X" is going to be. However the notion of con. in dist. gives us an easy way out. Given the marginal distribution functions  $F_{X_n}$  we can always see if the pointwise limit of these functions exists and if so is it a valid dist. function. This immediately gives  $F_X$ ! So this is an easy way to guess what "X" is.

On the other hand, unless X is known, the other notions of convergence are hard to verify. In fact, unless the X is a constant number, to verify (ii), (iii) it self is not straight forward because we will need the joint dist. of  $X_n$  & X!

(Note that  $X_n, X$  need to be r.v.s defined on <sup>the</sup> same probability space  $\mathcal{P} = (\mathcal{X}, \mathcal{F}, P)$ )



However in cases where  $X$  turns out to be a constant,  
in many cases most of the times, (iii) is the easiest to verify  
(we will see eq. of Law of large no. later)

ans to the  
So the question "which is the opt mode of convergence" depends  
on the need and convenience of the problem in hand.  
wrt.

Now in case mgf's of  $X_n$  &  $X$  all exist, then it is easy to  
see that convergence of mgf of  $X_n$  to mgf of  $X$  indeed  
(pointwise)

implies that  $\{X_n\} \xrightarrow{D} X$ . This is because, if mgf exists,  
mgf and dist. fun. uniquely determine each other.

III by if pdf of  $X_n \xrightarrow{\text{pointwise converges}} \text{pdf of } X \Rightarrow \{X_n\} \xrightarrow{D} X$

III by pmf of  $X_n \xrightarrow{\text{pointwise converges}} \text{pmf of } X \Rightarrow \{X_n\} \xrightarrow{D} X$

In summary, convergence (pointwise) of mgf's, pdf's or pmf's  
all imply convergence in distribution. However the converse  
is not true the reasons are simple:

(i) mgf may not always exist! So if  $\{X_n\} \xrightarrow{D} X$  we may not  
even have mgf existing so what is the question of their convergence?

(ii) We saw in prev. class that a sequence of conts. r.v can  
converge to a discrete r.v (in dist.) and in fact we can  
also show that sequence of discrete r.v converge (in dist.) to a  
conts. r.v!  
give examples where  
So  $\{X_n\} \xrightarrow{D} X \not\Rightarrow$  convergence in pdf/pmfs.

However if the mgf/pdf/pmfs exist for all  $X_n$  &  $X$  then convergence  
in dist and convergence in mgf/pdf/pmfs are equivalent.



All said and done, a rv is nothing but a function ~~also~~ from  $\Omega \rightarrow \mathbb{R}$ ! So we can as well talk abt notions of pointwise convergence of these functions:

### Convergence Everywhere of Sure Convergence

$$\left\{ X_n \right\} \xrightarrow{\text{sure}} X \quad \text{iff} \quad \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega) \quad \forall \omega \in \Omega$$

( $\{X_n\}$  is said to surely converge to  $X$ ) (pointwise convergence)

$$\left\{ X_n \right\} \xrightarrow{e} X$$

It is easy to see that this is the "strongest" notion of convergence, but rarely used in practice (because of difficulty in verification)

Now we can also relax above condition and say that the  $\omega$ 's at which pointwise convergence does not happen must be ~~not~~ negligible (i.e. prob. of those set of outcomes is zero). This leads to the following relaxed notion of convergence:

### Almost sure convergence (convergence with prob. 1)

or almost everywhere convergence

$$\left\{ X_n \right\} \xrightarrow{\text{a.s.}} X \quad \text{iff} \quad P(\{\omega \in \Omega \mid X_n(\omega) \not\rightarrow X(\omega) \text{ as } n \rightarrow \infty\}) = 0$$

( $\{X_n\}$  is said to almost surely converge to  $X$ ) or  $P(\{\omega \in \Omega \mid X_n(\omega) \rightarrow X(\omega) \text{ as } n \rightarrow \infty\}) = 1$

It turns out that there are convenient ways of verifying a.s. convergence but we will not venture into them. However, note that, this notion of convergence is strong and is also not very difficult to verify in many cases.

Finally we have following comparison of modes of convergence:



$$\{X_n\} \xrightarrow{a} X \Rightarrow \{X_n\} \xrightarrow{a.o.} X \Rightarrow \{X_n\} \xrightarrow{p} X \Rightarrow \{X_n\} \xrightarrow{d} X$$

$$\{X_n\} \xrightarrow{p} X \Rightarrow \{X_n\} \xrightarrow{d} X$$

In general, none of the converse implications are true. However, if  $X$  is a constant one can show that convergence in prob. & in dist. are equivalent. Also, in general, ~~the~~ convergence in a.o. and  $p$  are not comparable. (one may not imply the other)

## Law of Large Numbers

Suppose you are conducting an experiment and say you have measured a certain quantity. Since experimental procedures are prone to errors, you repeat the expt. (independently) many times and take average of the readings. Then we justify saying that this average for large number of readings is close to the true value.

Let's justify this using the notion of convergence of r.v.s.

Let  $\mu$  be the true quantity to be estimated. Let  $X_1, X_2, \dots, X_n, \dots$  be observed values of this true quantity.

$\downarrow$   
each are i.i.d r.v.s. represent the true value + noise

In other words,  $E[X_i] = \mu$ ; however due to some noise  $X_i$  are different from  $E[X_i]$  i.e. Let the variances  $\text{var}(X_i) = \sigma^2$  be non-zero.

Our intuition says  $\lim_{n \rightarrow \infty} \bar{S}_n = \frac{\sum_{i=1}^n X_i}{n} = \mu$  (i.e.  $E[X_i]$ )

$\downarrow$   
his large

average of  $n$  readings

$$\text{Let } \sum_{i=1}^n X_i = S_n, \frac{S_n}{n} = \bar{S}_n$$



Consider the sequence of r.v.s:

$$\bar{S}_1, \bar{S}_2, \dots, \bar{S}_n, \dots$$

$$\text{we know } E\{\bar{S}_n\} = E\left\{\frac{\sum_{i=1}^n X_i}{n}\right\} = \frac{n\mu}{n} = \mu$$

$$\begin{aligned} \text{Hence } \text{var}(\bar{S}_n) &= E\left\{\left[\frac{\sum_{i=1}^n X_i}{n} - \mu\right]^2\right\} \\ &= E\left\{\left[\frac{\sum_{i=1}^n (X_i - \mu)}{n}\right]^2\right\} = \frac{\sum_{i=1}^n E\{X_i - \mu\}^2}{n^2} \\ &= \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n} \end{aligned}$$

Now, we claim that  $\{\bar{S}_n\} \xrightarrow{n \rightarrow \infty} \mu$ .

Proof

$$\lim_{n \rightarrow \infty} E\{|\bar{S}_n - \mu|^2\} = \lim_{n \rightarrow \infty} \text{var}(\bar{S}_n) = \lim_{n \rightarrow \infty} \frac{\sigma^2}{n} = 0$$

$\therefore \{\bar{S}_n\}$  converges to  $\mu$  in mean squared sense and hence in probability i.e.  $\{\bar{S}_n\} \xrightarrow{p} \mu$

This statement is known as the "Weak law of large numbers".

Here is also the "strong" version of this:  $\{\bar{S}_n\} \xrightarrow{\text{a.s.}} \mu$   
(we won't prove this here)

→ The law of large numbers also helps us to link our axiomatic defn. of probability to classical notions:

Let  $X_i$  be indicator r.v. of an event  $E$ . ( $\forall i \in \mathbb{N}$ )  
 $\Rightarrow E\{X_i\} = P(E)$ .

Now  $\bar{S}_n \xrightarrow{\text{a.s.}} E\{X_i\} = P(E)$ . So if we want to estimate prob. of any event in a hard expt. then we can repeat the expt. (independently & identically) a large no. of times & avg. no. times event occurs is  $P(E)$ !

(Intro)

Recall the law of large numbers; it says that if  $X_1, X_2, \dots, X_n, \dots$  are a sequence of iid r.v.s with  $E\{X_i\} = \mu$ , then

$$\bar{S}_n \xrightarrow{\text{a.s.}} \mu \quad \left( \text{where } \bar{S}_n = \frac{\sum_{i=1}^n X_i}{n} \right).$$

Let  $X$  be a random quantity in a random experiment  $E$ . You can think about  $X_1, X_2, \dots, X_n, \dots$  as being the values ~~the~~ of taken by this random quantity in independent trials of the same random experiment  $E$ . Once this view is clear, the power of law of large numbers becomes evident:  $\bar{S}_n$  is a quantity easily computable ~~but~~ by just finding mean of the values the random quantity takes. What law of large no. says is that  $\bar{S}_n \xrightarrow{\text{a.s.}} \mu$ . Note that,  $\mu$ , on the other hand is a quantity which cannot be observed through repeats of the random exp.  $E$ ! In fact knowing  $\mu$  amounts to knowing partial information regarding the distribution of  $X$  whereas  $\bar{S}_n$  is simply a quantity easily computable in terms of observed values.

~~The law of large no. can also be employed to estimate quantity other than~~ This discussion (and that in prev. class) clearly show that this result can be employed to estimate mean of an unknown distribution.

In fact, we can try estimating other quantities of parameter of an unknown distribution by appropriately (or cleverly) choosing the r.v.s  $X$  such that  $E\{X\} = 0$ .

Here are some examples:



Suppose we are given a coin and asked to estimate the prob. of getting heads with it:

Let  $X$  be the (Bernoulli r.v.) indication of ~~heads~~ obtaining heads with the coin.  $P\{X=1\} = p \rightarrow$  to be determined.

$$P\{X=0\} = 1-p. \quad \text{Also } E\{X\} = p.$$

Now consider <sup>(independent)</sup> ~~repeats~~ repeats of coin tosses. Suppose  $X_1, X_2, \dots, X_n, \dots$  represent the indicators of observing heads.  $\{X_i\}$  are iid and have the <sup>same</sup> distribution as  $X$ .

So applying law of large no. we get:

$$\bar{S}_n \xrightarrow{\text{a.s.}} E\{X\} = p.$$

By defining iid Bernoulli r.v.s in a similar fashion allows us to ~~Now regard Bernoulli~~ ~~consider~~ estimate  $P(E)$  in any ~~of~~ random expt.   
  $\underbrace{\text{prob. of some event.}}$

$\rightarrow$  Another ex.

Suppose the true mean  <sup>$E\{X\} = \mu$</sup>  of a r.v.  $X$  is known. However  $\text{var}(X)$  is to be estimated. Again from school days knowledge we would say: ~~sample iid~~ estimate of  $\text{var}(X) = \frac{\sum_{i=1}^n (X_i - \mu)^2}{n}$ . ( $X_i$  are iid)

How is this justified?

Consider r.v.  $Y = (X - \mu)^2$  then  $E\{Y\} = \text{var}(X)$ .

Hence  $\bar{S}_n = \frac{\sum_{i=1}^n Y_i}{n} \xrightarrow{\text{a.s.}} E\{Y\} = \text{var}(X)$  (by law of large no.).   
  $\frac{\sum_{i=1}^n (X_i - \mu)^2}{n}$  (again iid)

In general, one may want to estimate a parameter  $\theta$  in an unknown distribution (of course some partial information regarding the unknown dist. may be available for eg. in the prev. situation we know the true mean of the dist.)

The idea is to come up with  $\hat{\theta}_n$  ( $\rightarrow$  r.v. which is a function of

Note that  $\hat{\theta}_n$  can be computed ~~using~~ <sup>using</sup> the observable values whereas  $\theta$  is something abstract.) such that  $\hat{\theta}_n \rightarrow \theta$ ,  
 $\leftarrow n \text{ iid samples from the unknown distribution.}$   
 $\downarrow$   
 (Convergence in some sense)

Such a r.v.  $\hat{\theta}_n$  is known as an estimator of  $\theta$ .

The following are two desirable properties of ~~any~~ <sup>an</sup> estimator:

(i) Unbiased estimator:  $E[\hat{\theta}_n] = \theta \quad \forall n$

(ii) Minimum variance Unbiased estimator  $\hat{\theta}_n^*$ :  $\text{var}(\hat{\theta}_n^*) \leq \text{var}(\hat{\theta}_n) \quad \forall n$   
 $\leftarrow$  optimal min var., unbiased estimator  
 $\leftarrow$  any unbiased estimator

(i) condition says that with any no. samples expected value of the estimator is the true value of the estimated quantity.

(~~help~~ Note that the estimators for mean,  $P(E)$ , variance discussed above are all unbiased estimators. Refer to assign. prob. for eg. of an estimator which is not unbiased).

(ii) condition says that at with any no. samples, the deviation ~~than~~ from true value  $\theta$  is lesser than with any other unbiased estimator.  
 This is the starting of Estimation Theory.



# Central limit theorem

Theorem: Let  $X_1, X_2, \dots, X_n, \dots$  be a sequence of iid r.v.s.

Then  $\tilde{S}_n \xrightarrow{D} N(0,1)$  (~~std. Normal~~) where  $\tilde{S}_n = \frac{S_n - E[S_n]}{\sqrt{\text{var}(S_n)}}$ .

(recall that  $S_n = \sum_{i=1}^n X_i$ )

→

Note that  $E[\tilde{S}_n] = 0$ ,  $\text{var}(\tilde{S}_n) = 1$  and is an affine transformation. In other words  $\tilde{S}_n$  is nothing but sum of  $n$  iid r.v.s normalized to zero mean & unit variance.

(In fact for any  $X$ ,  $Y = \frac{X - E[X]}{\sqrt{\text{var}(X)}}$  is the usual way of normalizing to zero mean & unit variance.)

Let  $E[X_i] = \mu$ ,  $\text{var}(X_i) = \sigma^2$ .

$$\begin{aligned}\tilde{S}_n &= \frac{S_n - E[S_n]}{\sqrt{\text{var}(S_n)}} = \frac{\sum_{i=1}^n X_i - n\mu}{\sqrt{n\sigma^2}} \\ &= \frac{\sum_{i=1}^n X_i - n\mu}{\sigma\sqrt{n}} \\ &= \frac{\sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma}\right)}{\sqrt{n}} = \frac{\sum_{i=1}^n Y_i}{\sqrt{n}}\end{aligned}$$

mean of variance

$Y_i$ 's are nothing but normalized versions of  $X_i$ . Also if  $X_i$  are iid then,  $Y_i$  are also iid.

Here the theorem above is proved if we show that:

$$\tilde{S}_n = \frac{\sum_{i=1}^n Y_i}{\sqrt{n}} \xrightarrow{D} N(0,1) \text{ for any } Y_1, Y_2, \dots, Y_n, \dots \text{ iid and mean 0, variance 1.}$$

Here is the proof:

Proof Suppose mgf of  $Y_i$  exists  $\rightarrow$  (In the generic proof we need not assume this. However, for the scope of this class, the assumption is OK).

Now,

$$M_{Y_i}(\lambda) = E[e^{\lambda Y_i}]$$

$$\Rightarrow \prod_{i=1}^n M_{Y_i}(\lambda) = \prod_{i=1}^n E[e^{\lambda Y_i}] = E\left[\prod_{i=1}^n e^{\lambda Y_i}\right] = E\left[e^{\lambda \sum_{i=1}^n Y_i}\right] \quad (\text{I})$$

$\underbrace{\qquad\qquad\qquad}_{(\because Y_i \text{ are independent})}$

( $M_{Y_i}(\lambda)$ )  
 $\swarrow$   
 $\because Y_i$  are identically distributed

Hence mgf of  $\frac{\sum_{i=1}^n Y_i}{\sqrt{n}}$  exists and is in fact  $E\left[e^{\frac{\lambda}{\sqrt{n}} \sum_{i=1}^n Y_i}\right] = \left(M_{Y_i}\left(\frac{\lambda}{\sqrt{n}}\right)\right)^n$

Our idea is to show that the sequence of mgf:

$$M_{S_1}(\lambda), M_{S_2}(\lambda), \dots, M_{S_n}(\lambda), \dots \xrightarrow{n \rightarrow \infty} M_X(\lambda) = e^{\lambda^2/2}$$

$\downarrow$   
X is a std Normal

$\hookrightarrow$  If we show this, then we are done.

in other words T.S.T.  $\lim_{n \rightarrow \infty} M_{S_n}(\lambda) = e^{\lambda^2/2} \quad (\forall \lambda)$

~~$\lim_{n \rightarrow \infty} M_{S_n}(\lambda) = \lim_{n \rightarrow \infty} \left(M_{Y_i}\left(\frac{\lambda}{\sqrt{n}}\right)\right)^n$~~  i.e. T.S.T.  $\lim_{n \rightarrow \infty} \left(M_{Y_i}\left(\frac{\lambda}{\sqrt{n}}\right)\right)^n = e^{\lambda^2/2} \quad (\forall \lambda)$

Let's expand  $M_{Y_i}(\lambda/\sqrt{n})$  using Maclaurin series:

$$M_{Y_i}\left(\frac{\lambda}{\sqrt{n}}\right) = 1 + \frac{(\lambda/\sqrt{n}) E\{Y_i\}}{1!} + \frac{(\lambda/\sqrt{n})^2 E\{Y_i^2\}}{2!} + \frac{(\lambda/\sqrt{n})^3 E\{Y_i^3\}}{3!} + \dots$$

$\downarrow$   $\downarrow$   $\downarrow$   
 $0$   $1$   $0\left(\left(\frac{\lambda}{\sqrt{n}}\right)^3\right)$

$$= 1 + \frac{\lambda^2}{2n} + o\left(\left(\frac{\lambda}{\sqrt{n}}\right)^3\right)$$

Now  $\lim_{n \rightarrow \infty} \left(1 + \frac{\lambda}{n}\right)^n = e^\lambda$  and  $\downarrow$  go faster to zero than  $\frac{1}{n}$ , we have that

$$\lim_{n \rightarrow \infty} \left(1 + \frac{\lambda^2}{2n} + o\left(\left(\frac{\lambda}{\sqrt{n}}\right)^3\right)\right)^n = e^{\lambda^2/2} \quad \text{Hence Proved.}$$



The key advantage of the Central Limit Theorem (CLT) is to approximate the distributions of  $S_n$  and  $\bar{S}_n$ !

$\downarrow$  num of iid rvs  $X_i$        $\downarrow$  mean of iid rvs  $X_i$

[Note that though the pmf/pdf of  $S_n$  is the convolution of pmf/pdf of  $X_i$  repeated  $n$  times, this convolution may not be easily computable. However the approximation provided by CLT for distribution is "easy" to compute and also is a good approximation for moderately large  $n$  (20-30 in practice!)]

CLT says that  $\tilde{S}_n \rightarrow N(0,1)$

$$\Rightarrow P[\tilde{S}_n \leq x] \approx \Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt.$$

$\underbrace{\hspace{10em}}_{\text{dist. fnc. of Std. Normal.}}$

$$\Rightarrow P\left[\frac{S_n - n\mu}{\sigma\sqrt{n}} \leq x\right] \approx \Phi(x)$$

$$P[S_n \leq x\sigma\sqrt{n} + n\mu] \approx \Phi(x)$$

$$P\left[\frac{S_n/n - \mu}{\sigma/\sqrt{n}} \leq x\right] \approx \Phi(x)$$

$$\Rightarrow P[S_n \leq y] \approx \Phi\left(\frac{y - n\mu}{\sigma\sqrt{n}}\right)$$

$$\Rightarrow P\left[\frac{\sqrt{n}(\bar{S}_n - \mu)}{\sigma} \leq x\right] \approx \Phi(x)$$

$$\Rightarrow P\left[\bar{S}_n \leq \frac{\sigma x}{\sqrt{n}} + \mu\right] \approx \Phi(x)$$

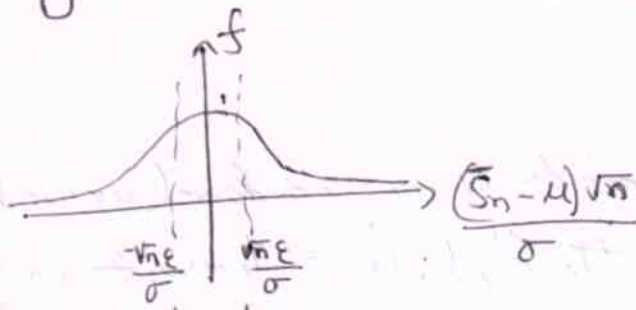
$$\Rightarrow \underline{\underline{F_{S_n}(y) \approx \Phi\left(\frac{y - n\mu}{\sigma\sqrt{n}}\right)}}$$

$$\Rightarrow \underline{\underline{F_{\bar{S}_n}(y) \approx \Phi\left(\frac{(y - \mu)\sqrt{n}}{\sigma}\right)}}$$

~~Dist~~  $S_n, \bar{S}_n$  can be approx. in terms of that of std. Normal r.v.!

$S_n, \bar{S}_n$  can be approx. in terms of that of std. Normal r.v.!

What this says alternatively (for large  $n$ ):



Area of this region  $\geq 1-\delta$

$1-\delta$  is usually known as the confidence

$\delta$  is " " " " the significance.

Many times we would wish that

$$|\bar{S}_n - \mu| < \epsilon$$

$\underbrace{\hspace{2cm}}_{\text{deviation from true mean}}$ 
 $\downarrow$  tolerance

$$\Rightarrow \frac{\sqrt{n}|\bar{S}_n - \mu|}{\sigma} < \frac{\sqrt{n}\epsilon}{\sigma}$$

but of course  $\bar{S}_n$  is r.v. so we can only satisfy this with high probability:

$$P\left[\frac{\sqrt{n}|\bar{S}_n - \mu|}{\sigma} < \frac{\sqrt{n}\epsilon}{\sigma}\right] \geq 1-\delta$$

$\delta$  is small  $\approx 0.05$

nothing but  $1 - 2\left(1 - \Phi\left(\frac{\sqrt{n}\epsilon}{\sigma}\right)\right)$

Here we want  $1 - 2\left(1 - \Phi\left(\frac{\sqrt{n}\epsilon}{\sigma}\right)\right) \geq 1-\delta$

$$\Rightarrow 1 - \Phi\left(\frac{\sqrt{n}\epsilon}{\sigma}\right) \leq \delta/2 \Rightarrow n \geq \left(\frac{\sigma \Phi^{-1}(1 - \delta/2)}{\epsilon}\right)^2 \quad \textcircled{II}$$

inequality is maintained since  $\Phi$  is monotonic

This shows that if we need to ensure that with high confidence ( $1-\delta=0.95$ ), the  $\bar{S}_n$  (sample mean) is within  $\epsilon$  ( $\approx 0.1$ ) tolerance of the true mean, then we need to have at least these many samples.

Here LLT can also give an idea about how many repeats of experiments are to be performed during estimation of a parameter



The same expression in (II) can be used to answer the following question:

Suppose  $\delta$  is fixed (confidence/significance level we work with is fixed). Also say we have repeated the <sup>ind.</sup> expt. for  $n$  times & we get  $\bar{S}_n$ . Now we can find critical value of  $\epsilon$  below which the value is acceptable as an approximation of  $\mu$ , i.e.

$$\epsilon \geq \frac{\sigma \Phi^{-1}(1-\delta/2)}{\sqrt{n}} \Rightarrow \epsilon_{\text{critical}} = \frac{\sigma \Phi^{-1}(1-\delta/2)}{\sqrt{n}}$$

$\therefore$  If  $|\bar{S}_n - \mu| < \epsilon_{\text{critical}}$  we accept  $\bar{S}_n$  as a valid approx. of  $\mu$  and otherwise reject it!

This notion leads to the theory of Hypothesis Testing.

RANDOM (OR) STOCHASTIC PROCESSES

Any collection of r.v.s (from the same prob. space  $\mathcal{P}$ ),  $\{X_t, t \in \mathcal{I}\}$ , (here  $\mathcal{I}$  is some countable or uncountable index set) is called as a stochastic or random process.

~~eg~~. The case of  $\mathcal{I} = \mathbb{N}$  gives back the case of sequences of r.v.s, ~~for~~ which we have been analysing since ~~couple~~ few lectures (in a limited sense).

Usually ~~if~~ the set  $\mathcal{I}$  has an interpretation of "time". Hence the name random process. Also till now we were concerned only with the case ~~of~~  $\mathcal{I} = \mathbb{N}$  and mostly looking at how the ~~data~~ "limiting" r.v. looks like. Now we are going to look at the collections of r.v.s in their entirety. i.e. ~~rather~~ specify all orders of joint-dist. functions:

i.e. specify  $f_{X_{t_1}, X_{t_2}, \dots, X_{t_k}}$   $\forall t_1, t_2, \dots, t_k$  and  $\forall k$ .

~~eg~~: The simplest eg. of stochastic process is Bernoulli process  
 $\rightarrow$  nothing but collection of independent Bernoulli r.v.s.

for instance: flipping of a coin is a Bernoulli process.

Note that it is easy to write down all orders of joint dist. func. in this case

$$f_{X_{t_1}, \dots, X_{t_k}} = f^k \quad \forall t_1, t_2, \dots, t_k \text{ and } k \text{ where } f \text{ is pmf of the Bernoulli p.v.}$$



Note that for ~~the~~ Bernoulli process, we can consider the "sub-process"  $\{X_t, t \in J, t \geq t_1\}$  where  $t_1$  is some  $t_1 \in J$ .

↳ This is again the same Bernoulli process  $\{X_t, t \in J\}$

↓  
in the sense that all orders of joint distributions are the same

In other words, it does not matter since when the process ~~originated~~ originated, the distribution (law) is the same.

Such processes are known as Stationary processes. We ~~will~~ formally define ~~the~~ stationary process as:

joint dist of  $X_{t_1}, \dots, X_{t_n}$  is same as that of  $X_{t_1+\Delta}, X_{t_2+\Delta}, \dots, X_{t_n+\Delta}$

and this happens for all  $\Delta$  and  $t_1, t_2, \dots, t_n$  and all  $n$

Some implications are:

(take  $n=1$ )  ~~$F_{X_{t_1}} = F_{X_{t_2}}$~~   $F_{X_{t_1}} = F_{X_{t_2}} \quad \forall t_1, t_2$ . (Marginal distributions are same)

(take  $n=2$ )  $F_{X_{t_1}, X_{t_2}} = F_{X_{t_1+\Delta}, X_{t_2+\Delta}}$  (Pairwise distributions only depend on  $|t_1 - t_2|$ )

... and so on.

→ classification of RPs (random processes)

<sup>(a)</sup> Depending on whether  $J$  is countable or uncountable:

If  $J$  is countably infinite → discrete time RP

$J$  is uncountable → continuous time RP

(b) Depending on whether the rv take on discrete or conts. values.

if  $X_t$  is discrete  $\therefore$  discrete-state r.p

$X_t$  is conts rv  $\therefore$  conts.-state r.p

Actually, the set of values the rvs take (in context of r.p) is known as state-space. Hence these names.

r.p.s with all  $2 \times 2 = 4$  combinations are possible.

One can have further classification in discrete-state: finite or infinite state space.

Two r.p we will be studying are (i) Markov chains  
(ii) Poisson process

(i) is a finite-state discrete time r.p.

(ii) is an infinite-state continuous time r.p.

Another r.p which is easy to describe is Gaussian Process (GP)

(analogous with defn. of Bernoulli process and relation with Bernoulli rv)

~~GP~~ ~~For now we are concerned with discrete time GP.~~

A r.p is called GP if any order joint dist is (multivariate) Gaussian distribution.

i.e.  $X_{t_1}, X_{t_2}, \dots, X_{t_k}$  are jointly Gaussian & jointly Normal.  
 $\forall t_1, t_2, \dots, t_k$  and  $k$

there is none as saying the matrix  
is pd!

$$\begin{bmatrix} \text{cov}(X_{t_1}, X_{t_1}) & \text{cov}(X_{t_1}, X_{t_2}) & \dots & \text{cov}(X_{t_1}, X_{t_k}) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(X_{t_k}, X_{t_1}) & \dots & \dots & \text{cov}(X_{t_k}, X_{t_k}) \end{bmatrix}$$

II



So in order to specify a GP we just need to specify

$E\{X_t\} \forall t$   $\xrightarrow{\text{known as}}$  mean function  
 &  $\text{cov}(X_{t_1}, X_{t_2}) \forall t_1, t_2$   $\xrightarrow{\text{known as}}$  auto covariance function.

(Imagine discrete and conts J)

In GP ~~terminology~~ notation, mean function is denoted by  $m_x(t) = E\{X_t\}$   
 auto covariance function by  $C_x(t_1, t_2) = \text{cov}(X_{t_1}, X_{t_2})$

We also know that  $\text{cov}(X_{t_1}, X_{t_2}) = E\{X_{t_1} X_{t_2}\} - E\{X_{t_1}\} E\{X_{t_2}\}$

↓  
~~called auto correlation~~  
 called auto correlation

auto correlation function is denoted by  $R_x(t_1, t_2) = E\{X_{t_1} X_{t_2}\}$ .

With this terminology we can ~~say~~ define GP with its mean function and auto covariance function (provided it is positive definite)  
 i.e. Condition **II**

We can ask when is GP stationary?

Ans **Major I**) first condition is Marginals same  $\Rightarrow E\{X_{t_1}\} = E\{X_{t_2}\} \forall t_1, t_2$

In other words  $m_x(t)$  is a constant  
 (mean function)

and  
 $\text{var}(X_{t_1}) = \text{var}(X_{t_2}) \forall t_1, t_2$

and  $C_x(t_1, t_1) = C_x(t_2, t_2)$

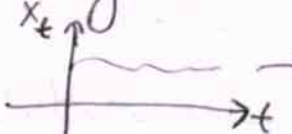
second condition is pairwise dist. only depend on  $|t_1 - t_2|$

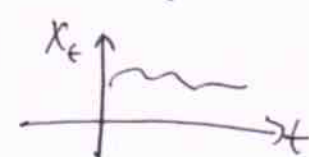
i.e.  $C_x(t_1, t_2)$  is function of  $|t_1 - t_2| \forall t_1, t_2 \in J$ .

In fact these are necessary & sufficient conds. for GP to be stationary  
 (why?)

Before going into the next example of Poisson Process, let us understand that we can view a r.p. in a slightly different way:

### Alternative View of r.p.

Suppose I conduct the random experiment of ~~coin~~<sup>coin</sup> toss repeated flip and observe the corresponding Bernoulli process. What I would register is a sequence (infinitely long) of 0 & 1. Now imagine a conts. time, conts. - state r.p. Then if  $\omega \in \Omega$ , plot the values taken by r.p. then we will get some  curve like this

Now when we repeat the expt we may get  and so on.

So we can also view a r.p. to be that which outputs a function of time for every outcome in a rad. expt. In other words:


$$X : \Omega \times J \rightarrow \mathbb{R} \quad \text{is a r.p.}$$

(compare with defn. of r.v. as  $X : \Omega \rightarrow \mathbb{R}$ )

i.e.  $X(\omega, t_1)$  is some number  $\forall \omega \in \Omega$ ,  $t_1 \in J$ . time domain.

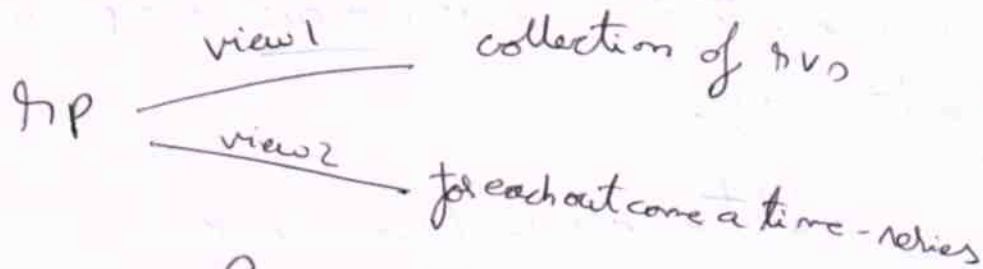
$\downarrow$   
 sample space  
 of  
 the rad. expt.

$X(\cdot, t_1)$  is nothing but the r.v.  $X_{t_1}$  i.e.  $X_{t_1} : \Omega \rightarrow \mathbb{R}$

$X(\omega, \cdot)$  is nothing but a curve  i.e.  $X(\omega, \cdot) : J \rightarrow \mathbb{R}$



To summarize



## POISSON PROCESS

Consider the P-Mart supermarket & consider a RV  $N_t$ : no. people who arrived at ~~visited~~ it before and till  $t$

~~HP~~ Consider the continuous time HP  $\{N_t\}$ .

Note that since values  $N_t$  can take is discrete (possibly countably infinite) the state space is countably infinite.

Such a HP which represent counts of no. events occurred since a reference '0' are known as counting processes. Here is formal defn.

(i)  $T \subseteq \mathbb{R}^+$  (set of positive reals)

(ii)  $N_0 = 0$  (iii)  $N_t$  takes on ~~countable values~~ values as ~~integers~~ whole numbers.

(iv) If  $t_1 \leq t_2$ , then  ~~$N(t_1) \leq N(t_2)$~~   $N_{t_1} \leq N_{t_2}$

(v) of course for a counting process no. events occurred between  $(t_1, t_2] = N_{t_2} - N_{t_1}$ .

~~It is~~ Often two kinds of assumptions are made in a counting process:

(a) No. events in disjoint intervals of time are independent  
~~from~~ called the independent increment assumption.

(Looks like for DMart example it is fair to assume this ~~less~~ unless people look at crowd inside and decide not to visit :D)

(b)  ~~$N_t$  is independent of  $t$~~  <sup>e.g.</sup>  $P\{N_{t_1} \leq x, N_{t_2} - N_{t_1} \leq y\} = P\{N_{t_1} \leq x\} P\{N_{t_2} - N_{t_1} \leq y\}$

(c)  ~~$N_{t_1}$  is independent of  $t_2$~~  (may depend on  $t_1$ )  $\forall t_1 \leq t_2$

(d) No. events occurring in a duration of time is ~~independent~~  
~~when~~ dependent (if at all) only on the duration rather than the exact start & end times.

$$\text{i.e. } P\{N_{t+s} - N_t \leq x\} = P\{N_s \leq x\} \quad \forall t, s$$

This is called the stationary increment assumption

(This may not be correct assumption for DMart as there are peak hours, market closed etc.)

However we might consider ~~an~~ <sup>analysis</sup> a duration ~~in~~ which this kind of ~~is~~ up & down behaviour is not there for e.g. analyze the peak hour itself!

So there are real-world eg. where these assumptions make sense.)



Now a Poisson process is a special counting process which makes both assumptions of i.e. stationary as well as independent increment.

One definition of Poisson process is here: (other equivalent defn. are possible refer Ross's Book)

A counting process satisfying independent increments:

(i) independent increments cond.

$$(ii) P[N_{t+s} - N_s = n] = \frac{e^{-\lambda t} (\lambda t)^n}{n!} \quad \forall t, s \geq 0.$$

Note that  $\downarrow$  says that Poisson process satisfies assumption of stationary increments. Next, No. events occurring in duration  $t$  is simply given by Poisson RV with parameter  $\lambda t$ .

$$\text{Now take } s=0 \Rightarrow P[N_t = n] = \frac{e^{-\lambda t} (\lambda t)^n}{n!}$$

$$\Rightarrow E[N_t] = \lambda t \rightarrow \text{nothing but the mean function.}$$

Since mean function is not a constant, the Poisson process is NOT a stationary process!

Note that  $\lambda = E[N_t]$ . In words,  $\lambda$  is the average rate of occurrence of events!

(Intuitively this is obviously true because it indeed matters from where I start observing the process)

Note however that  $X = \{N_t\}_{t \geq 0}$  is indeed a stationary process (only incrementally it is stationary and that is exactly stationary in defn.)

Now let's derive the auto ~~correlation~~ <sup>correlation</sup> function of  $N$ :

$$\begin{aligned}
 R_N(t_1, t_2) &= E[N_{t_1} N_{t_2}] = E[N_{t_1} (N_{t_2} - N_{t_1}) + N_{t_1}^2] \\
 (\text{Suppose } t_1 \leq t_2) &= E[N_{t_1}] E[N_{t_2} - N_{t_1}] + E[N_{t_1}^2] \\
 &= \lambda t_1 (\lambda (t_2 - t_1)) + \lambda t_1 + \lambda^2 t_1^2 \\
 &= \lambda^2 t_1 t_2 + \lambda t_1
 \end{aligned}$$

$\downarrow$  mean of Poisson  
 with parameter  $\lambda t_1$

$\downarrow$  mean of Poisson  
 with parameter  $\lambda (t_2 - t_1)$

$\downarrow$  second moment of  
 Poisson with  
 param. as  $\lambda t_1$

||| by in case  $t_1 > t_2$  we will get  $\lambda^2 t_1 t_2 + \lambda t_2$ .

~~$$R_N(t_1, t_2) = \lambda \min(t_1, t_2) + \lambda^2 t_1 t_2$$~~

$$\Rightarrow R_N(t_1, t_2) = \lambda \min(t_1, t_2) + \lambda^2 t_1 t_2$$

$$\begin{aligned}
 \Rightarrow C_N(t_1, t_2) &= R_N(t_1, t_2) - n_N(t_1) n_N(t_2) \\
 &= \lambda \min(t_1, t_2).
 \end{aligned}$$

⊗ Note that  $\lambda \min(t_1, t_2)$  is also not a function of  $|t_1 - t_2|$  again confirming that Poisson process is not stationary.

In assignment you will prove that:

$T_i$  is an exponential random variable with parameter  $\lambda$ , where  $T_i$  is waiting time for  $i^{\text{th}}$  event to occur from the instant of  $(i-1)^{\text{th}}$  event occurring.



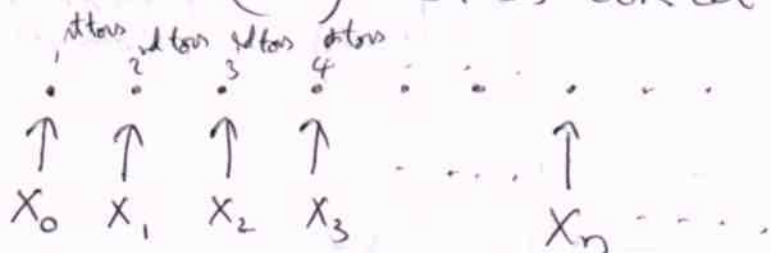
It is easy to see that  $T_1, T_2, T_3, \dots, T_n, \dots$  is a sequence of iid exponential rvs with parameter  $\lambda$ .

One can also show that  $W_i = \sum_{j=1}^i T_j$  [nothing but waiting time of event  $i$  (from the origin)] is gamma distributed. (sum of ~~independent~~ iid exponential rvs is  $\Gamma$  gamma)

↓  
none as Poisson  
Nocov's parameter

### Markov chains (MC)

To motivate (MC) let us look at the changeover problem again:



Consider the r.p.  $X$  which is collection of  $\{X_n\} \forall n \in \mathbb{W}$  <sup>set of whole numbers.</sup> where each  $X_n$  is indicator of changeover at  $(n+1)^{\text{th}}$  time.

Note that  $P\{X_0=0\}=1$

$$P\{X_n=1\} = 2p(1-p) \forall n, \quad P\{X_n=0\} = 1 - 2p(1-p) \forall n.$$

$$\text{Also, } P\{X_2=1, X_1=1\} = p^2(1-p) + p(1-p)^2 = p(1-p)$$

(probabilities as  $\begin{matrix} HTH \\ THH \end{matrix}$ )

$$\Rightarrow \text{III by } P\{X_n=1, X_{n-1}=1\} = p(1-p)$$

~~admitted~~  $\Rightarrow X_n$  and  $X_{n-1}$  are NOT independent

III by we can show that  $X_n$  &  $X_m$  are independent if  $|n-m| > 1$

A generalization of this notion that neighbouring rvs are not independent & far off rvs are independent is what is known as the Markov Property!

In these couple of lectures (left in Prob. theory) we will study Markov chains (in fact overview is a better word than study). Before venturing into that, we will introduce few useful notions regarding independence of r.v.s:

Suppose we are concerned with Discrete r.v.s only.  
We know,

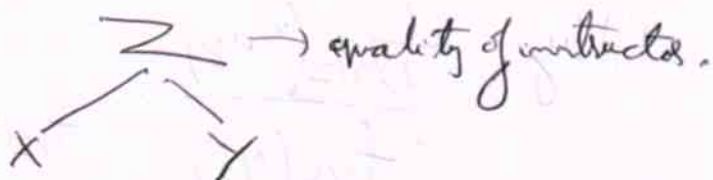
$$X, Y \text{ are independent} \iff f_{XY}(x, y) = f_X(x) f_Y(y) \quad \forall x, y$$

(joint-pmf factorizes)

Now, in some situations it may so happen that two r.v.s are not independent but given information abt another r.v., they ~~can~~ behave like independent r.v.s!

eg Let  $X, Y$  be marks obtained ( & reference) <sup>by</sup> of two students in IPL class. Initially one might be tempted to say that  $X, Y$  are independent. But a bit more thought will convince us that reference of students (in general) also depends on the quality of the instructor.

For instance, if the instructor is poor (like me), then the reference of students is <sup>more</sup> likely to be poor and vice-versa. I can represent this relation:





However, given the quality of instructor, then the performances are pretty much independent!



given  $Z=z$

$X, Y$  are linked thru  $Z$  and when  $Z$  is given, they are no more linked

In such cases, we usually talk about the notion of conditional independence:

We say  $X$  is cond. ind. of  $Y$  given  $Z$  iff:

$$\textcircled{I} \quad f_{X,Y/Z}(x,y/z) = f_{X/Z}(x/z) f_{Y/Z}(y/z) \quad (\text{conditional pmf's factorize})$$

$\forall x, y, z$

Note that  $X, Y$  are ind  $\implies X, Y$  are cond. ind. given  $Z$

$\forall Z$

However  $X, Y$  are cond. ind. given  $Z \not\implies X, Y$  are independent.

$\hookrightarrow$  (we already gave eg. of student's performance).

Now  $\textcircled{I}$  gives:

$$\frac{f_{X,Y/Z}(x,y/z)}{f_{X/Z}(x/z)} = \frac{f_{X,Y,Z}(x,y,z)/f_Z(z)}{f_{X,Z}(x,z)/f_Z(z)} = \frac{f_{Y/Z}(y/z)}{f_Z(z)} \stackrel{\text{by } \textcircled{I}}{=} f_{Y/Z}(y/z)$$

$$\text{||| by } \frac{f_{X,Y/Z}(x,y/z)}{f_{Y/Z}(y/z)} = \frac{f_{X,Y,Z}(x,y,z)/f_Z(z)}{f_{Y,Z}(y,z)/f_Z(z)} = f_{X/Y/Z}(x/y,z) \stackrel{\text{by } \textcircled{I}}{=} f_{X/Z}(x/z)$$

In other words, if  $X, Y$  are cond. ind. given  $Z$ , then:

$$f_{Y/XZ} = f_{Y/Z} \quad ; \quad f_{X/YZ} = f_{X/Z}$$

$X$  provides no more information given  $Z$  abt  $Y$

$Y$  provides no more information given  $Z$  abt  $X$ .

We are just dropping the arguments of cond. pmf to ease notation

Also,  
both these conditions imply each other.

III) My notion can be extended to 4 r.v.s:

$X$  is cond. ind. of  $Y, Z$  given  $W$  means

$$f_{X/YZW} = f_{X/W} \quad \overset{\text{of course}}{\iff} \quad f_{Y/Z/XW} = f_{Y/Z/W}$$

In words,  $Y$  and  $Z$  is cond. ind. of  $X$  given  $W$

be aware that this does not mean

that  $X, Y, Z$  are cond. ind. given  $W$ .

$$\iff f_{XYZ/W} = f_{X/W} f_{Y/W} f_{Z/W}$$

$$\iff f_{X/YZW} = f_{X/W} \quad \text{and} \quad \underline{f_{Y/ZW} = f_{Y/W}}$$



III) My notion can be extended <sup>n gvs</sup> to like this:

~~$X_n$  is cond. ind. of  $X_1, \dots, X_{n-1}$~~

$$f_{X_n / X_{n-1}, X_{n-2}, \dots, X_2, X_1} = f_{X_n / X_{n-1}} \quad \text{II}$$

This in words is name as saying  $X_n$  is cond. ind. of  $X_{n-2}, X_{n-3}, \dots, X_1$  given  $X_{n-1}$ .

In fact, this property is known as Markov property and leads us to the defn. of Markov chain:

### MARKOV CHAINS

A Markov Chain (MC) is a discrete time, discrete-state random process where the Markov property II holds for all  $n$ .

i.e.

$$X_0, X_1, X_2, \dots, X_{n-1}, X_n, X_{n+1}, \dots$$

This collection (countable collection) of gvs is called a MC iff:

(i) ~~the~~ all the gvs ( $X_n \forall n$ ) take on discrete values i.e. values from a discrete set  $S = \{s_1, s_2, \dots\}$

state-space  $\swarrow$   
states  $\searrow$

(ii) Markov property II holds  $\forall n (\geq 1)$

i.e. for eg  $f_{X_3 / X_2, X_1, X_0} = f_{X_3 / X_2}$  ;  $f_{X_2 / X_1, X_0} = f_{X_2 / X_1}$  non.

Now,

$$f_{X_0 X_1 \dots X_n} = f_{X_0} f_{X_1/X_0} f_{X_2/X_1 X_0} f_{X_3/X_2 X_1 X_0} \dots f_{X_{n-1}/X_{n-2} X_{n-3} \dots X_1 X_0} f_{X_n/X_{n-1} X_{n-2} \dots X_1 X_0}$$

(Note that this relation holds for any r.v.s  $X_0, \dots, X_n$  and NOT particular to MC).

However in case of MC; the above relation boils down to:

$$f_{X_0 X_1 \dots X_{n-1} X_n} = f_{X_0} f_{X_1/X_0} f_{X_2/X_1} f_{X_3/X_2} \dots f_{X_{n-2}/X_{n-3}} f_{X_{n-1}/X_{n-2}} f_{X_n/X_{n-1}} \text{ (III)}$$

(You can view this as first step towards forgoing assumption of independence of r.v.s  $f$ ; as for ind r.v.s we would have  $f_{X_0 \dots X_n} = f_{X_0} f_{X_1} \dots f_{X_n}$  and of course the weaker notion of independence i.e. the cond. ind. helped us in doing this)

Note that (III) implies that just by specifying one marginal (i.e.  $f_{X_0}$ ) and specifying  $n$  <sup>pairwise</sup> cond. for neighbours ( $f_{X_n/X_{n-1}} \forall n \in \mathbb{N}$ ), one can specify the MC (as all orders of joint pmf's would have then been specified).

Here is some notation:



Usually we denote:

$$f_{X_n}^{(n)} \text{ as } T_n(x) \text{ where } x \in S$$

$$\text{and } f_{X_n/X_{n-1}}(x/y) \text{ as } P(n, y, x) \text{ where } x, y \in S \text{ \& } n \in \mathbb{N}$$

In other words, to specify a MC we just need:

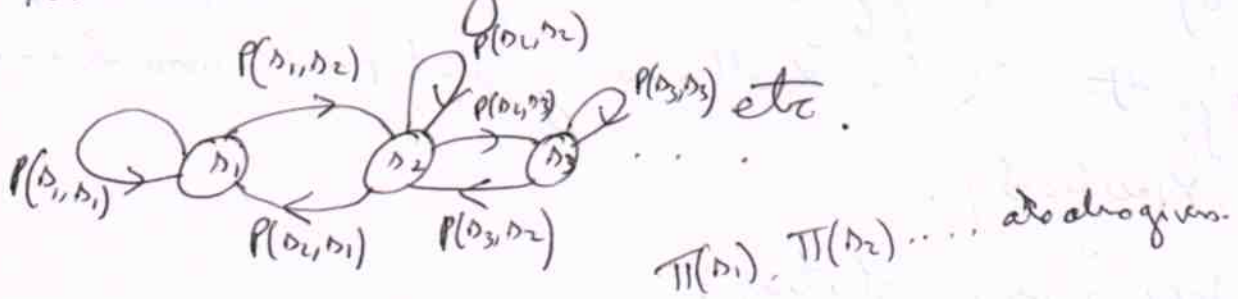
$$T_0(x) \forall x \in S \quad \text{i.e. the initial state prob. dist.}$$

$$P(n, y, x) \forall x, y \in S, \forall n \in \mathbb{N} \quad \text{i.e. the state transition probabilities for all pairs of states and for all stages } n$$

If it so happens that  $P(n, y, x)$  is independent of  $n$ ,

i.e. if  $P(n, y, x) = P(y, x)$ , then these MCs are known as homogeneous MC.

One can also view a MC as a state transition process, where initial state prob. dist. is given and the law of transitions (i.e. state transition prob.) are given. And hence MC can also be represented using a state transition diagram:



eg 1 Consider the (countably infinite) collection of rvs which represent "indicator of seeing a changeover at ~~the~~  $i$ th toss".

$\begin{matrix} \text{1st toss} & \text{2nd toss} & \text{3} & \dots & \text{n} & \text{n+1} & \text{n+2} & \dots \\ \dot{X}_0 & \dot{X}_1 & \dot{X}_2 & \dots & \dot{X}_{n-1} & \dot{X}_n & \dot{X}_{n+1} & \dots \end{matrix}$

$\downarrow$   
always zero

Each  $X_n$  indicator of a changeover at  $(i+1)$ th toss.

In last class we argued that

$X_n$  &  $X_m$  are  $\begin{cases} \text{independent if } |n-m| > 1 \\ \text{dependent if } |n-m| \leq 1 \end{cases}$

Hence the Markov property is indeed satisfied. ~~How~~ Hence ~~this~~ this tp is a MC.

However note that, MC only requires that:

$X_n$  is cond. ind. of  $X_{n-2}, X_{n-3}, \dots, X_1, X_0$  given  $X_{n-1}$

but here a stronger condition than this i.e.:

$X_n$  is ind. of  $X_{n-2}, X_{n-3}, \dots, X_1, X_0$  is being satisfied.

So this tp is a MC, but something "more" than a MC.

Now one can also consider the (problem in Mid-Sem) collection of rvs which indicate presence of HHH, in that case it will turn out

that  $X_n$  &  $X_m$  are  $\begin{cases} \text{ind. if } |n-m| > 2 \\ \text{not ind. otherwise.} \end{cases}$



These tips motivate us to define 2<sup>nd</sup> order Markov property and soon... And in general we can write down a k<sup>th</sup> order Markov property (and hence define a k<sup>th</sup> order MC):

$X_n$  is cond. ind. of  $X_{n-k-1}, X_{n-k-2}, \dots, X_1, X_0$  given  $X_{n-1}, \dots, X_{n-k}$

eg 2 Consider the gambler's ruin problem (refer prob. 4a of first assignment)  
~~It is easy to see that even in this case, the Markov property is satisfied.~~ Let  $X_n$  be the amount of money player A has when he plays the n<sup>th</sup> round of gambling.

~~$X_n$~~   $X_1 = n_1 \rightarrow$  given he starts with  $n_1$  rupees.

$$X_2 = \begin{cases} n_1 + 1 & \text{with prob } p_1 \\ n_1 - 1 & \text{" } 1 - p_1 \end{cases}$$

$$X_3 = \begin{cases} n_1 + 2 & \text{with prob } p_1^2 \\ n_1 & p_1(1-p_1) + (1-p_1)p_1 \\ n_1 - 2 & \text{with prob } (1-p_1)^2 \end{cases}$$

no on.

Note that given  $X_{n-1} = n$ , then  $P\{X_n = n+1\} = p_1$ , in other words  $P\{X_n = n-1\} = 1 - p_1$

words distribution of  $X_n$  is determined completely given  $X_{n-1}$  alone (and additional knowledge of  $X_{n-2}, X_{n-3}, \dots, X_1$  are not required)

Hence  $\{X_n, n \in \mathbb{N}\}$  is indeed a MC.

Also note that  $X_n$  &  $X_m$  are not independent for any  $n, m$ !! (why?)

So the Markov property is "tightly" satisfied (unlike eg 1).

Now, a little thought shows that this MC is not a homogeneous MC:

$$P(2, n_1, n_1+1) = p, \quad P(2, n_1, n_1-1) = 1-p,$$

$$\text{and } P(2, y, x) = 0 \quad \forall y \neq n_1, \quad \text{and } x \neq n_1+1 \text{ or } n_1-1$$

$$\text{However } P(3, n_1+1, n_1+2) = p, \quad P(3, n_1+1, n_1) = 1-p,$$

$$P(3, n_1-1, n_1) = p, \quad P(3, n_1-1, n_1-2) = 1-p,$$

$$\text{and } P(3, y, x) = 0 \quad \forall y \neq n_1+1 \text{ or } n_1-1 \quad \text{or } x \neq n_1+2 \text{ or } n_1 \text{ or } n_1-2$$

In particular  $P(2, n_1, n_1+1) = p$ , whereas  $P(3, n_1, n_1+1) = 0$

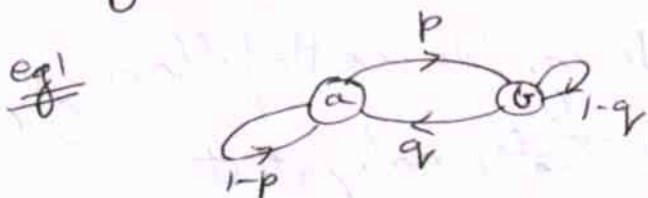
hence this MC is not homogeneous.

In next lecture we will analyze a homogeneous MC of simple nature & provide some results on MCs.



Rev.

In ~~last~~ lecture we learnt that MC can also be represented using a state-transition diagram. Consider the following eg:



(assume  $p, q \neq 0$  &  $p, q \neq 1$ )  
(considers non-trivial cases)

Now,  $\pi_n(a) = \pi_{n-1}(a)P(a,a) + \pi_{n-1}(b)P(b,a)$

Similarly,  $\pi_n(b) = \pi_{n-1}(a)P(a,b) + \pi_{n-1}(b)P(b,b)$

(I)

Since there are two states in this MC we get two eqns and in each two terms

If there were  $m$  states, we would end-up with  $m$  equations and in each eqn. we would have  $m$  terms (as when added give LHS)

Sometimes it is convenient to represent eqns. like (I) which relate  $\pi_n(a)$  with  $\pi_{n-1}(a)$  in a vectorial form:

define  $\pi_n = \begin{bmatrix} \pi_n(a_1) \\ \pi_n(a_2) \\ \vdots \\ \pi_n(a_m) \end{bmatrix}$  &  $P = \begin{bmatrix} P(a_1, a_1) & P(a_1, a_2) & \dots & P(a_1, a_m) \\ P(a_2, a_1) & P(a_2, a_2) & \dots & P(a_2, a_m) \\ \vdots & \vdots & \ddots & \vdots \\ P(a_m, a_1) & P(a_m, a_2) & \dots & P(a_m, a_m) \end{bmatrix}$

column vector of size  $m \times 1$  (min no. states)      state transition Prob. matrix size  $m \times m$

It is easy to see that:  $\pi_n^T = \pi_{n-1}^T P$  (II)   
 (this is nothing but vectorial representation of (I)).

$P$  is known as stochastic matrix or eqn. matrix steady

Of course, when  $S$  is countably infinite (i.e. no. states is not finite) then form (II) is not convenient representation and we can have this relation:

(realize that in each row of  $P$ , the row sum is 1. Such matrices are known as Stochastic Matrices.)

$$\pi_n(x) = \sum_{y \in S} \pi_{n-1}(y) P(y, x)$$

$\forall x \in S$

III

~~III & IV are the same~~ →

Nevertheless, let's come back to the eq. We have the following by writing  $P(a, b) = p$  &  $P(b, a) = q$  and realizing that

$$\pi_n(a) + \pi_n(b) = 1$$

$$\begin{aligned} \pi_n(a) &= \pi_{n-1}(a) (1-p) + (1-\pi_{n-1}(a)) q \\ &= \pi_{n-1}(a) (1-p-q) + q \end{aligned} \quad (\text{recursive relation})$$

$$\text{III by } \pi_{n-1}(a) = \pi_{n-2}(a) (1-p-q) + q$$

$$\Rightarrow \pi_n(a) = (\pi_{n-2}(a) (1-p-q) + q) (1-p-q) + q$$

= ...

$$= \pi_0(a) (1-p-q)^n + q(1-p-q)^{n-1} + q(1-p-q)^{n-2} + \dots + q(1-p-q) + q$$

$$= \pi_0(a) (1-p-q)^n + \frac{q(1-(1-p-q)^n)}{p+q}$$

$$\Rightarrow \pi_n(a) = (1-p-q)^n \left( \pi_0(a) + \frac{q}{p+q} \right) + \frac{q}{p+q} \quad \text{IV}$$

$$\text{III by we get } \pi_n(b) = (1-p-q)^n \left( \pi_0(b) - \frac{p}{p+q} \right) + \frac{p}{p+q}$$

Note that IV can also be obtained from the recursive relation I:

$$\begin{aligned} \pi_n^T &= \pi_{n-1}^T P \\ &= (\pi_{n-2}^T P) P = \pi_{n-2}^T P^2 \\ &= \dots \end{aligned}$$

$$\Rightarrow \underline{\underline{\pi_n^T = \pi_0^T P^n}} \quad \text{IV}$$



Now, a closer look at IV reveals many important things regarding the MC in the example, and thereby opens up many questions regarding generic Markov chains.

→ The first observation is that if  $\pi_0(a)$  was chosen as  $\frac{q}{p+q}$   
 $\& \pi_0(b) = \frac{p}{p+q}$ .

(Note that  $\pi_0(a) + \pi_0(b) = 1$ , no OK)

$$\begin{aligned} \text{then, } \pi_n(a) &= \pi_0(a) = \frac{q}{p+q} \\ \pi_n(b) &= \pi_0(b) = \frac{p}{p+q} \quad \forall n \end{aligned} \quad (\text{from } \textcircled{\text{IV}})$$

This says that all marginals  $f_{x_0}, f_{x_1}, \dots, f_{x_n}, \dots$  are the same.

In fact, for a HMC (homogeneous Markov chain)  $\pi_n(x)$  is independent of  $n$   $\forall n \in \mathbb{S}$

is none as saying that the HMC is a stationary  $\pi$ .

Here is why:

We know for HMC:

$$\begin{aligned} f_{x_0 x_1 \dots x_{n-1} x_n} &= f_{x_0}(x_0) f_{x_1/x_0}(x_1/x_0) \dots f_{x_n/x_{n-1}}(x_n/x_{n-1}) \\ &= \pi_0(x_0) P(x_0, x_1) \dots P(x_{n-1}, x_n) \end{aligned}$$

Now let's shift time by  $+s$  and consider:

$$f_{x_0 x_{s+1} \dots x_{n+1} x_{n+s}} = \pi_0(x_0) P(x_0, x_1) \dots P(x_{n-1}, x_n).$$

VI

By the very defn. of stationary process:

$$f_{x_0 \dots x_n} = f_{x_0 x_{n+1} \dots x_{n+D}} \quad \forall n$$

$\Rightarrow$   $\Uparrow$  by (VI)

$$\underline{\pi_0(x_0) = \pi_0(x_0)} \quad \forall x_0 \text{ (od for all } x_0)$$

In other words iff  $\pi_n(x) = \pi_0(x) \quad \forall x \in S$  then the HMC is stationary

In terms of (V):  $\pi_0^T = \pi_n^T = \pi_0^T P^n \quad \forall n \Leftrightarrow \pi_0^T = \pi_0^T P$

In terms of (III):  $\pi_0(x) = \sum_{y \in S} \pi_0(y) P(y, x)$  (VII)

Such a distribution  $\pi$  which satisfies  $\pi_0(x) = \sum_{y \in S} \pi_0(y) P(y, x)$  i.e.

$$\pi(x) = \sum_{y \in S} \pi(y) P(y, x) \quad \forall x$$

is defined as the stationary dist. of the HMC. And if  $\pi_0 = \pi$  then the HMC is stationary.

In case of eq1,  $\pi$  exists and is unique  $\pi = \begin{bmatrix} \frac{q}{p+q} \\ \frac{p}{p+q} \end{bmatrix}$

This observation obviously raises the following question:

① Will a HMC always have  $\pi$ ? If not characterize HMC based on existence of  $\pi$ .



~~Any~~ In general, the answer to question (1) is not all HMC will have a  $\pi$ .

But it is easy to see that if HMC has finite states then  $\pi$  will exist!

Proof: In case of finite states  $\pi$  is that vector which satisfies:

$$\pi^T = \pi^T P \quad (\text{by (VII)})$$

$\downarrow$                        $\downarrow$                        $\downarrow$   
 $1 \times m$                        $1 \times m$                        $m \times m$

$$\pi^T = \pi^T P \Leftrightarrow \pi^T (P - I) = 0 \Leftrightarrow \exists \text{ a linear comb. of rows of } P - I \text{ giving row of zeros}$$

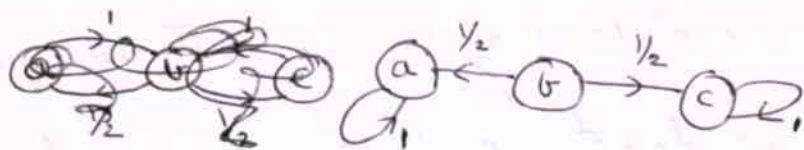
$$\Leftrightarrow \det(P - I) = 0$$

$\det(P - I)$  is indeed zero for any stochastic matrix  $P$ .

(because, if we add all columns of  $P - I$  we will get zero as each row in  $P$  sums to 1)

$\therefore$  In case of finite HMC,  $\pi$  always exists. Let's look at another example:

eg 2



By inspection of this HMC we can immediately tell that

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \text{ and } \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \text{ both satisfy (VII)}$$

Because  $\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$  says we start in (a) and if no we are bound to be in (a) because  $P(a,a)=1$ .

III) by  $\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$  says we start in (c) & we are bound to be in (c) as  $P(c,c)=1$ .

Now lets solve  $\pi^T = \pi^T P$  and see if  $\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$  and  $\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$  are the only two possible solutions or are there more?

$$\begin{pmatrix} \pi(a) & \pi(b) & \pi(c) \end{pmatrix} = \begin{pmatrix} \pi(a) & \pi(b) & \pi(c) \end{pmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 1/2 & 0 & 1/2 \\ 0 & 0 & 1 \end{bmatrix}$$

also  $\pi(a) + \pi(b) + \pi(c) = 1 \rightarrow$  (it must be valid dist. after all)

~~⇒~~ You will realize the ~~only~~ condition required is  $\pi(b) = 0$   
 $\pi(a), \pi(c)$  can be chosen anything such that they add to 1.

Hence  $\pi = \begin{bmatrix} p \\ 0 \\ 1-p \end{bmatrix} \forall 0 \leq p \leq 1$  are all stationary distributions.

So this says the simple HMC in eg 2 has uncountable number of stationary distributions! This raises an important question:

② ~~When is~~ For what kind of HMC, is  $\pi$  unique?

Again ~~into~~ we do not provide here a complete answer, however eg 2 shows that  ~~$\pi$  can be realized as~~ if  $\pi$  is not unique then it can be realized as convex combination of few fundamental stationary distributions! Here is what we mean:



Realize that  $\begin{bmatrix} p \\ 0 \\ 1-p \end{bmatrix} = p \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + (1-p) \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$

generic soln.  
for  $\Pi$

fundamental stationary dist.

convex comb.  
meaning ~~is~~ combined with ~~positive~~ <sup>non-negative</sup> weights adding to 1.

This insight turns out to be the generic nature of HMCs!  
i.e. Everytime a HMC has more than one  $\Pi$ , then it can be realized as convex comb. of few fundamental  $\Pi$ 's.

→ Next from **IV** of (eq 1) one will also find that:

$$\lim_{n \rightarrow \infty} \Pi_n(a) = \Pi(a) = \frac{q}{p+q} \quad \left( \begin{array}{l} \text{since } \lim_{n \rightarrow \infty} (1-q)^n = 0 \\ \text{as } 1-q \text{ is in } [0,1) \end{array} \right)$$

$$\lim_{n \rightarrow \infty} \Pi_n(b) = \Pi(b) = \frac{p}{p+q}$$

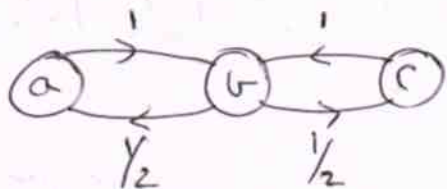
and this happens irrespective of whether  $\Pi_0 = \Pi$ !

In other words if we run this hp (in eq 1) then after a long time, it will attain asymptotic stationarity! This raises an important question:

③ ~~With~~ Suppose  $\Pi$  exists and moreover is unique. Then can we always say  $\lim_{n \rightarrow \infty} \Pi_n = \Pi$ ? (converge in pointwise sense)

In general, it turns out that the answer to above question is NO. Here is example:

eg 3



$$P = \begin{bmatrix} 0 & 1 & 0 \\ 1/2 & 0 & 1/2 \\ 0 & 1 & 0 \end{bmatrix}$$

The  $\pi$  turns out to be  $\begin{bmatrix} 1/4 \\ 1/2 \\ 1/4 \end{bmatrix}$  in this case. ( $\pi$  exists & is unique)

However it turns out that  $P^n = \begin{cases} P & n \text{ is odd} \\ P^T & n \text{ is even} \end{cases}$

~~$\pi_n = \pi_0 P^n$  if n is odd~~

$$\Rightarrow \pi_n^T = \pi_0^T P^n = \begin{cases} \pi_0^T P & \text{if } n \text{ is odd} \\ \pi_0^T P^T & \text{if } n \text{ is even} \end{cases}$$

& You can also convince yourself that unless  $\pi_0 = \pi$ , in general,  $\pi_0^T P \neq \pi_0^T P^T$ . In other words,  $\pi_n$  oscillates between two distinct vectors. Hence  $\lim_{n \rightarrow \infty} \pi_n$  (pointwise) does not even exist !!

It turns out that such things happen because of periodic nature of HMCs. For instance in (eg 3) each state can return to itself only in multiples of 2 steps! For ~~aperiodic~~ aperiodic HMC (where this kind of patterns are not there) where  $\pi$  is unique, it does happen that  $\lim_{n \rightarrow \infty} \pi_n = \pi$  (pointwise)