In the last class we observed that whenever $P[X=x] \neq 0$, the distribution function $F_X(\cdot)$ is not continuous at $x$. Now we want to study such class of distribution functions which are discontinuous at countable number of points. Here goes the formal definition:

## DISCRETE RANDOM VARIABLE:

A random variable $X$ is called a discrete random variable if there exists a countable set $E$ such that:

$$P[X \in E] = 1$$

In other words, $P[X \notin E] = 0$.

Let the set $E$ be $\{x_1, x_2, \ldots\}$. Now without loss of generality we can assume $x_1 \le x_2 \le \cdots$ and also assume $P[X = x_i] \neq 0 \ \forall \ x_i \in E$.

Now by the very defn. of discrete r.v. we have: $\sum_{x_i \in E} P[X = x_i] = 1$.

It is easy to see that the distribution function of $X$ can now be written in terms of $P[X = x_i]$ as follows:

$$F_X(x) = \sum_{\substack{x_i : x_i \le x, \\ x_i \in E}} P[X = x_i]$$

This shows that given the values of $P[X = x_i]$, the distribution function gets uniquely determined.

Also, $P[X=x_i] = P[X\le x_i] - P[X < x_i]$

$$\underset{\downarrow}{\quad} \quad \underset{\downarrow}{\quad}$$
$$F_X(x_i) \qquad F_X(x_i^-) \to (\text{the left limit})$$

Hence specifying $P[X=x_i] \ \forall \ x_i \in E$ is equivalent to specifying the distribution function $F_X(x)$ and vice-versa

We give a name to the $P[X=x_i]$ as probability mass function :

$$f_X(x) = \begin{cases} P[X=x_i] & \text{if } x = x_i \in E \\ \\ 0 & \text{if } x \notin E \end{cases}$$

In case of discrete r.v. we always specify the prob. mass function (p.m.f.) i.e. $f_X(x)$ instead of distribution function $F_X(x)$.

Now, the only constraints on $f_X$ are as follows:

(i) $\quad 0 \le f_X(x_i) \le 1 \quad \forall \ x_i \in E \qquad\qquad (\because \text{ defn. of } f_X)$

(ii) $\quad \underset{x_i \in E}{\sum} f_X(x_i) = 1 \qquad\qquad\qquad (\because P[X \in E]=1)$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad (\text{I})$

Recall that this is exactly what we did when we looked for "valid" probability functions on countable sets ! (So we already know some eg. like geometric series etc. do the job)

We explore some special discrete r.v. now:

0 Our first eg. is as follows:

Consider the set $E = \{1, 2, \ldots, n\}$ & the pmf: $f_X(i) = \frac{1}{n}$, $\forall i \in E$

It is easy to verify $f_X$ is a valid p.m.f.     (n is a parameter)

Now this can be applied to any situation where we know the outcomes are "equally likely". This basically "models" the classical probability. eg. cointoss, throwing die etc.

## Bernoulli R.V.

Consider the set $E = \{0, 1\}$ & the pmf: $f_X(1) = p$, $f_X(0) = 1 - p$

Here $0 \leq p \leq 1$ is a parameter. Again $f_X$ satisfies (I) and hence is a valid p.m.f.

This r.v. models all random expts. with two outcomes. For eg. coin toss, manufacture of good/bad parts etc. Such expts. are known as Bernoulli trials (i.e. expts. with two outcomes). Usually in a Bernoulli trial one of the two outcomes is called success $(x=1)$ and the other failure $(x=0)$. Hence prob. of success is $p$.

## Binomial R.V.

Consider the set $E = \{0, 1, 2, \ldots, n-1, n\}$ & the pmf: $f_X(i) = {}^{n}c_i \, p^i (1-p)^{n-i}$ $\forall i \in E$

Here $n \in \mathbb{N}$ & $0 \leq p \leq 1$ are two parameters to the Binomial random variable. Now lets check if $f_X$ is a p.m.f:

It is obvious that $f_X(i) \geq 0$ $\forall i \in E$

So we need to verify if $\sum_{i=0}^{n}{}^nC_i \, p^i (1-p)^{n-i} = 1$. This is indeed true because $\hookrightarrow = (p+(1-p))^n = 1^n = \nearrow$. Hence $f_x$ is a valid p.m.f.

Binomial r.v. can be employed to "model" probability of 'k' successes in 'n' independent Bernoulli trials. Recall the defn. of Bernoulli trial: it is a rand. expt with two outcomes: success (prob. $p$) and failure (prob. $1-p$). Now let ~~us~~ denote ~~us too~~ ~~why a binomial r.v.~~ the prob. space of the $i^{th}$ Bernoulli trial by $(\Omega_i, \mathcal{F}_i, P_i)$. Here $\Omega_i = \{ \text{Success}(S), \text{Failure}(F) \}$. $P_i(\{S\}) = p$ & $P_i(\{F\}) = 1-p$. $\forall$ $i = 1$ to $n$ Bernoulli trials.

Now consider the combined expt of all the $n$ Bernoulli trials. The sample space of this is $\Omega_1 \times \Omega_2 \times \cdots \times \Omega_n$. Now consider this singleton ~~set~~ event of the combined expt.: $\{ (\underbrace{S, S, \ldots, S}_{k}, \underbrace{F, F, \ldots, F}_{n-k}) \} = A$

In words this event is nothing but the event where first 'k' trials ~~Now, P(A)~~ were a success & the remaining '$n-k$' trials were failures.

Now consider the event : | We want to calculate prob. of events in combined expts. using probabilities $P_1, P_2, \ldots, P_n$.

$$A_1 = \{ (S, \omega_2, \omega_3, \ldots, \omega_n) \mid \omega_2 \in \Omega_2, \omega_3 \in \Omega_2, \ldots, \omega_n \in \Omega_n \}$$

In words, this event is the event of a success in 1st trial. Similarly define $A_i$ for $i = 1$ to $n$. Note that $A_i$ $i=1$ to $k$ represent success in $i^{th}$ trial. $A_i$ $i = k+1$ to $n$ " failures in $i^{th}$ trial.

It is easy to see that $A = A_1 \cap A_2 \cap \cdots \cap A_n$

Now let $P$ be a prob. function of the sample space $\Omega_1 \times \Omega_2 \times \cdots \times \Omega_n$

(4)

So, $P(A) = P(A_1 \cap A_2 \cap \ldots \cap A_n)$

$\qquad = P(A_1) P(A_2) \ldots P(A_n) \longrightarrow$ we assume each trial is independent of others so the event $A_i$ are independent of each other

This is the assumption of Independent trials.

Now we take $P(A_1) = P_1(\{S\})$, $P(A_2) = P_2(\{S\}) \ldots$, ~~$P(A_n) = P_n(\{S\})$~~ and so on

i.e. we construct the prob. in the combined evnt. with $\Omega = \Omega_1 \times \Omega_2 \times \ldots \times \Omega_n$ such that it is "consistent" with the probalities $P_1, P_2, \ldots, P_n$ in the individual trials!.

~~Doing similar argument~~

Hence we hae $P(A) = P(A_1 \cap A_2 \ldots \cap A_n)$ ⟍ Assumption of Independent trials

$\qquad\qquad = P(A_1) P(A_2) \ldots P(A_n)$ ← Assumption of consistency

$\qquad\qquad = P_1(\{S\}) P_2(\{S\}) \cdots \underbrace{P_k(\{S\})}_{k} \underbrace{P_{k+1}(\{F\}) \cdots P_n(\{F\})}_{n-k}$

Assumption of identical Trials $\left\{ \quad = p^k (1-p)^{n-k}. \right.$

Now we have that prob. of first $k$ trials a success & next '$n-k$' trials failure $= p^k (1-p)^{n-k}$. Now actually prob. of any '$k$' trials success & remaining failures is again $p^k (1-p)^{n-k}$. But there are exactly ${}^n C_k$ ways '$k$' successes can happen in '$n$' trials.

Hence prob. of '$k$' successes in '$n$' trials $= \underbrace{p^k(1-p)^{n-k} + \cdots + p^k(1-p)^{n-k}}_{{}^n C_k}$

$\qquad\qquad\qquad = {}^n C_k \, p^k (1-p)^{n-k}$

⑤

Now let us ask a slightly different but related question: "what is prob. that the number of trials ~~that need to be done to~~ for realizing the first success is 'k'." As we shall see below a geometric r.v. helps us ~~to do~~ model this:

## Geometric R.V.

Consider the set $E = \{1, 2, \ldots\}$ i.e. $E = \mathbb{N}$. (Note that this is the first eg. for the countably infinite discrete r.v). Define pmf as $f_X(x_i) = p(1-p)^{x_i - 1}$ $\forall x_i \in E$. It is routine to verify this $f_X$ is indeed a valid pmf.

This random variable models the trial at which the first success occurs in a sequence of ~~Ber~~ identical and independent Bernoulli trials. (IIB trials)

$X$ = trial at which the first success occured

It is easy to $X$ can take values $\{1, 2, \ldots\}$ which is exactly the set $E$ for geometric r.v. Also, using the ideas discussed in the previous section for analyzing IIB trials, we have:

$$P[X = k] = p(1-p)^{k-1}.$$

Hence a geometric variable is suitable to model "trail at which first success occurs". Now, let us look at:

$$P[X > m] = 1 - P[X \leq m] = 1 - \sum_{i=1}^{m} p(1-p)^{i-1} = 1 - p\frac{1-(1-p)^m}{1-(1-p)}$$
$$= (1-p)^m.$$

Also, $P[X > k+m / X > k] = P_X((k+m, \infty)) / (k, \infty)) = \dfrac{P_X((k+m, \infty) \cap (k, \infty))}{P_X((k, \infty))}$

$= \dfrac{P_X((k+m, \infty))}{P_X((k, \infty))} = \dfrac{(1-p)^{m+k}}{(1-p)^k} = (1-p)^m = P[X > m]$

This important observation that, $P[X > k+m / X > k] = P[X > m]$ is called the 'memory less' property. In words it says that prob. of acheiving success after $m$ steps is same immaterial of how many trials have been performed.

(6)

In this lecture we will first complete our discussions on discrete r.v. by presenting the Poisson r.v. Then we will move on to continuous r.v. — their definition, examples and applications.
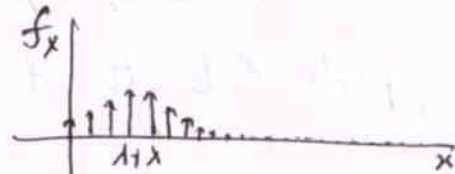
## POISSON R.V.

$$E = \{0, 1, 2, \ldots \}, \text{ the set of whole numbers.}$$

pmf: $f_X(x_i) = e^{-\lambda} \dfrac{\lambda^{x_i}}{x_i!} \quad \forall \ x_i \in E$, where $\lambda > 0$ is a parameter.

Its easy to verify $f_X$ is non-negative & $\displaystyle\sum_{x_i=0}^{\infty} e^{-\lambda}\left(\frac{\lambda^{x_i}}{x_i!}\right) = e^{-\lambda}\underbrace{\left(1 + \lambda + \frac{\lambda^2}{2!} + \cdots\right)}_{e^{\lambda}} = 1$

Hence $f_X$ is indeed a valid pmf. Now lets plot this pmf. for that lets try to look at the ratio of pmf values at two consecutive numbers:

$$\frac{f_X(k+1)}{f_X(k)} = \frac{\lambda}{k+1}$$

In other words, $f_X$ increases till $\lambda-1$ & then decreases. This plot is similar to the binomial case; the difference being that this extends to all (whole) numbers. The distribution function again is an (infinite) staircase of equal length steps by heights proportional to $f_X$ values.

Infact, after the study of concept of convergence of r.v., one can show that the binomial distribution "converges" to the poisson distribution in the case $n \to \infty$, $p \to 0$ such that $np = \lambda$,
(no. bernoulli trials is large) (prob. of success in each trial is low)

In other words,

$$P[X_b = k] \qquad\qquad P[X_p=k]$$
$$\downarrow \qquad\qquad\qquad\qquad \downarrow$$
$$^nC_k \, p^k (1-p)^{n-k} \xrightarrow[np=\lambda]{n\to\infty,\ p\to 0} e^{-\lambda}\frac{\lambda^k}{k!}$$

$X_b \to$ binomial r.v.
$X_p \to$ poisson r.v.

①

Recall that Bionomial r.v. can model 'no. successes in n Bernoulli trials'.

Hence, by the above relation, we can say that the Poisson r.v. can be used to model 'no. of ~~successes~~ occurances of a rare event in large no. B. trials'.

eg: A person keeps ~~buy~~ buying lottery tickets. The no. times he wins a lottery follows Poisson distribution (why?)

♧ No. words written by IPL instructor which are perfectly legible :)
(on board & this notes)

→————————→

Till now we have looked at random variable which took discrete values and had discontinuous distribution functions. Now lets turn our attention to r.v. whose distribution functions are continuous (infact absolutely conts.) ~~the~~ we already discussed:

$$P[X \leq x] = P[X < x] + P[X = x]$$

$$\underset{\text{(right limit)}}{F_X(x)} \quad = \quad \underset{\text{(left limit)}}{F_X(x^-)} \iff P[X = x] = 0 \quad \forall \, x \in \mathbb{R}.$$

In other words r.v. with continuous ~~&~~ distribution functions cannot have $P[X=x] \neq 0$ for any $x \in \mathbb{R}$! Hence we cannot have a "pmf" function in this case. The idea is to have a prob. density function (pdf) instead ~~which~~ and finding area under the density fuction would give probabilities.

More formally we define Continuous r.v. as follows:

(2)

# CONTINUOUS R.V.

A r.v. $X$ is called a continuous r.v. if there exists a function $f_X : \mathbb{R} \to \mathbb{R}$, called the probability density function (pdf), such that:

$$P_X(B) = \int_B f_X(x)\, dx \qquad \forall B \in \mathcal{B}$$

induced prob. function with r.v. $X$ ↙        ↳ borelo-algebra.

↓ we know how to calculate integrals when $B$ are intervals etc.

→

Here goes an intuition why $f_X$ is called a pdf:

Suppose we consider $B = (x - \varepsilon/2, x + \varepsilon/2)$ where $\varepsilon$ is tiny, i.e. $B$ is a small interval around $x$. Then, $P_X\left((x - \varepsilon/2, x + \varepsilon/2)\right) = \int_{x - \varepsilon/2}^{x + \varepsilon/2} f_X(y)\, dy = \varepsilon f_X(x)$

In other words $f_X(x) = \dfrac{P_X((x - \varepsilon/2, x + \varepsilon/2))}{\varepsilon}$ for $\varepsilon \to 0$. Since $f_X$ is ratio   ↓ since $\varepsilon$ is tiny $f_X$ does not change.

of prob. & lengths it is called as 'prob. density'.

→

Now, $F_X(x) = P_X\left((-\infty, x]\right) = \int_{-\infty}^{x} f_X(y)\, dy$     Ⓘ

In other words, given the p.d.f, the dist. func. $F_X(x)$ is fixed. Functions like $F_X$ which are expressible as integral over functions like $f_X$ are known as absolutely continuous functions. Absolute continuity is a stricter condition than continuity. In fact, we even know that $F_X$ is differentiable:

$$\frac{dF_X(x)}{dx} = f_X(x) \quad \forall x \text{ at which } f_X \text{ is continuous.} \quad Ⅱ$$

Now since one can obtain the dist. function $F_X$ given $f_X$ (pdf) and vice-versa, we characterize continuous r.v. using pdfs.

Lets look at some properties of the pdf:

③

pdf ($f_X$) ratisfies:

(i) Non-negativity: i.e. $f_X(x) \geq 0 \ \forall \ x \in \mathbb{R}$. This follows from the monotonicity of $F_X$. Since $F_X$ is are under $f_X$, there is no way the $F_X$ (area) can montonically increase if $f_X < 0$. Mathematically:

$$x_1 \leq x_2 \implies F_X(x_2) - F_X(x_1) = P_X([x_1, x_2]) = \int_{x_1}^{x_2} f_X(y) \, dy \geq 0 \quad \forall \ x_1 \leq x_2$$
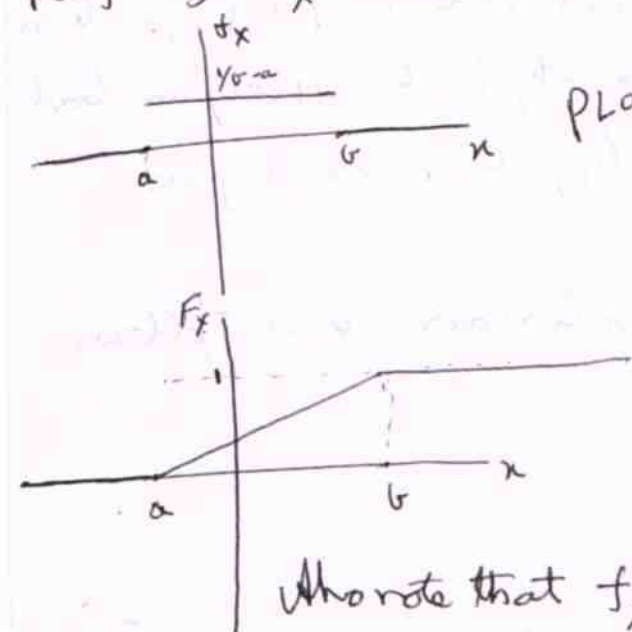
$$\implies f_X(x) \geq 0 \ \forall \ x.$$

(ii) Unit-area: We have, $1 = P_X(\mathbb{R}) = \int_{-\infty}^{\infty} f_X(x) \, dx$, Hence the area under $f_X$ must be unity.

Any fuction which ratisfies there two conditions we called it a prob. density fuction (pdf). Let us book at rome eg. of conts. r.v.

## (conts.) Uniform R.V.

$$\text{pdf:} \quad f_X(x) = \begin{cases} \dfrac{1}{b-a} & \forall \ x \in [a, b] \\ 0 & \forall \ x \notin [a, b] \end{cases}$$

· Here $a < b$ are two Parameters.

Its trivial to check $f_X$ is indeed a pdf. Now the plots of pdf & $F_X$ are:



Note the relations (I), (II) from those PLOTS graphs. The points while $F_X$ is not differentiable is exactly while $f_X$ is discontinuous.

Observe that $\frac{1}{b-a}$ can be $> 1$. So there is no reason to belive in general that $f_X(x) \leq 1$. (So $f_X(x)$ need not be $\leq 1$)

Also note that $f_X(a), f_X(b)$ can be changed to arbitrary values without changing $F_X$!
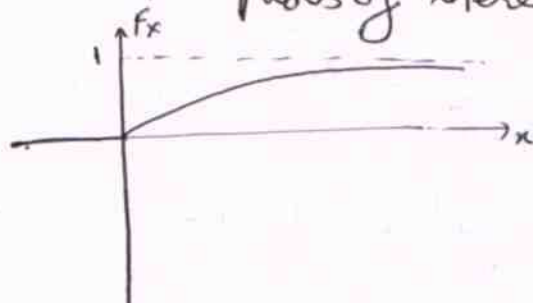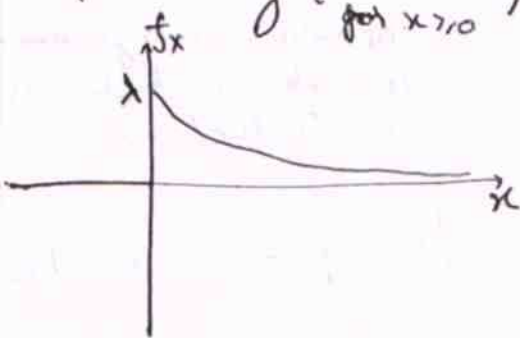
(4)

# Exponential R.V.

**. pdf :** $f_X(x) = \begin{cases} \lambda e^{-\lambda x} & x \geqslant 0 \\ 0 & x < 0 \end{cases}$    $\lambda > 0$ is a parameter

Again $f_X(x) \geqslant 0$. Also, $\int_{-\infty}^{\infty} f_X(x)dx = \int_0^{\infty} \lambda e^{-\lambda x}dx = -e^{-\lambda x}\Big|_0^{\infty} = 1$. Hence $f_X(x)$ is a valid pdf.

Now, $F_X(x) = \int P_X((-\infty, x]) = \int_{-\infty}^{x} f_X(y)dy = \int_0^x \lambda e^{-\lambda y}dy = e^{-\lambda y}\Big|_0^x = 1 - e^{-\lambda x}$.

$f_X(x)$ is an exponential decay function for $x \geqslant 0$, whereas $F_X$ is a negative exp. decay (so concave for $x \geqslant 0$). Here are the plots of these functions:  (so convex) on $x \geqslant 0$



Now, similar to the geometric r.v. in discrete case, exponential r.v. satisfies the memory-less property:

**TST** $P[X > x+y \mid X > y] = P[X > x]$     $\forall x, y \geqslant 0$.

**Proof:** $P[X > x] = 1 - F_X(x) = e^{-\lambda x}$    $\forall x \geqslant 0$.

$P[X > x+y \mid X > y] = \dfrac{P[X > x+y \cap X > y]}{P[X > y]} = \dfrac{P[X > x+y]}{P[X > y]} = \dfrac{e^{-\lambda(x+y)}}{e^{-\lambda y}} = e^{-\lambda x} = P[X > x]$

$\longrightarrow$

One can show the converse also i.e if $X$ is a conts. r.v. which is non-negative and satisfies memory-less property then $X$ MUST be exponential r.v.

**Proof:** From above proof we have that memory lessproperty is same as:

$$P[X > x+y] = P[X > x]P[X > y] \quad \forall x, y \geqslant 0$$

Now we use the fact that $X$ is a r.v., this gives:

$$[1 - F_X(x+y)] = [1 - F_X(x)][1 - F_X(y)]$$

Note this does not say the events $X > x$, $X > y$ are "independent"

⑤

Now call $G_X(x) \equiv 1 - F_X(x)$. We know that $F_X(x)$ is a continuous function, hence $G_X(x)$ is a conts. function which satisfies:

$$G_X(x+y) = G_X(x) \, G_X(y) \qquad \forall \, x, y \geq 0. \quad —①$$

Now,
$$G_X\left(\frac{m}{n}\right) = G_X\left(\underbrace{\frac{1}{n} + \cdots + \frac{1}{n}}_{m \text{ times}}\right) = \overset{m}{G_X}\left(\frac{1}{n}\right) \qquad (\because \text{ repeated application of ①})$$

Also, if $m = n$, we have $G_X(1) = \overset{n}{G_X}\left(\frac{1}{n}\right) \Rightarrow G_X\left(\frac{1}{n}\right) = \left(G_X(1)\right)^{\frac{1}{n}}$

Hence, $\qquad G_X\left(\frac{m}{n}\right) = \left(G_X(1)\right)^{\frac{m}{n}}.$

So we proved that $G_X$ is a power function for all rationals $\frac{m}{n}$ $(\geq 0)$. By continuity of $G_X$, $G_X$ must be a power function for all reals $(\geq 0)$

$$\Rightarrow \quad G_X(x) = \left(G_X(1)\right)^{x} \qquad \forall \, x \geq 0.$$

Now $G_X(1) = P_X[X > 1]$. Hence $0 \leq G_X(1) \leq 1$. Because of this I can choose a $\lambda > 0$ such that $\lambda = -\log(G_X(1))$.

$$\Rightarrow \quad G_X(x) = e^{-\lambda x}, \quad \lambda > 0. \quad \forall \, x \geq 0.$$

$$\Rightarrow \quad F_X(x) = 1 - e^{-\lambda x} \qquad \forall \, x \geq 0 \quad (\lambda > 0). \text{ which is nothing but}$$
the distribution function of an exponential r.v. Hence proved.

Thus in non-negative conts. r.v., memory-less property is unique to Exponential r.v. Now, ~~we easily~~ encouraged by this, we can apply exp. r.v. to model all cases (conts. versions) for which geometric r.v. was applicable. (Recall that geometric r.v. also is the only memory-less discrete r.v.)

$\therefore$ Expn. v. can model waiting time for a successful event.

However care needs to taken (as in care of geometric r.v.) that the physical situation make not support the memory-less property.

**Foleg:** Suppose we model the time to failure of a T.V. by an exponential random variable. Then we will be saying an absurd statement as follows:

"Let prob. of T.V. working for 10yrs. be 0.7. Then given that it already worked (new) for 10 years, the prob. that it works for 10 more yrs. is again 0.7."

### Normal R.V.

pdf: $\qquad f_X(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \qquad x \in \mathbb{R}$

Appears in many many applications. And needs no introduction.

$f_X$ is indeed non-negative. To show that $f_X$ is a valid pdf, we need to show that:

$$I \equiv \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \, dx = 1.$$
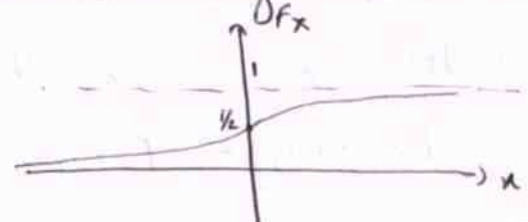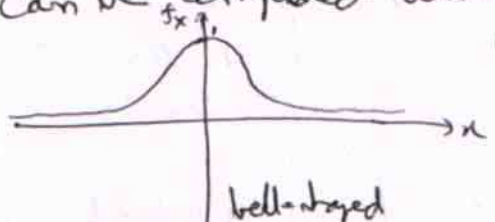
**Proof:** Consider the integral,

$$\int_{-\infty}^{\infty}\int_{-\infty}^{\infty} \frac{1}{2\pi} e^{-(x^2+y^2)/2} \, dx\,dy = \int_{-\infty}^{\infty} \left[ \underbrace{\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx}_{I} \right] \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy = I^2.$$

Now transform the integration $I^2$ using polar coordinates:

$$I^2 = \int_0^{2\pi} \left[ \int_0^{\infty} e^{-r^2/2} \, r\,dr \right] \frac{1}{2\pi} \, d\theta = \int_0^{\infty} e^{-r'} \, dr' = -e^{-r'}\Big|_0^{\infty} = 1$$

$\Rightarrow I = 1$ ~~because~~ ($I \neq -1$ because $I$ is integral of non-negative fraction).

Unfortunately $F_X(x) = \int_{-\infty}^{x} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$ cannot be computed in closed form. However can be computed using numerical integration. Here are the graphs:



bell-shaped

Consider the following random experiment: "Choose a random circle centered at origin and having radius between 0 & 1. Assume all radii are equally likely." Let the probability space for this expt. be $\mathbb{P} = (\Omega, \mathcal{F}, P)$.

Now, consider a random variable "R" defined on this probability space, which in words is "radius of the circle". In other words, R is a r.v following uniform distribution between [0,1].

Consider another mapping from $\Omega$ onto $\mathbb{R}$, which is "A": "area of the circle". ~~Is this~~ Now Note that for each circle $\omega \in \Omega$, we have:

$$A(\omega) = \pi (R(\omega))^2.$$

~~Getting~~ Defining a new function $g: \mathbb{R} \to \mathbb{R}$ such that $g(x) = \pi x^2$, it is easy to see that $A = g \circ R$, where 'o' denotes composition of functions. In other words, $A(\omega) = g\circ R(\omega) = g(R(\omega))$ $\forall \omega \in \Omega$.

We denote this as $\underline{A = g(R)}$ (abuse of notation?)

Now obvious questions are:

① Given that R is a r.v., Is $A = g(R)$ ~~always~~ a "valid" r.v?

② If so, what is the distribution of the new random variable A, which is defined in terms of ~~another~~ r.v R? (In particular, what is the distribution of the r.v "Area of circle", given that "R: radius of circle" follows uniform distribution in [0,1]?)

In this lecture, we try & answer the above questions and in general, study the notion of functions of R.V,."

The only thing we need to check is the whether :

$$A^{-1}(B) \in \mathcal{F} \quad \forall B \in \mathcal{B}$$

$$\Leftrightarrow (g \circ R)^{-1}(B) \in \mathcal{F} \quad \forall B \in \mathcal{B}$$

$$\Leftrightarrow R^{-1}(g^{-1}(B)) \in \mathcal{F} \quad \forall B \in \mathcal{B}$$

Now suppose $g^{-1}(B) \in \mathcal{B} \; \forall B \in \mathcal{B}$, then since $R^{-1}(B) \in \mathcal{B} \; \forall B \in \mathcal{B}$, by the very fact that $R$ is a r.v, we have that

$$\boxed{I} \quad \underline{g^{-1}(B) \in \mathcal{B} \; \forall B \in \mathcal{B}} \quad \text{is sufficient for } A \text{ being a}$$

valid r.v.

Now $g : \mathbb{R} \to \mathbb{R}$. One can show (not in this class) that if $g$ is conts. then the above contd. is meet. In other world if $g$ is conts then $A$ is assured to be r.v.

(In our case, $g(x) = \pi x^2$, so indeed $A$ is a r.v.)

Also note that contd. ① itself implies that $g$ is a valid r.v. with initial probability space as $\mathbb{P}_R = (\mathbb{R}, \mathcal{B}, P_R)$!

So now we exploit this fact and attempt defining prob. wrt. $A^{i.e. P_A}$ using $P_R$ (same thing as we did in case of while defining r.v!)

We define $\quad P_A(B) = P_R(g^{-1}(B)) \quad \forall B \in \mathcal{B}$

Note that $P_A$ is a valid prob. function because $g$ is itself a r.v. on $\mathbb{P}_R$ (as noted above).

To give an overall picture:

$$\mathbb{P} = (\Omega, \mathcal{F}, P)$$

$$\downarrow X$$

$$\mathbb{P}_X = (\mathbb{R}, \mathcal{B}, P_X)$$

$$\downarrow g$$

$$\mathbb{P}_Y = (\mathbb{R}, \mathcal{B}, P_Y)$$

$Y = g(X)$

$\longrightarrow$ something like picking circles at random

$\longrightarrow$ X is something like "radius of circle"
X is a r.v.

$\longrightarrow$ Y is something like "area of circle"
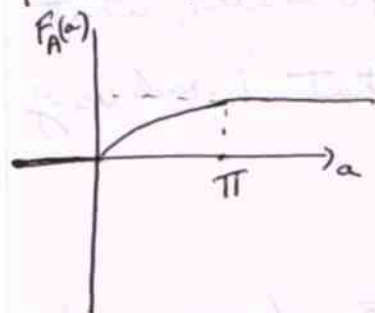g is like $\pi r^2$
g relates Y & X thru: $Y = g(X)$.

and the idea is to compute $P_Y(B) \equiv P_X(g^{-1}(B)) \longrightarrow$ this is already known!

---

Now lets answer the question "what is the distribution of "Area of circle" given radius follows uniform dist. between $[0,1]$?"

Let dist. function of A be $F_A$ & that of R be $F_R$.

$$F_A(a) = P[A \leq a] = P[g(R) \leq a]$$
$$= P[\pi R^2 \leq a] \quad \left[ = P_R(\{r \in \mathbb{R} \mid \pi r^2 \leq a\}) \right]$$
$$= \begin{cases} P\left[ -\sqrt{\frac{a}{\pi}} \leq R \leq \sqrt{a/\pi} \right] & a \geq 0 \\ 0 & a < 0 \end{cases}$$

Here's the plot:



$F_A(a)$

$\pi$ , a

$$= \begin{cases} F_R(\sqrt{a/\pi}) & a \geq 0 \quad (\because R \geq 0) \\ 0 & a < 0 \end{cases}$$

$$= \begin{cases} 0 & a < 0 \\ \sqrt{a/\pi} & 0 \leq a < \pi \\ 1 & a \geq \pi \end{cases}$$

③

Now $f_A(a) = \dfrac{d\, F_A(a)}{da}$  (∀ $a$ at which $f_A$ is conts)



$$= \begin{cases} 0 & a < 0 \\ \frac{1}{2\sqrt{\pi a}} & 0 \le a < \pi \\ 0 & a \ge \pi \end{cases}$$

Note that the distribution of "$A$: area of circle" is nowhere near uniform distribution. Also, from the pdf it looks like the values near zero are "preferred" i.e. have more prob. density. This is also intuitive as $R$ is uniform & more importantly $\le 1$! (~~By now~~ the Bertrand's Paradox also must be resolved!)

Note that the only trick is in writing the dist. function of $Y=g(x)$ in terms of dist. function of $X$. The above example would have showed that the ~~the~~ care in doing this really depends on ~~both~~ how "simple" is $g^{-1}((-\infty, x])$ for any $x \in \mathbb{R}$.

This immediately hints on considering ~~g to~~ cases where '$g$' is monotonic; because if $g$ is monotonic, then: ~~eig~~

$$g(x) \le a$$
$$\Leftrightarrow \begin{cases} x \le g^{-1}(a) & \text{if } g\uparrow \text{ (monotonically increasing } g) \\ x \ge g^{-1}(a) & \text{if } g\downarrow \quad ( \text{ " decreasing } g) \end{cases}$$

The following result is immediate:

the ~~dar following is true~~

Result 1: If $X$ is a r.v & $g$ is conts, monotonic, then ~~$\text{for} \ g(x)$~~
for the r.v $Y = g(x)$:

$$F_Y(y) = P[\, Y \le y\,] = P[\, g(x) \le y\,]$$
$$= \begin{cases} \cancel{x \le (g^{-1}(y))}\ P[\, X \le g^{-1}(y)\,] & \text{if } g\uparrow \\ P[\, X \ge g^{-1}(y)\,] & \text{if } g\downarrow \end{cases}$$

④

$$F_Y(y) = \begin{cases} F_X(g^{-1}(y)) & \\ 1 - F_X(g^{-1}(y)) + P[\,X = g^{-1}(y)\,] & \text{if } g\downarrow \end{cases}$$

Also, the following result is true:

<u>Result2</u>: Suppose further that $g$ is diff. & $X$ is conts. r.v., then:

$$f_Y(y) = \frac{dF_Y(y)}{dy}$$

(* $y$ at which $f_Y$ is conts.)

$$= \begin{cases} \dfrac{d}{dy} F_X(g^{-1}(y)) & \text{if } g\uparrow \\[2mm] -\dfrac{d}{dy} F_X(g^{-1}(y)) & \text{if } g\downarrow \end{cases}$$

($\because X$ is conts) $P[X = g^{-1}(y)] = 0$

$$= \begin{cases} f_X(g^{-1}(y)) \dfrac{d\,g^{-1}(y)}{dy} & \text{if } g\uparrow \\[2mm] -f_X(g^{-1}(y)) \dfrac{d\,g^{-1}(y)}{dy} & \text{if } g\downarrow \end{cases}$$

($\because$ chain rule)

$$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{d g^{-1}(y)}{dy} \right|$$

One can apply these results to various $g$. Let us take the care of $g(x) = ax + b$ $(a>0)$, which is a monotonically increasing diff. function.

$$Y = aX + b$$

(i) $X$ is uniform between $(0,1)$ $\longrightarrow f_Y(y) = f_X\left(\dfrac{y-b}{a}\right)\dfrac{1}{a} = \begin{cases} 0 & y < b \\ \frac{1}{a} & b \le y \le a+b \\ 1 & y \ge a+b \end{cases}$

Again models "equally likely". So we can call $Y$ as uniform distribution between $[b, a+b]$.

(ii) If $X$ is Normal, then $f_Y(y) = \dfrac{1}{a\sqrt{2\pi}} e^{-\frac{(y-b)^2}{2a^2}}$.

⑤

(iii) $X$ is uniform between $[0,1]$ & $g = H^{-1}$, where $H$ is any distribution ~~of could~~ function. ~~of a continuous~~

Note that, indeed $g$ is monotonically increasing (since $H$ is dist.func.) ~~and in fact differentiated (at most points)~~ and $g$ is also conts.

$$Y = H^{-1}(X)$$

From result 1, $F_Y(y) = F_X(g^{-1}(y)) = F_X(H(y))$

$$= H(y) \qquad \left( \begin{array}{l} \because X \text{ is uniform between} \\ \& \quad [0,1] \\ \& \quad H \text{ is a distr. function} \\ \& \text{ hence} \\ 0 \le H(x) \le 1 \ \forall x \end{array} \right)$$

In other words, if $Y = H^{-1}(X)$ where $H$ is a distr. func., then dist. of $Y$ is itself $H$! This is way of generating random ~~numbers~~ variables from a uniform r.v. itself. This can be used to generate random numbers with diff. distributions using random ~~no~~ numbers from uniform distribution. This technique of random number generation is called "~~Inverse~~ Prob. Inverse" Technique.

However it may not be useful in practice always because $H^{-1}$ may not be easily computable (for eg. for Normal dist.)

Here is an eg. where $H^{-1}$ ~~has~~ has closed form solution :

Consider $H$ as dist. of exponential r.v.

$$H(y) = 1 - e^{-\lambda y} \qquad y \ge 0 \qquad (\lambda > 0).$$

Now ~~$\cancel{X}$~~ $H^{-1}(x) = \frac{-1}{\lambda} \log(1-x)$

$\therefore$ If one takes $Y = H^{-1}(x) = -\frac{1}{\lambda} \log(1-X)$ & $X$ is uniform between $[0,1]$, ~~then~~

$$F_Y(y) = H(y) = 1 - e^{-\lambda y} \qquad !$$

⑥

This lecture introduces the concept of expectation of a r.v. or expected value or mean value of a r.v. It is denoted as $E[x]$ for a r.v X.

Intuition: We know that, in a rand. expt one cannot predict which exact value a random variable takes. However one can talk about an average value or an expected value for the r.v. The concept of expected value will help us to relate to notions like " if we toss a fair coin repeatedly on an average we will see heads for half no. times" etc.

Here goes the definition.

$$E[x] = \sum_{x_i \in E} x_i \, f_X(x_i) \qquad \text{(Discrete case)}$$

$$= \int_{-\infty}^{\infty} x \, f_X(x) \, dx \qquad \text{(Conts. case)}$$

⓵

Note that, $E[x]$ is either a series sum (probably infinite sum) of numbers which are not necessarily +ve or an improper integral over functions which are " " +ve. In such cases, the value of sum/improper integral might depend on the way we compute them. For eg. consider the Cauchy r.v. defined in Assignment problem ⑦a. There we showed that if one computes $\int_{-\infty}^{\infty} x f_X(x) dx$ by splitting it into $\int_{-\infty}^{0} x f_X(x) dx + \int_{0}^{\infty} x f_X(x) dx$, then the value is undefined. Whereas if we compute it taking $\int_{-\infty}^{\infty} x f_X(x) dx = \lim_{a \to \infty} \int_{-a}^{a} x f_X(x) dx = 0$.

Hence one additionally puts the condition that $E[x]$ is defined

①

if the corresponding sum/integral is absolutely convergent.

In other words, $\begin{cases} \text{if } \sum\limits_{x_i \in E} |x| f_X(x) \text{ converges then } E[X] = \sum\limits_{n_i \in E} x_i f_X(x_i) \\ \text{if } \int\limits_{-\infty}^{\infty} |x| f_X(x) dx \text{ converges then } E[X] = \int\limits_{-\infty}^{\infty} x f_X(x) dx. \end{cases}$

Once a sum/integral is absolutely convergent many properties satisfied by "usual" sums/integrals also get satisfied. We will indicate these as and when they are used.

eg: Consider $f_X(x) = \dfrac{1}{x^2}$, $x \geq 1$.

Here $\int\limits_{-\infty}^{\infty} x f_X(x) dx = \int\limits_{1}^{\infty} \dfrac{1}{x} dx = \infty$.

In this case the improper integral is defined and is equal to $\infty$ (unlike the case of Cauchy rv, where integral itself was undefined!) However we may choose to consider whether to include the case $E[X] = \infty$ as "well-defined" or not. For the purposes of this class we can choose $E[X] \overset{=\infty}{}$ as being "well-defined".

It is a straight-forward exercise to show that: (pls. do this exercise)

(i) $E[X] = np$ for binomial rv. (iv) $E[X] = \frac{1}{2}$ for Uniform $(0,1)$

(ii) $E[X] = \dfrac{1}{p}$ for geometric rv (v) $E[X] = \dfrac{1}{\lambda}$ for Exponential

(iii) $E[X] = \lambda$ for Poisson rv. (vi) $E[X] = 0$ for Normal rv.

In all cases, note the intuition behind each value: (let success prob. be $p$)

(i) Avg. no. of successes in 'n' trials is $np$ (iv) Center of gravity of uniform rod is at midpoint

(ii) Avg. no. trials needed for a success is $1/p$ (v) ~~Avg waiting time~~ $\frac{1}{\lambda}$, $\lambda$ is the success rate

(iii) $\lambda$ is the avg. no. successes (vi) avg. error in measurements are zero. ②

Now say $Y = g(X)$. One way to compute $E[Y]$ is to use the defn ① after computing the pmf/pdf of $Y$. Assignment prob. 9⑥ shows this can be tedious (and unnecessary).

One can directly compute $E[Y]$ using the distri. of $X$ itself:

Theorem: If $X$ is discrete, $\cancel{E[X]}$ $E[g(X)] = \sum_{x_i \in E} g(x_i) f_X(x_i)$

$X$ is cont., $\qquad E[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x) \, dx$

The proof is simple in the discrete case:

$$E[Y] = E[g(X)] = \cancel{P}\sum_{y_i} y_i \, f_Y(y_i) \qquad (\text{let } y_i \text{ be the values } Y \text{ takes on})$$

$$= \sum_{y_i} y_i \, P[g(X) = y_i]$$

$$= \sum_{y_i} y_i \sum_{x_i : g(x_i) = y_i} P[X = x_i]$$

$$= \sum_{y_i} \sum_{x_i : g(x_i) = y_i} g(x_i) \overset{f_X(x_i)}{\cancel{P[X=x_i]}}$$

$$= \sum_{x_i} g(x_i) f_X(x_i)$$

∴ We assume $E[g(x)]$ whenever defined, the sum is absolutely convergent ⟶ no order in which summation is done, does not matter!

Properties of $E[X]$: (we show them for conts r.v. but also true for discrete r.v.)

① $E$ is a linear operator: $E[aX + b] = a E[X] + b$

$$E[aX+b] = \underbrace{\int_{-\infty}^{\infty} (ax+b) f_X(x) \, dx}_{\text{by above thm.}} = a \underbrace{\int_{-\infty}^{\infty} x f_X(x) \, dx}_{E[X]} + b \underbrace{\int_{-\infty}^{\infty} f_X(x) \, dx}_{\downarrow} = a E[X] + b$$

abs. convergence of improper integral & integral is linear operator

③

②  $L \leq X \leq U$  ~~Ab~~ $\Rightarrow$  $L \leq E[X] \leq U$

$P[L \leq X \leq U] = 1$

**Proof of**

$X \leq U \Rightarrow E[X] \leq U$ :

$$E[X] = \int_{-\infty}^{\infty} x \, f_X(x) \, dx \leq \int_{-\infty}^{\infty} U f_X(x) \, dx = U$$

$\underbrace{\qquad}$ abs. conv. of implied integral & monotonicity of integral prop.

||ly $L \leq X \Rightarrow L \leq E[X]$.

here $L \leq X \leq U \Rightarrow L \leq E[X] \leq U$

$\rightarrow$ Now its easy to prove :

$P[L \leq X \leq U] = 1 \Rightarrow L \leq E[X] \leq U$

$\underbrace{\qquad}$ X is almost surely between $[L, U]$ ✓

③

✱ The "best" constant value approximation of a rv $X$ which
(optimal)

minimizes the ~~sq~~ "average squared error" in approximation is $E[X]$

In other words  $E[X] = \underset{C}{argmin} \; E\left[(X-c)^2\right]$

**Proof**

$\underset{C}{argmin} \; E\left[(X-c)^2\right] = \underset{C}{argmin} \; E\left[(X^2 - 2cX + c^2)\right]$

. repeated
application of
linearity
property of E

$= \underset{C}{argmin} \; E[X^2] - 2cE[X] + c^2$

$= \underset{C}{argmin} \; \left(c - E[X]\right)^2 + E[X^2] - \left(E[X]\right)^2$

$= E[X]$

Now, the minimized error is called as variance of $X$ $\underset{\wedge}{\{}$var$(X)$.

denoted by

i.e.  $var(X) = \underset{C}{min} \; E\left[(X - \cancel{c})^2\right] = E\left[(X - E[X])^2\right]$   $\}$ by above proof

$= E[X^2] - \left(E[X]\right)^2$

Now one can compute $var(X)$ of various rvs discussed in this course:

ⓘ $var(X) = np(1-p)$ for binomial

ⓘⓘ $var(X) = \frac{1-p}{p^2}$ for geometric

ⓘⓘⓘ $var(X) = \lambda$ for Poisson

ⓘⓥ $var(X) = 1/3$ for Uniform $[0,1]$

ⓥ $Var(X) = 1/\lambda^2$ for exponential rv

ⓥⓘ $var(X) = 1$ for Normal rv.

# Properties of var (x):

① By the very definition of var(x) it is $E$ of a non negative r.v. (which is $(X - E[x])^2$). Hence $\underline{var(x) \geqslant 0}$.

Now, $var(x) = E[(X - E[x])^2] = E\{x^2\} - (E[x])^2 \geqslant 0$

$$\Rightarrow \quad \underline{E\{x^2\} \geqslant (E[x])^2}. \qquad \text{—— Ⓐ}$$

This is a very important inequality and is a specific case of the following inequality:

$$E[g(x)] \geqslant g(E[x]) \qquad \forall\ g \ \underline{convex}$$

which is far known as the Jensen's inequality. $\left(\text{In Ⓐ, } g(x) = x^2\right)$

(fundamental inequalities like $AM \geqslant GM$ are special cases of this inequality)

② $var(X + b) = E\left[(X+b - E[x+b])^2\right] = E\{(X - E[x])^2\} = var(x)$

$\underbrace{\qquad\qquad}_{\text{linearity of } E}$

③ $var(ax) = E\{(ax - E[x])^2\} = E[a^2(x - E[x])^2] = a^2\, var(x)$

$\underbrace{\qquad\qquad\qquad}_{\text{linearity of } E}$

$\left.\begin{array}{l}\\ \\ \\ \end{array}\right\}$ Variance is not a linear operator

Now let $Y = \sigma X + \overset{\mu}{b}$ where $X$ is a Normal r.v. $(\sigma, \mu$ are some numbers$)$

$E[Y] = E\{\sigma X + \mu\} = \sigma E\{x\} + \mu = \mu \qquad (\because E[x] = 0)$

$var(Y) = var(\sigma X + \mu) = var(\sigma X) = \sigma^2 var(x) = \sigma^2 \quad (\because var(x) = 1)$

We also know, $\qquad f_Y(y) = \dfrac{1}{\sqrt{2\pi}\ \sigma}\ e^{-\frac{(x-\mu)^2}{2\sigma^2}} \qquad x \in \mathbb{R}.$

This is our new defn of Normal (Gaussian) r.v. with mean $\mu$ & variance $\sigma^2$. We call the case $\mu = 0$ & $\sigma^2 = 1$ as std. Normal r.v.

⑤

Median of a r.v. is that number "for which $P[X \le M] = \frac{1}{2}$".

Assignment problem shows that

$$M = \underset{c}{\text{argmin}} \; E[|X-c|]$$

In other words median is that value that minimizes the absolute error in approximating a r.v. by a constant.

Mode of a r.v. is the "most frequently taken value" of a r.v.

Let Mode of X be m. Mathematically,

$$m = \underset{c}{\text{argmin}} \; E[1_{\{X \ne c\}}].$$

$$1_{\{X \ne c\}} = \begin{cases} 1 & \text{if } X \ne c \\ 0 & \text{if } X = c \end{cases}$$

mode minimizes the average number of times X does'nt take its value. In other words, maximizes the avg. no. times X takes the particular value.

Now, $E[1_{\{X \ne c\}}] = 1 \; P[X \ne c] + 0 \; P[X=c] = P[X \ne c]$

In discrete case,

$$m = \underset{c}{\text{argmin}} \; P[X \ne c] = \underset{c}{\text{argmax}} \; P[X = c]$$

i.e. Mode is the value with highest prob. of occuring.

Analogously for the conts. case we have $m = \underset{c}{\text{argmax}} \; f_X(c)$.

For eg. std. Normal r.v. is "Unimodal" with $m = 0$.
(has one max $f_X(c)$)

Each such peak in pdf/pmf is called mode (loosely).
Mean, median need not be values taken by X, whereas mode must be a value taken by X.

⑥

This lecture completes our discussion on single svs by discussing concepts of generic moments, moment generating function and few important inequalities like Jensen's, Markov and Chebyshev's inequalities

## Moments & Moment Generating function

Encouraged by the notions of mean & variance etc. we now define some higher order moments of r.v. as follows:

$E[x^n]$ is called the $n$'th moment of $X$ ( first moment is mean (central) second " is ~~mean~~ moment of inertia etc. )

$E\left[(X - E[X])^n\right]$ is called the $n$'th central moment of $X$ ( $n=2$ is variance (central moment of inertia etc.) )

$E\left[(X - a)^n\right]$ is the $n$'th generalized moment of $X$ about 'a'. ( eg. moment of inertia abt some axis )

We can define the absolute value versions of these:

$E[|X|^n] \longrightarrow n^{th}$ absolute moment

$E[|X - E[X]|^n] \longrightarrow n^{th}$ absolute central moment

$E[|X - a|^n] \rightarrow n^{th}$ absolute moment about 'a'.    and so on ...

Now consider a function $M_x$ defined as follows:

$$M_x(s) = E[e^{sX}]$$

This function is known as the moment generating function (wherever it exists!)

assuming $E[e^{sX}]$ exists. (eg. of conts. case:)

$$M_x(s) = \int_{-\infty}^{\infty} e^{sx} f_x(x) dx = \int_{-\infty}^{\infty} \left(1 + sx + \frac{s^2 x^2}{2!} + \cdots\right) f_x(x) dx \quad \text{because of abs. convergence of } E[e^{sX}]$$

$$= \int_{-\infty}^{\infty} f_x(x) dx + s \int_{-\infty}^{\infty} x f_x(x) dx + \frac{s^2}{2!} \int_{-\infty}^{\infty} x^2 f_x(x) dx + \cdots$$

$$\Rightarrow \quad M_x(s) = 1 + sE[x] + \frac{s^2}{2!}E[x^2] + \dots + \frac{s^n}{n!}E[x^n] + \dots \quad \text{(2)}$$

(This is why it is called as moment generating function!)

Also it is easy to see that $\dfrac{d}{ds}M_x(s)\Big|_{s=0} = E[x]$

$$\dfrac{d^2}{ds^2}M_x(s)\Big|_{s=0} = E[x^2] \quad \text{and so on} \dots$$

In general, $\dfrac{d^n}{ds^n}M_x(s)\Big|_{s=0} = E[x^n]$. $\quad\longleftarrow$

In other words, m g f is a <u>MacLaurin</u> series with diff. given by

In general, m g f may not exist ( for eg. take care of Cauchy distribution where we know first moment itself does'nt exist!). But whenever it exists, by the above relations all moments exist (and are finite).

However the converse statement that if all moment exist then then mgf also exists may not be true in general.
(Take the log Normal dist defined as $X = e^X$ where $X$ is std. Normal and try!)
(You will notice that all moments exist but mgf does not.

Lets compute mgf for Poisson r.v:

$$M_x(s) = E[e^{sx}] = \sum_{k=0}^{\infty} e^{sk}\, e^{-\lambda}\frac{\lambda^k}{k!} = e^{-\lambda}\sum_{k=0}^{\infty} \frac{(e^s\lambda)^k}{k!} = e^{-\lambda}\, e^{e^s\lambda}$$

$$= e^{e^s\lambda - \lambda}$$

mgf for std. Normal r.v:

$$M_x(s) = E[e^{sx}] = \int_{-\infty}^{\infty} e^{sx}\frac{1}{\sqrt{2\pi}}e^{-x^2/2}\,dx = \underbrace{\frac{1}{\sqrt{2\pi}}\int_{-\infty}^{\infty} e^{\frac{s^2}{2} - \frac{(s-x)^2}{2}}\,dx}_{\text{completing quadratic}} = \frac{e^{s^2/2}}{\sqrt{2\pi}}\underbrace{\int_{-\infty}^{\infty} e^{-t^2/2}\,dt}_{t = s - x} = e^{s^2/2}$$

$$\therefore \quad M_x(s) = e^{s^2/2}$$

for std. Normal

$$\left(\because \frac{1}{\sqrt{2\pi}}\int_{-\infty}^{\infty} e^{-t^2/2}\,dt = 1\right)$$

mgf for Normal r.v: $Y = \sigma X + \mu$

$$M_y(s) = E[e^{sY}] = E[e^{\sigma s X + s\mu}] = E[e^{s\sigma X}]e^{s\mu} = e^{\frac{\sigma^2 s^2}{2}}e^{s\mu} = e^{s\mu + \frac{1}{2}s^2\sigma^2}$$

$$\Rightarrow \boxed{M_y(s) = e^{s\mu + \frac{1}{2}s^2\sigma^2}} \qquad \text{(I)}$$

$\underbrace{\phantom{e^{\frac{\sigma^2 s^2}{2}}}}$ mgf of std. Normal

Apart from fact that mgf "generates" moments there is an important application of it: mgf also characterize r.v ! In other words if somebody ~~proves that~~ proves/ascertains that a certain alien r.v has mgf for eg. as $e^{s\mu + \frac{1}{2}s^2\sigma^2}$, then certainly that alien r.v must be a Normal r.v.

[ Intuitively, here's the reason why mgf characterizes a r.v :

Take $s = j\omega$ then mgf is nothing but the Fourier transform of $f_x(x)$! $\longrightarrow M_X(j\omega) = E[e^{j\omega X}] = \int_{-\infty}^{\infty} e^{j\omega x} f_x(x)\, dx$. Hence charaterizing $f_x$ is equivalent to characterizing mgf's ]

Let's calculate $n^{th}$ moment of a log-Normal r.v : $Y = e^X$, X is Normal r.v.

$$E[Y^n] = E[(e^X)^n] = E[e^{nX}] = e^{n\mu + \frac{1}{2}n^2\sigma^2}$$

$\underbrace{\phantom{E[e^{nX}]}}_{\text{by (I)}}$

This shows that log-Normal has all moments! ( but as said earlier does'nt have an mgf)

## Jensen's Inequality

$$E[f(x)] \geqslant f(E[x]) \qquad \forall\, X,\, f \text{ convex on } \mathbb{R}.$$
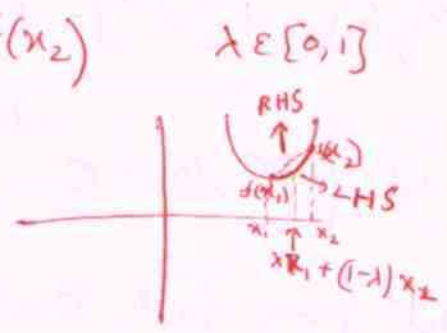
In assignment we saw a lengthily and restricted proof of this inequality following from the very defn. of a convex fuction. Now lets look at diff characterization of a convex fuction which leads to a simple proof.

$\longrightarrow f(\lambda x_1 + (1-\lambda) x_2) \leq \lambda \cdot f(x_1) + (1-\lambda) f(x_2) \qquad \lambda \in [0,1]$

This is the definition. $\xrightarrow{\text{intuition}}$

$\longrightarrow f(x) \geq f(x_0) + \dfrac{d}{dx} f(x_0) (x - x_0)$

$\frac{\text{holds, the linear approx. at } x_0 \text{ always}}{\text{under-estimates the function}}$

But this view is limited to differentiable convex functions only.

$\longrightarrow$ Above intuition holds in all convex function (even though not differentiable) i.e. At every pt. there is a supporting line

Mathematically,

$f$ is convex $\iff \exists$ a $\lambda(x)$ such that in $\mathbb{R}$

$$f(x) \geq f(x_0) + \lambda(x_0)(x - x_0) \qquad \forall x \in \mathbb{R}.$$

$\downarrow$

This is the characterization which we use now:

$X$ is a r.v. $\longrightarrow$ in other words a mapping from $\Omega: \mathcal{R} \to \mathbb{R}$. Hence take $x = X(\omega)$.

Take $x_0 = E[X]$.

$\Rightarrow \quad f(X(\omega)) \geq f(E[X]) + \lambda(E[X])(X(\omega) - E[X]) \qquad \forall \omega \in \mathcal{R}.$

Now we saw that expectation maintains order relations

$\Rightarrow E[f(X)] \geq E[f(E[X]) + \lambda[E[X]](X - E[X])]$

$\qquad\qquad = f(E[X]) + \lambda E[X](E[X] - E[X]) = f(E[X])$

Hence Proved.

# Applications of Jensen's inequality:

① Take $f(x) = -\log(x) \to$ convex. Take $X$ as discrete r.v taking values $x_1 \cdots x_n$.
& Uniform distribution. $\underbrace{\phantom{xx}}_{\geqslant 0}$

$$f(E[X]) \leq E[f(X)]$$

$$\implies -\log\left(\frac{x_1 + x_2 + \cdots + x_n}{n}\right) \leq \frac{-\log(x_1) + \cdots + \log(x_n)}{n}$$

$$= -\frac{1}{n}\log(x_1 x_2 \cdots x_n)$$

$$= -\log \sqrt[n]{x_1 x_2 \cdots x_n}$$

$$\implies \frac{x_1 + x_2 + \cdots + x_n}{n} \geqslant \sqrt[n]{x_1 x_2 \cdots x_n} \qquad (\because -\log \text{ is a monotonically} \atop \text{decr. function})$$

$$\text{AM} \quad \geqslant \quad \text{GM}.$$

So Jensen's inequality is a generalization of AM, GM inequality.

② **TST** $|\mu - M| \leq \sigma$ ($\mu$ is mean, $\sigma$ is std.div., $M$ is median).

$$|\mu - M| = |E[X] - M| = |E[X - M]| \leq \underbrace{E[|X - M|]}_{\substack{\text{Jensen's Ineq.} \\ \text{with } |\cdot| \text{ as convex} \\ \text{function!}}} \quad \cancel{E[\sqrt{(X-M)^2}]}$$

$$\cancel{E[(X-M)^2]} \Big\downarrow$$

$$= \min_c E[|X - c|] \to \text{we saw this in assignment}$$

$$\leq E[|X - c|] \quad \forall c \text{ in particular} \atop \text{take } c = \mu = E[X]$$

$$= E[|X - E[X]|]$$

$$= E\left[\sqrt{(X - E[X])^2}\right] \leq \underbrace{\sqrt{E[(X - E[X])^2]}}_{\substack{\text{Jensen's Inequality} \\ \text{with } \sqrt{\cdot} \text{ as concave function!}}} = \sigma$$

Hence Proved.

It occurs frequently in many places for eg: Information theory, Cross-entropy etc.

Till now we have been looking at random variables which take on real values. In other words, the range of r.v. was always ℝ. Now, we will generalize the notion of r.v.s to include ones taking on __vectorial__ values i.e. r.v.s for which the range is $\mathbb{R}^n$, the n-dimensional Euclidian Space.

i.e. We define functions of the form $X : \Omega \to \mathbb{R}^n$. Such functions from some sample space to $\mathbb{R}^n$ (with additional restrictions as we shall see as we proceed) are called as __Multi-variate Random variables.__ Other names are random vector, multi-valued random variable so on ...

Intuitively a (usual) random ~~vector~~ variable quantifies outcomes in terms of numbers (scalars) whereas a ~~random vector~~ multivariate r.v. (m.r.v.) quantifies outcomes in terms of n-tuples (vectors). Once this view is clear, the applications where an m.r.v. ~~is~~ can be employed are obvious; eg. whereaver the outcome of a random experiment can be defined in terms of vectorial values rather than scalar values.

To give a ~~red~~ realistic example, let us consider the random experiment where people in IITB are clinically examined for presence of Swine-Flu. Here it is immediate the each person (an outcome in our case) ~~could be~~ health cannot be described by a single quantity such as temperature or cough etc., but ~~can~~ be described using a collection of all these data!

Let us run through this example:

⓵

Let $\Omega$ be the set of all people (living) in IITB (say $N$ of them).
~~Theset~~ An event in $\Omega$ is nothing but groups of people (take $\mathcal{F}=2^{\Omega}$).
Now define $P(\{x_i\}) = \frac{1}{N} \quad \forall x_i \in \Omega$ ($x_i$ is the $i^{th}$ person).
This gives a valid prob. space $\mathbb{P} = (\Omega, \mathcal{F}, P)$.

Now define a r.v. $X_1$ which is nothing but a 'thermometer'.
(Thermometer takes input as a patient gives output as a number, specifically
the body temperature of that patient). Let $B_1$ be the set of all
"high temperatures" i.e. $B_1 = \{x \in \mathbb{R} / x \geqslant 103\}$

Now we saw, that, $P_{X_1}(B_1) = P(X_1^{-1}(B_1)) = P(\underbrace{\{\omega \in \Omega / X_1(\omega) \geqslant 103\}}_{\text{set of all people with high temp.}})$

This gives us the induced prob. space $\mathbb{P}_{X_1} = (\mathbb{R}, \mathcal{B}, P_{X_1})$.

1)) ~~by~~ for each symptom of swine flu $\left(\begin{array}{c}\text{which is of course}\\ \text{quantifiable}\\ \text{as a number}\end{array}\right)$ we can
represent it with a r.v.

Let $X_1, X_2, \ldots, X_{n-1}$ be r.v.s representing "clinical data"
quantifying each symptom of swine-flu disease. Now consider
an expert doctor who looks at the diagnostic report of patient
$\omega$ (i.e. looks at $X_1(\omega), \ldots, X_{n-1}(\omega)$) and certifies presence of
swine flu or not. ~~In other words,~~ Let $X_n$ r.v. represent the
expert doctor (again he takes as input a patient $\omega \in \Omega$ & gives as
output a number 1 (if swine-flu presence) or 0 (~~else~~ for normal patient)).
In other words $X_n$ is the indicator function of presence of disease.

(Note that $X_n$ depends "implicitly" on all $X_1, \ldots, X_{n-1}$). ②

Let $B_n \in \mathcal{B}$ be an event of drawing swine flu i.e. $B_n = $ ~~~~ $\{1\}$.

Now let us define a m.r.v. $X$ as follows:

$$X: \mathcal{R} \to \mathbb{R}^n \text{ such that,}$$

$$X(\omega) = \left( X_1(\omega), X_2(\omega), X_3(\omega), \ldots, X_{n-1}(\omega), X_n(\omega) \right) \quad \forall \omega \in \mathcal{R}.$$

*Note that intuitively $X(\omega)$ is nothing but an analyzed diagnostic report of patient $\omega$).

Note that$\quad X$ (which is an m.r.v.) not only helps in representing a "complicated" outcomes likes health of a person, but also helps in analyzing relationships between $X_i$, $X_j$ ~~r.v.'s~~ r.v.'s !!

~~Now let us see~~ ~~of there is a concept of~~ ~~reduced probability~~ ~~for a r.v.?~~ Let us see how do events in $\mathbb{R}^n$ look like:

~~In order to do that let~~ An event in $\mathbb{R}^n$ looks like:

$$B = B_1 \times B_2 \times \ldots \times B_n \qquad \text{where } B_i \in \mathcal{B} \quad \forall i$$

$$= \left\{ (x_1, x_2, \ldots, x_n) \mid x_i \in B_i, B_i \in \mathcal{B} \right\}$$

In our medical eg: $B$ is nothing but ~~high~~ temperature values in first coordinate, $\ldots$ , presence of urine fever i.e. 1 in the last coordinate.

Now collection of all such events $B$ in $\mathbb{R}^n$ is the Borel-$\sigma$-algebra in $\mathbb{R}^n$.

$$\mathcal{B}^n = \left\{ B \mid B = B_1 \times B_2 \times \ldots \times B_n , B_i \in \mathcal{B} \forall i \right\}$$

(Borel $\sigma$-algebra in $\mathbb{R}^n$)

Let us now see if there exists a concept of induced probability for a m.r.v. :

i.e. can we define $P_X(B)$ ?

~~Let~~ Following a strategy similar to the case of (usual) r.v. we have:

$$P_X(B) \equiv P\left(\{\omega \in \Omega \mid X(\omega) \in B\}\right) \quad \xrightarrow{\text{Nothing but}} \quad P\left(X^{-1}(B)\right) = P[X \in B]$$

$$= P\left(\{\omega \in \Omega \mid (X_1(\omega), X_2(\omega), \dots, X_n(\omega)) \in B\}\right)$$

$\underbrace{\qquad}_{\longrightarrow B_1 \times B_2 \times \dots \times B_n}$ , square bracket notation

$$= P\left(\{\omega \in \Omega \mid X_1(\omega) \in B_1, X_2(\omega) \in B_2, \dots, X_n(\omega) \in B_n\}\right)$$

$$= P\left(\bigcap_{i=1}^{n} \{\omega \in \Omega \mid X_i(\omega) \in B_i\}\right) \quad \xleftarrow{\text{Notation}}$$

$$= P\left[ X_1 \in B_1, X_2 \in B_2, \dots, X_n \in B_n \right]$$

(I)

In words, (in case of medical eg.) $P_X(B)$ is nothing prob. of observing high temp (first coordinate), ...., runin flu (last co-ordinate).

$P[\text{...}]$ $P\left(\bigcap_{i=1}^{n} \{\omega \in \Omega \mid X_i(\omega) \in B_i\}\right)$ is nothing but the prob. of a person having high temp., ...., & having swine flu. This is indeed intuitive. Since $P$ models classical prob' in our medical eg., this is exactly the "fraction of people in IITB having all symptoms of swine-flu & also have swine flu!". So in future classes, we will try to ~~see~~ answers~~to~~ ~~which~~ questions like which symptoms are crucial for swine-flu, how to predict presence of swine flu given a raw diagnostic report (i.e. guess values of $X_n$ given say values of $X_1, \dots, X_{n-1}$) and so on & so forth!

Math. detail
$\longrightarrow$ ~~Tech. question~~ :

Q: Why does $\bigcap_{i=1}^{n} \{\omega \in \Omega \mid X_i \in B_i\} \in \mathcal{F}$ ? (unless this happens my induced prob. defn. is invalid!)

A: I know each of $\in \mathcal{F}$ for fixed $i \in \mathcal{F}$ ($\because X_i$ is a r.v.) So intersection of them also $\in \mathcal{F}$ ($\because \mathcal{F}$ is a $\sigma$-algebra)

(4)

Now that we have correctly defined Induced prob. of a m.r.v., let us extend the concept of "distribution fuction" to an m.r.v. :

## Distribution fuctions of m.r.v.

In case of r.v. we defined $F_X : \mathbb{R} \to \mathbb{R}$ such that

$$F_X(x) = P[X \leq x].$$

Now we will define in an analogous way for an m.r.v.
(Note that range of m.r.v is $\mathbb{R}^n$):

$$f_X : \mathbb{R}^n \to \mathbb{R} \ni \quad F_X\left(\underset{\underset{\in \mathbb{R}^n}{\downarrow}}{\underline{x}}\right) = P[\underline{X} \leq \underline{x}]$$

$$\hookrightarrow \underline{x} = (x_1, x_2, \ldots, x_n)$$

Using $\boxed{I}$ we get :

$$F_X(\underline{x}) = P\{\underline{X} \leq \underline{x}\} = P[X_1 \leq x_1, X_2 \leq x_2, \ldots, X_n \leq x_n].$$

$\hookrightarrow$ called the prob. dist. func. of m.r.v. $\underline{X}$ & it is also called as joint prob. dist. func. of $X_1, X_2, \ldots, X_n$ (in this case it is

represented as $F_{X_1 X_2 \cdots X_n}\overset{\downarrow}{(x_1, x_2, x_3, \ldots, x_n)})$

Now as in case of $\overset{F_X \ of \ a}{r.v.}$ we can show the following 4 prop. for $F_X$ of a m.r.v. also :

i) $F_X(\underline{x}) \geq 0 \quad \forall \underline{x}$ (afterall its a prob.)

ii) $F_X(\underline{\infty}) = P[X_1 \leq \infty, X_2 \leq \infty, \ldots, X_n \leq \infty] = P[\underline{X} \in \mathbb{R}^n] = P(\underline{x}) = 1.$

$F_X(-\underline{\infty}) = P[X_1 \leq -\infty, X_2 \leq -\infty, \ldots, X_n \leq -\infty] = P(\{\phi\}) = 0.$

$\underline{\hookrightarrow}$ represents vector with all values $\infty / -\infty$.

⑤

This not only shows that ~~the joint~~ prob. of events can be computed in terms of dist. function but also the extra cond. that $F_{X_1 X_2}(b_1, b_2) - F_{X_1 X_2}(a_1, b_2) - F_{X_1 X_2}(b_1, a_2) + F_{X_1 X_2}(a_1, a_2) \geq 0$.

Now ⓘⓘ ~~can be proved three (cumbersome) induction using these ideas~~. is a generalization of this to a n-dimensional case.

Now using set algebra wee can show (not in this class) that $P[X \in B]$ can be written in terms of $F_X(\cdot)$. So we from

$\downarrow$
$\mathcal{B}^n$

now onwards characterize a m.r.v. $X = [X_1, X_2, \dots X_n]$ using $\left(F_X \text{ or } F_{X_1 X_2 \cdots X_n}\right)$ dist. function.

---

Lets look at the case where all $X_i$'s are discrete. Then we can define a prob. mass function (pmf) for the m.r.v. $X$: (analogous to discrete r.v's case):

$$f_X(\underline{x}) = \overset{\text{pmf of } X}{f_X(x_1, x_2, \dots, x_n)} = \overset{\text{joint pmf of } X_1, X_2, \dots, X_n}{f_{X_1 X_2 \cdots X_n}(x_1, x_2, \dots, x_n)}$$

$$\equiv P[X_1 = x_1, X_2 = x_2, \dots, X_n = x_n]$$

Now let $\mathbb{E}$ be the set of (discrete) values in $\mathbb{R}^n$ taken by $X$, then following two properties of pmf are immediate:

$$f_X(\underline{x}) \geq 0 \qquad \left( \because \overset{f_X \text{ is a pm}}{\cancel{P[\text{ }]}} \text{ all a prob} \right)$$

$$\sum_{\underline{x} \in \mathbb{E}} f_X(\underline{x}) = 1 \qquad (\because P(\Omega) = 1)$$

Now again any function ~~satisfying~~ $f_X : \mathbb{R}^n \to \mathbb{R}$ satisfying

⑦

Here, two properties is called a pmf. $\&$ also given a pmf, of $x$, dint. of $X (F_x)$ is fixed and vice-versa. So we can characterize discrete m.r.v. using pmf. (from now-onwards). Lets look at an eg. of a discrete m.r.v. {which is a generalization of the binomial r.v.) :

## Multinomial R.V.

Suppose we $\&$ define the following $f_X : \mathbb{R}^n \to \mathbb{R}$, $\bullet$

$$f_X(\underline{x}) = P[X_1 = x_1, X_2 = x_2, \ldots, X_n = x_n] = \frac{n!}{x_1! x_2! \ldots x_n!} p_1^{x_1} p_2^{x_2} \cdots p_n^{x_n} \quad \text{—ⓘ}$$

$$\forall \; \bullet \; x_i \geq 0, \; \sum_{i=1}^{n} x_i = n$$

and $(n, p_1, \ldots, p_n)$ are parameters such that $n \in \mathbb{N}$ and $p_i \geq 0$, $\sum_{i=1}^{n} p_i = 1$.

(It is an exercise to first check if this is a valid pmf!)

At first look this might look weird but consider the following random expt : Suppose I throw a die $n$ times. In each throw I $\cancel{\text{of}}$ have $(n = 6)$ outcomes. Suppose $p_i$ is probability of seeing no. '$i$' ($i=1$ to $6$). Now the answer to the question: "what is the prob. of seeing $x_1$ 1s, $x_2$ 2s, $\ldots$, $x_6$ 6s" is exactly given by ⓘ ! (why?)

⑧

Now that we know some 'physical' interpretation of multinomial $r.v.$ Lets see (if at all) what kind of $r.v.$'s are $X_1, X_2, \ldots, X_6$ individually?

In the die throwing case, $X_i = \#$ throws in which 'i' was observed

Now it is easy to see that $X_i$ is a binomial $r.v.$ with parameters $(n, p_i)$ !.

So $X = [X_1, X_2, \ldots, X_6]$ when follows multinomial distribution then each of $X_i$ $(i = 1 \text{ to } n)$ follow binomial distri. (with diff. parameters)

Now we can compute prob like
$$P[X_i = x_i] \quad \text{using the pmf of } X_i.$$

But note that $P[X_1 = x_1, \ldots, X_n = x_n]$ (which is nothing but the joint pmf of $X_1, X_2, \ldots X_n$), cannot be computed merely from the knowledge of $P[X_i = x_i]$'s. (at most you can give bounds on joint pmf using inequalities like Berferroni's inequality you proved in the assigns.)

However the pmf of $X_i$ can be computed given the pmf of $X$ (i.e. the joint pmf of $X_i$'s):

$$\sum_{\substack{x_2, x_3, \ldots, x_n \\ \ni x_i \geq 0, \sum_{i=1}^{n} x_i = n}} f_X(x_1, x_2, \ldots x_n) = f_{X_i}(x_i).$$

fixed    all varied for all allowed values

$$\text{So specifying joint pmf is a "richer" information! (than specifying pmf of } X_i's \text{ alone)}$$

(9)

(iii) $F_X$ is right conts. & has left limit"

(we dont show this here)

(iv) Monotonicity: (but in all variables)

$$F_X(x_1+\Delta x_1, x_2+\Delta x_2, \ldots, x_n+\Delta x_n) - F_X(x_1, x_2, \ldots, x_n)$$
$$\underset{\geq 0}{\downarrow} \quad \underset{\geq 0}{\downarrow} \quad \underset{\geq 0}{\downarrow}$$

$$= P[X_1 \leq x_1 + \Delta x_1, \ldots, X_n \leq x_n + \Delta x_n]$$
$$- P[X_1 \leq x_1, \ldots, X_n \leq x_n]$$

$$= P[x_1 \leq X_1 \leq x_1 + \Delta x_1, \ldots, x_n \leq X_n \leq x_n + \Delta x_n]$$

$$\geq 0 \quad (\because \text{ it is a prob. of some event}).$$

$\longrightarrow$ Till this all prop. are ~~common~~ analogous to those in case of $r.v.$

But in case of $m.r.v.$ an <u>extra condition</u> needs to be satisfied:

(v)

$$F_X(x_1+\varepsilon_1, x_2+\varepsilon_2, \ldots, x_n+\varepsilon_n) - \sum_i F_X(x_1+\varepsilon_1, \ldots, x_i, \ldots x_n+\varepsilon_n)$$

$$+ \sum_i \sum_{j>i} F_X(x_1+\varepsilon_1, \ldots, x_i, x_{i+1}+\varepsilon_{i+1}, \ldots, x_j, x_{j+1}+\varepsilon_{j+1}, \ldots x_n+\varepsilon_n)$$

$$\vdots$$

$$(-1)^n F_X(x_1, x_2, \ldots, x_n) \geq 0 \qquad \forall x_i, \varepsilon_i > 0. \qquad \boxed{II}$$

$\cancel{*}$ We can get an <u>intuition</u> for this by looking at a $r.v.$ taking values ~~from~~ in $\mathbb{R}^2$:

Let $X = [x_1 \quad x_2]$  Let $F_{X_1 X_2}(x_1, x_2)$ be the ~~di~~ joint prob. dist. function at $(x_1, x_2)$.

Now suppose I want to compute

$$0 \leq P[a_1 < X_1 \leq b_1, a_2 < X_2 \leq b_2] \text{ in terms of } \left.\right\} \text{ then:}$$

$$= F_{X_1 X_2}(b_1, b_2) - F_{X_1 X_2}(a_1, b_2) - F_{X_1 X_2}(b_1, a_2) + F_{X X_2}(a_1, a_2)$$

(how we get this was explained in class)

6

Let $X = [x_1, x_2, \ldots, x_n]$ be a m.r.v. $X$ is called a continuous m.r.v. (& equivalently $x_1, x_2, \ldots, x_n$ are said to be "jointly continuous" r.v.s) iff there exists a function $f_X : \mathbb{R}^n \to \mathbb{R}$ such that:

$$P[X \in B] = \int_B f_X(\underline{x}) \, d\underline{x} \qquad \forall \, B \in \mathcal{B}^n .$$

$\underbrace{\phantom{\int_B f_X(\underline{x}) \, d\underline{x}}}_{\text{multidimensional integral}}$

Such a function $f_X$ is called the prob. density function of $X$ & joint prob. density function of $X_1, X_2, \ldots, X_n$.

For eg. if $n = 2$,
$$P[X \in B] = \iint_B f_{X_1 X_2}(x_1, x_2) \, dx_1 \, dx_2 .$$

Now take $B = \{(a_1, a_2, \ldots a_n) \mid a_i \in (-\infty, x_i]\}$

$$\Rightarrow \quad P[X_1 \leq x_1, X_2 \leq x_2, \ldots, X_n \leq x_n] = \int_{-\infty}^{x_1} \int_{-\infty}^{x_2} \cdots \int_{-\infty}^{x_n} f_X(x_1, x_2, \ldots, x_n) \, dx_1 \, dx_2 \cdots dx_n$$

$$\underset{\downarrow}{\phantom{=}}$$
$$F_X(x_1, x_2, \ldots, x_n)$$

Hence $F_X$ is fixed if $f_X$ is known. ~~One can show the same also in~~

Also we have:

$$\frac{\partial^n F_X(x_1, x_2, \ldots, x_n)}{\partial x_1 \, \partial x_2 \cdots \partial x_n} = f_X(x_1, x_2, \ldots, x_n) \qquad \left(\begin{array}{l}\text{wherever} \\ f_X \text{ is ~~cnts~~} \\ \text{continuous}\end{array}\right)$$

As in case of 1-d r.v., the values of $f_X$ (where $f_X$ is discont.) doesn't matter (they do not account for the idea), not here) that such points are "few", so we can make the statement ~~$f_X$~~ "given $F_X$, we have $f_X$ fixed and vice-versa".

So from now on we characterize $X$ by $f_X$ (pdf)

(1)

Now lets look at some properties of $f_X$:

we have, ① $1 = P(\Omega) = P[\underline{X} \in \mathbb{R}^n] = \int_{-\infty}^{\infty}\int_{-\infty}^{\infty}\cdots\int_{-\infty}^{\infty} f_X(x_1, x_2, \ldots x_n)\,dx_1\,dx_2\ldots dx_n$

② we can now $f_X(x_1, x_2, \ldots, x_n) \geqslant 0 \quad \forall\ \underline{X} \in \mathbb{R}^n$. Recall that an analogous statement in case of $1$-v. followed from "monotonicity property of $F_X$". Here it follows from the "⑤$^{th}$ prop. of $F_X$" which in $2$-d case is illustrated below:

We know, $F_{X_1 X_2}(b_1, b_2) - F_{X_1 X_2}(a_1, b_2) - F_{X_1 X_2}(b_1, a_2) + F_{X_1 X_2}(a_1, a_2) \geqslant 0$

$\qquad \forall\ a_1 \leq b_1,\ a_2 \leq b_2$

$\Leftrightarrow \int_{-\infty}^{b_1}\int_{-\infty}^{b_2} f_X(x_1, x_2)\,dx_1\,dx_2 - \int_{-\infty}^{a_1}\int_{-\infty}^{b_2} f_X(x_1, x_2)\,dx_1\,dx_2 - \int_{-\infty}^{b_1}\int_{-\infty}^{a_2} f_X(x_1, x_2)\,dx_1\,dx_2 + \int_{-\infty}^{a_1}\int_{-\infty}^{a_2} f_X(x_1, x_2)\,dx_1\,dx_2 \geqslant 0$

$\qquad \forall\ a_1 \leq b_1,\ a_2 \leq b_2$

$\Leftrightarrow \int_{a_1}^{b_1}\int_{-\infty}^{b_2} f_X(x_1, x_2)\,dx_1\,dx_2 - \int_{a_1}^{b_1}\int_{-\infty}^{a_2} f_X(x_1, x_2)\,dx_1\,dx_2 \geqslant 0$

$\Leftrightarrow \int_{a_1}^{b_1}\int_{a_2}^{b_2} f_X(x_1, x_2)\,dx_1\,dx_2 \geqslant 0 \quad \forall\ a_1 \leq b_1,\ a_2 \leq b_2 \Leftrightarrow f_X(\underline{x}) \geqslant 0\ \forall\ \underline{x}$

Hence, pdf is any function that satisfies:

① $f_X(\underline{x}) \geqslant 0 \quad \forall\ \underline{x}$

② $\int_{-\infty}^{\infty}\int\cdots\int f_X(\underline{x})\,d\underline{x} = 1$

Now lets look at a particular eg. of a conts. m.r.v.:

$$f_X(\underline{x}) = \frac{1}{(2\pi)^{n/2}}\, e^{-\frac{1}{2}\underline{x}^T\underline{x}} \quad \forall\ \underline{X} \in \mathbb{R}^n. \qquad ⓘ$$

②

First lets check if it is pdf?

$f_X$ is indeed non-negative. Only non-trivial thing to verify is of it integrates to unity:

$$\int_{-\infty}^{\infty}\int_{-\infty}^{\infty}\cdots\int_{-\infty}^{\infty} \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2}\sum_{i=1}^{n} x_i^2} \, dx_1 \, dx_2 \cdots dx_n$$

$$\underbrace{\prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x_i^2}}$$

$$= \left( \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x_1^2} \, dx_1 \right) \left( \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x_2^2} \, dx_2 \right) \cdots \left( \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x_n^2} \, dx_n \right) = 1$$

(each of integral is 1).

Now, $F_{X_1 X_2 \cdots X_n}(x_1, x_2, \ldots, x_n) = \int_{-\infty}^{x_1}\int_{-\infty}^{x_2}\cdots\int_{-\infty}^{x_n} \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2}\sum x_i^2} \, dx_1 \, dx_2 \cdots dx_n$

$$= \left( \int_{-\infty}^{x_1} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x_1^2} \, dx_1 \right) \cdots \left( \int_{-\infty}^{x_n} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x_n^2} \, dx_n \right)$$

$$= \underbrace{F_{X_1}(x_1) \, F_{X_2}(x_2) \cdots \cdots F_{X_n}(x_n)}$$

each is the distribution function
of the std. Normal r.v !!

There are two things to note about ① :

① Its distribution fn. is product of dist. fn. of individual r.v.s

② Dist. fn. of each individual r.v.
is the std. Normal dist.

(later on we will see that
such r.v.'s are called as
independent r.v.'s)

③

Now both in case of discrete and conts. m.r.v., we saw that the distribution functions of pmfs of pdfs of individual r.v.'s can be obtained from their joint distribution. This leads to the notion of Marginal distributions:

Now $F_{X, X_2}$ From now onwards to simplify the notation we will consider collections of two r.v.s. However keep in mind that the analogous results do hold in the generic n-d case also.

So from now onwards consider two r.v. X, Y., $F_{XY}$ ~~( )~~ is the joint dist. function of X and Y.

Now, $F_X(x) = P[X \le x] = P\{X \le x, \; Y \le \infty\} = F_{XY}(x, \infty) \quad \forall x$

lll-ly. $F_Y(y) = F_{XY}(\infty, y) \quad \forall y.$

The dist. functions of X/Y are also known as the marginal dist. of X/Y wrt the joint dist. function of X and Y.

Now if X, Y are discrete,

$$f_X(x) = P[X = x] = \sum_{\forall y} P[X = x, Y = y] = \sum_{\forall y} f_{XY}(x, y) \quad \forall x$$

lll-ly $\quad f_Y(y) = \sum_{\forall x} f_{XY}(x, y) \quad \forall y$

Again $f_X, f_Y$ are known as the marginal pmfs of X and Y wrt. to the joint pmf of X, Y i.e. $f_{XY}$.

Now suppose X, Y are jointly conts:

④

$$\Rightarrow F_{XY}(x,y) = \int_{-\infty}^{x}\int_{-\infty}^{y} f_{XY}(x',y')\,dx'dy'.$$

Now $F_{XY}(x,\infty) = \int_{-\infty}^{x}\int_{-\infty}^{\infty} f_{XY}(x',y')\,\underbrace{\cancel{dx'}}\,dy'dx'$ $\left.\phantom{\int}\right\}$ (II)

$\|$

Also, $F_X(x) = \int_{-\infty}^{x} f_X(x')\,dx'$

Since the pdf is fixed (except at "few" points) given the dist. fun.,

we have that $\qquad f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x,y')\,dy'$ $\qquad$ from (II)

$\qquad$ IIIly $\quad f_Y(y) = \int_{-\infty}^{\infty} f_{XY}(x',y)\,dx'$

Again $f_X, f_Y$ are known as the marginal pdfs of $X$, and $Y$
w.r.t. to the joint-pdf of $X$ and $Y$ ie. $f_{XY}$.

Now, it is easy to see that in example (I), the following holds:

$$f_{X_1 X_2 \dots X_n}(x_1, x_2, \dots, x_n) = \underbrace{f_{X_1}(x_1)\,f_{X_2}(x_2)\dots f_{X_n}(x_n)}_{}$$

$$\underbrace{\downarrow \qquad\qquad \downarrow \qquad\qquad \downarrow}_{\substack{\text{each is pdf of} \\ \text{std. Normal.}}}$$

Now lets run through the calculation of marginal pdfs using
a tony-example:

$$\text{Let } f_{XY}(x,y) = \begin{cases} 24xy & \text{if } 0<x,\ 0<y,\ 0<x+y<1 \\ 0 & \text{otherwise} \end{cases}$$

be the joint pdf of $X$ and $Y$. Let us compute the marginals: (5)

$$f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x,y)\,dy \qquad \text{(we know that } 0<x,\ 0<y \atop 0<x+y<1 )$$

$$= \begin{cases} \int_{0}^{1-x} 24xy\,dy & \text{if } x \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$

$$= \begin{cases} 12x(1-x)^2 & x \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$

Since expression of $f_{xy}$ is symmetric in terms of $x,y$, we will get that

$$f_Y(y) = \begin{cases} 12y(1-y^2) & y \in [0,1] \\ 0 & \text{otherwise.} \end{cases}$$

Also, 

$$\int_{-\infty}^{\infty}\int_{-\infty}^{\infty} f_{XY}(x,y)\,dx\,dy = \int_{-\infty}^{\infty} f_X(x)\,dx = \int_{0}^{1} 12x(1-x)^2\,dx$$

$$= \int_{0}^{1} 12x^2(1-x)\,dx = 4x^3 - 3x^4 \Big|_{0}^{1}$$
$$= 1$$

This verifies that both the marginals and in turn the joint pdf are indeed "valid" pdfs!

Now recall the example of "swine-flu" analysis done in a previous lecture. Suppose I've want to evaluate the truth in the statement that given 'n' symptoms of swine-flu are indeed good indicators of presence of disease". To answer such question we would (say) (absence)

consider the set of all ppl. who have the symptoms and then look at frac. of ppl. in that set who also have high swine-flu. (if this value is high, then symptoms are indeed indicators). In other words, we need to ask questions abt. probabilities as if the original reduced Ω is shrinked to the set of all people having the symptoms. As we explained in one of the early lectures, such conditional probability is a mechanism which facilitates this "shrinkage of Ω".

Defn : → Let's now look at the concept of cond. prob. transfere. defined in terms of usual conditional prob. over events

Consider two discrete RVs X, Y with joint pmf given by $f_{XY}$. Suppose $p_Y(y_i) = P[Y = y_i] \neq 0$ for some fixed value of $y_i$. Now define a new prob. mass function :

$$f_{X/Y}(x_i/y_i) = \frac{f_{XY}(x_i, y_i)}{f_Y(y_i)} \quad \forall x_i \rightarrow \text{values which random variable } X \text{ takes.}$$

notation for conditional probability mass function of X given Y = $y_i$

First of all, let's check if (this) is a valid pmf ?

①

First of all its a ratio of values of pmf's and hence is $\geq 0$. Secondly,

$$\sum_{\forall x_i} f_{X|Y}(x_i/y_i) = \sum_{\forall x_i} \frac{f_{XY}(x_i, y_i)}{f_Y(y_i)} = \frac{\sum_{\forall x_i} f_{XY}(x_i, y_i)}{f_Y(y_i)} = \frac{f_Y(y_i)}{f_Y(y_i)} = 1$$

marginal pmf
definition.

Now for each value of $y_i$ such that $f_Y(y_i) \neq 0$, we can define a different pmf. Hence we have a family of pmf. all derived from the joint pmf of $X, Y$.

Note that, $f_{X|Y}(x_i/y_i) = \frac{f_{XY}(x_i, y_i)}{f_Y(y_i)} = \frac{P[X=x_i, Y=y_i]}{P[Y=y_i]} = P[X=x_i | Y=y_i]$

defn. of cond.
prob. of over events
which is familiar to
all of us.

In other words, the conditional pmf is defined in terms of cond. prob. over events.

cond.

Now once we have a pmf we can also define conditional distribution function : $F_{X|Y}$

$$F_{X|Y}(x_0/y_i) = \sum_{\forall x_i \leq x} f_{X|Y}(x_i/y_i) = \sum_{\forall x_i \leq x} P[X=x_i | Y=y_i]$$

$$= P[X \leq x_0 | Y=y_i]$$

Lets try to put down the cond. prob. for a toy example.
(Note that, instead we could have started by defining cond. distri. and Later discovered the name defn. of cond. pmf!)

②

Consider the ~~too~~ usual prob. space associated with tossing of a coin (with prob. of getting a head=p) for 'n' times (identical & independent Bernoulli trials).

$\searrow$ say (n>2)

Now define two r.v.s

  X : trial at which first head appears  (X takes values 1 to n)

  Y : no. of heads in the 'n' tosses $\uparrow$ (Y takes values 0 to n)

$\longrightarrow$ Note that X is not a valid r.v. as per the $\Big/$ defn. since there is "no trial id" which handles the case of all tails in n tosses. As a correcting factor lets ~~also~~ include a dummy value (say "0") which X takes on to represent the case of all tails.

$\longrightarrow$ Here is the marginal pmf of X :

$$f_X(x_i) = \begin{cases} (1-p)^n & x_i = 0 \ (\text{dummy value representing all tails}) \\ (1-p)^{x_i-1}\, p & x_i = 1 \text{ to } n \qquad \to \text{Now this is a valid pmf} \\ 0 & \text{otherwise} \end{cases}$$

marginal pmf of Y : $\to$ Binomial r.v.

$$f_Y(y_i) = \begin{cases} {}^nC_{y_i}\, p^{y_i} (1-p)^{n-y_i} & y_i = 0 \text{ to } n. \end{cases}$$

~~joint pmf of~~ lets write down the conditional prob. mass function of X | Y = i

(a) i=0 :  $f_{X/Y}(x_i/0) = \begin{cases} 1 & \text{if } x_i = 0 \qquad \to \text{valid pmf} \\ 0 & \text{otherwise} \end{cases}$

(b) i=1 :  $f_{X/Y}(x_i/1) = \begin{cases} \dfrac{f_{XY}(x_i,1)}{f_Y(1)} = \dfrac{(1-p)^{n-1} p}{{}^nC_1\, p(1-p)^{n-1}} = \dfrac{1}{n} & \text{if } x_i = 1 \text{ to } n \\ 0 & \text{otherwise} \end{cases}$

(3)

(c) $\underline{i=2}$ : $f_{X/Y}\left(x_i/2\right) = \begin{cases} \dfrac{^{n-x_i}C_1 \, p^2(1-p)^{n-2}}{^nC_2 \, p^2(1-p)^{n-2}} & \text{of } x_i = 1 \text{ to } n-1. \\ \\ 0 & \text{otherwise} \end{cases}$

$\rightarrow$ This is also a valid pmf.

no on...........

~~Now we be has already put down values of joint pmf for few and~~ ~~values~~

In the process of writing down the cond. pmf we have also put down the joint pmf for few value pair of $(x_i, y_i)$.

In the next lecture we will look at the case of cond.prob. for jointly conts. r.v. etc.

(c) $i=1$ : $f_{X/Y}\left(x_i/2\right)$

Suppose X, Y are jointly conts. r.v. Now we want to define notion of say cond. dist. and cond. pdf (if possible), moreover using the familiar notion of cond. prob. over events (which is very familiar to us). Now pdf has "no direct link" with prob. So may be its better to start from by defining cond. distr. function using cond. prob. over events.

Suppose we attempt the following: $F_{X|Y}\left(\dfrac{x}{y}\right) = P\left[X \le x \,/\, Y = y\right]$

note this was the defn for discrete r.v. case.

We are bound to fail since $P\{Y = y\} = 0 \;\forall\, y$. The work around is to define as follows: (which intuitively means the same!)

$$F_{X|Y}\left(\dfrac{x}{y}\right) = \lim_{\Delta y \downarrow 0} P\left[X \le x \,/\, y \le Y \le y + \Delta y\right]$$ ⟶ (This is the definition we go with)

$$= \lim_{\Delta y \downarrow 0} \dfrac{P\left[X \le x, \; y \le Y \le y + \Delta y\right]}{P\left[y \le Y \le y + \Delta y\right]}$$ ⟶ (familiar notion of cond. prob. over events)

$$= \lim_{\Delta y \downarrow 0} \dfrac{\displaystyle\int_{-\infty}^{x}\int_{y}^{y+\Delta y} f_{XY}(x', y')\, dy'\, dx'}{\displaystyle\int_{y}^{y+\Delta y} f_Y(y')\, dy'}$$ ⟶ (Since X, Y are jointly conts. there ∃ a pdf)

$$= \lim_{\Delta y \downarrow 0} \dfrac{\displaystyle\int_{-\infty}^{x} f_{XY}(x', y)\, \Delta y\, dx'}{f_Y(y)\,\Delta y}$$ ⟶ (over small intervals we assume pdf doos'nt change) so area is height × interval length

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

↓
this must be cond. pdf!!

$$\Rightarrow \boxed{\begin{array}{l} f_{X/Y}(x/y) = \dfrac{f_{XY}(x, y)}{f_Y(y)} \\ \text{(Conditional pdf)} \end{array}}$$

→ Note the similarity in the expression even in the discrete rv.s case!

↓

Now for this defn. is first of all valid if $y$ is such that $f_Y(y) \neq 0$. (i.e. prob. density of $Y$ at $y$ is non-zero).

Now ~~suppose Defn.~~ let us check $y$ for a fixed $y$, the $f_{X/Y}$ is indeed a pdf or not. (This check will complete the defn.)

① firstly $f_{X/Y}$ is $\geq 0$.

② $\displaystyle\int_{-\infty}^{\infty} f_{X/Y}(x/y)\, dx = \int_{-\infty}^{\infty} \frac{f_{XY}(x, y)}{f_Y(y)}\, dx = \frac{\int_{-\infty}^{\infty} f_{XY}(x, y)\, dx}{f_Y(y)} = \frac{f_Y(y)}{f_Y(y)} = 1$

Hence for a fixed value of $y$ such that $f_Y(y) \neq 0$, $f_{X/Y}$ is indeed a pdf and with different values of $y$ (satisfying $f_Y(y) \neq 0$) we get different pdf's!

Let look at an eg. given in one of prev. lectures!

eg 1. $f_{XY}(x,y) = \begin{cases} 24xy & x>0, y>0, \ x+y<1 \\ 0 & \text{otherwise} \end{cases}$

We already saw that

$f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x,y)\,dy = \begin{cases} 12x(1-x)^2 & 0<x<1 \\ 0 & \text{otherwise} \end{cases}$

III y $f_Y(y) = \begin{cases} 12y(1-y)^2 & 0<y<1 \\ 0 & \text{otherwise.} \end{cases}$

Now lets compute

$f_{X/Y}(x/y) = \begin{cases} \dfrac{24xy}{12y(1-y)^2} = \dfrac{2x}{(1-y)^2} & \text{if } 0<x<1-y \\ 0 & \text{otherwise} \end{cases}$

↓
for some $y \in (0,1)$
where I am
rule $f_Y(y) \neq 0$      (9 It is an easy exercise to check if 1 is valid pdf)

For eg. $f_{X/Y}(x/0.25) = \begin{cases} \dfrac{32x}{9} & 0<x<0.75 \\ 0 & \text{otherwise} \end{cases}$

$f_{X/Y}(x/0.75) = \begin{cases} 32x & 0<x<0.25 \\ 0 & \text{otherwise} \end{cases}$

no we can get a family of pdf using different values of $y \in (0,1)$.

Now lets take another eg of a joint's cts pdf we saw in last class

$$\text{eg} \quad f_X(\underline{x}) = f_{X_1 X_2 \ldots X_n}(x_1, x_2, \ldots, x_n) = \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2}\sum_{i=1}^{n} x_i^2} \qquad \underline{x} \in \mathbb{R}^n.$$

We also say that (i) each of $X_1, X_2, \ldots X_n$ have std. Normal dist.

(ii)
$$F_{X_1 X_2 \ldots X_n}(x_1, x_2, \ldots, x_n) = F_{X_1}(x_1) F_{X_2}(x_2) \cdots F_{X_n}(x_n)$$

(why?) $\Longleftrightarrow$

$$f_{X_1 X_2 \ldots X_n}(x_1, x_2, \ldots, x_n) = f_{X_1}(x_1) f_{X_2}(x_2) \cdots f_{X_n}(x_n)$$

Now lets take $n=2$, then calculate

$$f_{X_1 \backslash X_2}(x_1 / x_2) = \frac{f_{X_1 X_2}(x_1, x_2)}{f_{X_2}(x_2)} = \frac{f_{X_1}(x_1) f_{X_2}(x_2)}{f_{X_2}(x_2)} = f_X(x_1).$$

In other words knowledge abt. $X_2$ is not effect pdf of $X_1$ !

Similar
This notion was discussed as while discussing notion of independent events !

Lets formalize this notion of Independence:

$$\longrightarrow \qquad \underline{\text{In dependence of } RVs}$$

$X, Y$ are said to be independent

$\Longleftrightarrow \forall B_1, B_2 \in \mathcal{B}$ it happens that: $\qquad$ (I)

$$P[X \in B_1, Y \in B_2] = P[X \in B_1] P[Y \in B_2]$$

in other words the events
$$[X \in B_1] \quad \& \quad [Y \in B_2]$$
$$||| \qquad\qquad ||| $$
$$\{\omega \in \Omega / X(\omega) \in B_1\} \qquad \{\omega \in \Omega / Y(\omega) \in B_2\}$$
are independent

Now lets choose $B_1 = (-\infty, x]$ & $B_2 = (-\infty, y]$

then ① $\implies$ ~~(⊕ ⊗)~~ $F_{XY}(x,y) = F_X(x) F_Y(y)$

$\rightarrow$ This means the joint dint factorizes into marginals of in case of independent rvs the marginals completely determine the joint-dint.!!

(the converse is also true & beyond scope of this course)

$\longrightarrow$ Now if $X, Y$ are discrete rvs, then:

Take $B_1 = \{x\}$, $B_2 = \{y\}$ : ① $\implies f_{XY}(x,y) = f_X(x) f_Y(y)$

$\longrightarrow$ If $X, Y$ are jointly conts. rvs, then we anyway have:

$$F_{XY}(x,y) = F_X(x) F_Y(y)$$

$\implies \dfrac{\partial F_{XY}(x,y)}{\partial y} = \dfrac{F_X(x) d F_Y(y)}{dy} = F_X(x) f_Y(y)$

$\implies f_{XY}(x,y) = \dfrac{\partial^2 F_{XY}(x,y)}{\partial x \partial y} = \dfrac{d}{dx} F_X(x) f_Y(y) = f_X(x) f_Y(y).$

$\longrightarrow f_{XY}(x,y) = f_X(x) f_Y(y).$

So joint . functions, pmf, pdf's factorize !!

We ~~can extend~~
Also for independent rvs : $f_{X|Y}(x/y) = \dfrac{f_{XY}(x,y)}{f_Y(y)} = \dfrac{f_X(x) f_Y(y)}{f_Y(y)}$

$\underbrace{\qquad}_{\text{conditional}} = \underbrace{\qquad}_{\text{marginal !}}$

$= f_X(x)$

The notion of independence of r.v's can be extended to any $X_1, X_2, \ldots, X_n$:

We say $X_1, \ldots, X_n$ are independent if for all sub-collections $X_i \ldots X_j$ are independent (in other words $X_i X_j$ are independent, $X_i X_j X_k$ are independent, so on $\ldots$) and $F_{X_1 X_2 \ldots X_n}(x_1, x_2, \ldots, x_n)$ factorizes into marginals.

(So there are $2^n - n - 1$ conditions to be checked!)

again One of the conditions is $F_{X_1 X_2 \ldots X_n}(x_1, x_2, \ldots, x_n)$
$$= F_{X_1}(x_1) F_{X_2}(x_2) \cdots F_{X_n}(x_n).$$

Note that the r.v's in eg 1 are not independent and those in eg 2 are independent r.v's!

Moreover the r.v's in eg 2 are also identically distributed

In other words, each of the $X_1, X_2, \ldots, X_n$ has the same std. Normal distribution.

Such a collection of r.v's which are independent and are identically distributed are called as iid r.v's.

Let's look at a quick eg:

# Suppose $X, Y$ are iid. & conts. r.v's.

Calculate $P[X > Y]$.

Intuitive answer is:

$1 = P\{(X,Y) \in \mathbb{R}^2\} = P[X > Y \cup X < Y \cup X = Y]$

$= P\{X > Y\} + P\{X < Y\} + P\{X = Y\}$

↓ each are mutually exclusive events

$\longrightarrow P\{X = x, Y = x\}$

↓ zero since $X, Y$ are conts. r.v's!

Now there is no reason to believe $X > Y$ or $Y < X$ (since they are independent & imitations of same distribution) so $P\{X > Y\}$ must be equal to $P\{Y < X\}$

$\Rightarrow P[X > Y] = \frac{1}{2}$

→ nay all are iid.

||| A similar argument shows $P\{X_1 > X_2 > \cdots > X_n\} = \frac{1}{n!}$

because this is just one ordering among all $n!$ orderings !!

⟶ Now lets look at a more rigorous answer and we will realize that our intuition is right!

$P[X > Y] = \int_{-\infty}^{\infty} \int_{-\infty}^{x} f_{XY}(x,y) \, dy \, dx = \int_{-\infty}^{\infty} \int_{-\infty}^{x} f_X(x) f_Y(y) \, dy \, dx$

$X, Y$ are independent

⑦

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{x} f_X(x) f_X(y) \, dy \, dx \quad \left( \because X, y \text{ are identically distributed} \right)$$

$$= \int_{-\infty}^{\infty} f_X(x) \left[ \int_{-\infty}^{x} f_X(y) \, dy \right] dx \longrightarrow F_X(x) \, !$$

$$= \int_{-\infty}^{\infty} f_X(x) \, F_X(x) \, dx$$

(two ways of computing it)

(by parts method)

$$= F_X(x) \, F_X(x) \Big|_{-\infty}^{\infty} - \int_{-\infty}^{\infty} F_X(x) f_X(x) \, dx$$

$$\Rightarrow \int_{-\infty}^{\infty} f_X(x) \, dF_X(x) \, dx = \frac{F_X(x) \, F_X(x) \Big|_{-\infty}^{\infty}}{2}$$

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

(substitution method)

Put $x \longrightarrow F_X(x)$

$$= \int_{0}^{1} F_X(x) \, d \, F_X(x) = \frac{F_X(x)^2}{2} \Big|_{0}^{1} = \frac{1}{2}$$

$\longrightarrow$

## BAYE's Theorem (extension to r.v).

Let $X, Y$ are discrete / jointly conts. r.v's

Let $f_{XY}, f_X, f_Y$ represent their joint & marginals pmf/pdf. (whichever is the case)

We know:
$$f_{X|Y}(x/y) = \frac{f_{XY}(x,y)}{f_Y(y)} \Rightarrow f_{XY}(x,y) = f_{X|Y}(x/y) f_Y(y)$$

also, $$f_{Y|X}(y/x) = \frac{f_{XY}(x,y)}{f_X(x)} = \frac{f_{X|Y}(x/y) f_Y(y)}{f_X(x)}$$

Now if X, Y are discrete r.vs, then:

$$f_X(x) = \sum_{\forall y'} f_{XY}(x, y') = \sum_{\forall y'} f_{X/Y}(x/y') f_Y(y')$$

Now substituting back we get:

$$f_{Y/X}(y/x) = \frac{f_{X/Y}(x/y) f_Y(y)}{\sum_{\forall y'} f_{X/Y}(x/y') f_Y(y')}$$   } Baye's theorem for discrete r.vs.

If X, Y are jointly conts., then:

$$f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x, y') dy' = \int_{-\infty}^{\infty} f_{X/Y}(x/y') f_Y(y') dy'$$

Substituting back we get:

$$f_{Y/X}(y/x) = \frac{f_{X/Y}(x/y) f_Y(y)}{\int_{-\infty}^{\infty} f_{X/Y}(x/y') f_Y(y') dy'}$$   } Bayes theorem for jointly conts r.vs.

The advantage ~~In the Past class~~ is that without knowing joint pmf ~~step~~ of joint pdf we are able to calculate cond. prob on one side using cond. prob on the other side i.e. $f_{X/Y}$

$^{+}f_{Y/X}$
and ofcourse we also are using $f_Y(y)$ or $f_X(x)$ i.e. marginals.

∴ Use wherever joint -pmf/pdf is difficult to compute !!
(unecessarly)

This lecture begins with some applications of the Bayes's theorem

eg1 Suppose we are told that a person picks up a coin at random from a set of 'm' coins. It is also given to us that that the prob. of picking the $i^{th}$ coin is $q_i$ (i.e. $q_i \geq 0$, $\sum_{i=1}^{m} q_i = 1$). Now just given this information suppose we are asked to guess what coin was picked up by the person, intuitively our answer would be $argmax\ q_i$. i.e. the coin which has the maximum prob. of being picked. Note that we are using absolutely no information regarding the particular coin picked up but some "generic" information about $q_i$.

Now suppose the person is generous to reveal some partial information regarding the coin in his hand and then asks us to guess. In particular, suppose he reveals the number of heads he got by tossing the coin in his hand for 'n' times and he also reveals the prob. of getting heads with each of the coin (i.e. $p_i \geq 0$, $p_i \leq 1$, $i=1$ to $m$).

Now, a little bit of thinking will show that given the partial information, we can come up with a better "guess". (Think abt two extreme cases where all $q_i$ are same & $q_i$ highly distinct)

Let us formalize our ideas: define two rvs

$X$: # heads in 'n' tosses (with the coin picked up)
$\longrightarrow \in \{0,1,2,\ldots,n\}$

$Y$: 'id' of the coin picked up $\longrightarrow \in \{1,2,\ldots,m\}$

Now ~~Since~~ pmf of $Y$ is given: $f_Y(y) = \begin{cases} q_y & y \in \{1,2,\ldots,m\} \\ 0 & \text{otherwise} \end{cases}$

~~Since there is an information given prior to the partial~~ ~~regarding the coin picked up, this configuration is~~

Also, the following cond. pmf is given:

$$f_{X/Y}\left(\frac{x}{y}\right) = \begin{cases} {}^{n}C_x \, P_y^x \, (1-P_y)^{n-x} & x \in \{0,1,\ldots,n\} \\ 0 & \text{otherwise} \end{cases}$$

$y \in \{1,2,\ldots,m\}$

Now we wish to ~~end~~ calculate the prob. that the coin picked up is 'y' given that there were 'x' heads:

$$f_{Y/X}\left(\frac{y}{x}\right) = \frac{f_{X/Y}\left(\frac{x}{y}\right) f_Y(y)}{\sum\limits_{\forall y'} f_{X/Y}\left(\frac{x}{y'}\right) f_Y(y')} \qquad (\because \text{Bayes theorem})$$

$x \in \{0,1,\ldots,m\}$

$$= \frac{{}^{n}C_x \, P_y^x \, (1-P_y)^{n-x} \, q_y}{\sum\limits_{\forall y'} {}^{n}C_x \, P_{y'}^x \, (1-P_{y'})^{n-x} \, q_{y'}}$$

Now given a 'x' by the person, for different values of $y$ we can calculate this quantity. Again (intuitively) the guess is to pick the coin which maximizes $f_{y/x}$!

Note that if all $V_y = \frac{1}{m}$ (for $y=1$ to $m$), then $f_{y/x}$ is same as $f_y$ (and our guess would'nt change!) In other words, if all coins have same prob. of getting heads, the partial information regarding the coin picked up does not give any ~~more~~ additional insight to it! $\rightarrow$ i.e. no. heads in 'n' tosses
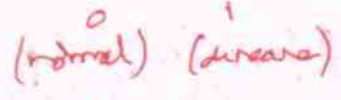
Since $f_{y/x}$ is prob. of ~~picking~~ a coin being picked after looking ~~after~~ at the partial information abt $X$, it is usually called as posterior prob. & $f_y$ is called as prior probability. (The same idea is extensively used in prob. model known as ~~How to Bum through the~~ Hidden Markov model (HMM) which are ⊗ state of the art models for automated speech recognition systems!)

Hence the idea of cond. prob. & Baye's theorem have far reaching consequences.

Now lets look at another eg. illustrating the utility of the Baye's theorem. The reader is encouraged to see ⊗ at every step the analogy between these two examples.

eg2 Suppose X represents the diagnostic report of a patient ④
(for the sake of simplicity assume it represents the body temp.
of the patient) and suppose Y represents whether he has a
disease or not.

(normal) (disease)

Now (again) the task is to predict (guess) whether
a patient has disease or not! Assume the following
information is given:

① $f_Y$ (pmf) is given. (This is the prior information). In words, the
fraction of normal & diseased people in a population is given.

Ⓘ

② $f_{X/Y}$ (pdf) is given. In words, the body temp. distribution
of normal as well as of patients with disease are given. Note
that X/Y is a conts. r.v. $\Rightarrow F_{X/Y}(x/y) \equiv P\{X \le x/y=y\} = \int_{-\infty}^{x} f_{X/Y}(x'/y)dx'$.

Again we wish to compute $f_{Y/X}$ (pmf) which in words
is the prob. of the patient is normal or has disease given his
diagnostic report (body temp.)

Before this let us answer a simpler question "what is $f_X(x)$?".
In words, what is the body temp. dist. of the entire
population? Since we have not assumed anything abt X, conts. or
discrete,
let us compute its dist. function:

$$F_X(x) = P[X \le x] = \sum_{\forall y} P[X \le x, Y = y] \qquad (\because \text{marginal fundn})$$

$$= \sum_{\forall y} P[X \le x / Y = y] \, P[Y = y] \qquad (\because \text{cond. prob.})$$

$$= \sum_{\forall y} F_{X/Y}(x/y) \, f_Y(y) \qquad (\because \text{defn. of cond. distn.})$$

$$= \sum_{\forall y} \int_{-\infty}^{x} f_{X/Y}(x'/y) \, dx' \, f_Y(y) \qquad \left(\because \frac{X/Y \text{ is a}}{\text{conts rv}}\right)$$

$$= \int_{-\infty}^{x} \left[ \sum_{\forall y} f_{X/Y}(x'/y) \, f_Y(y) \right] dx' \qquad \left(\because \frac{\text{interchange}}{\int \ \& \ \Sigma}\right)$$

$$\Rightarrow \quad f_X(x) = \sum_{\forall y} f_{X/Y}(x/y) \, f_Y(y) \qquad (\because \text{defn. of } \int \text{pdf})$$

Recall that this ↓ resembles the total prob. rule (for the case X, Y are both discrete). However note that $f_X$ and $f_{X/Y}$ are pdf's and $f_Y$ is a pmf. Also, it looks like $f_X$ is a convex combination of conditional pdf's ($f_{X/Y}$). In other words, it looks like X is a 'mixture' of two kinds of rvs ($X/Y=0$, $X/Y=1$ here!). $f_X$ is also sometimes called as mixture density. $f_Y$ are called as mixing prob. & $f_{X/Y}$ as class conditional density!

Models satisfying ① ② in Ⓘ here are called as Mixture Models.

Now lets try to compute:

$$f_{Y/X}(y/x) \equiv \lim_{\Delta x \downarrow 0} P\left[Y = y / X \in (x, x+\Delta x]\right]$$

(III $\downarrow$ to defn. in case of X, Y conts. r.v.)

$$\downarrow \atop \text{pmf for fixed } x \atop \text{such that } f_X(x) \neq 0} = \lim_{\Delta x \downarrow 0} \frac{P[x \leq X \leq x+\Delta x, Y=y]}{P[x \leq X \leq x+\Delta x]}$$

$$= \lim_{\Delta x \downarrow 0} \frac{P[x \leq X \leq x+\Delta x / Y=y] P[Y=y]}{P[x \leq X \leq x+\Delta x]}$$

$$= \lim_{\Delta x \downarrow 0} \frac{\int_{x}^{x+\Delta x} f_{X/Y}(x'/y)\, dx' \; f_Y(y)}{\int_{x}^{x+\Delta x} f_X(x')\, dx'}$$

($\because$ X/y and therefore X are conts r.v's)

$$= \lim_{\Delta x \downarrow 0} \frac{f_{X/Y}(x/y)\,\Delta x\, f_Y(y)}{f_X(x)\,\Delta x}$$

$$\Rightarrow f_{Y/X}(y/x) = \frac{f_{X/Y}(x/y)\, f_Y(y)}{f_X(x)} = \frac{f_{X/Y}(x/y)\, f_Y(y)}{\sum_{\forall y'} f_{X/Y}(x/y')\, f_Y(y')}$$

Again, this looks like Baye's theorem in case of X, Y both discrete r.v's! But here $f_{Y/Y}$ and $f_Y$ are ~~pdfs~~ pmfs wheleas $f_{X/Y}$ is a pdf!

   Again $f_{Y/X}$ is posterior pmf (pmf after looking at partial information i.e. diagnostic report of the patient)

   $f_Y$ is simply the prior information.

In general, there are two ways to guess:

$$\arg\max_y f_Y(y) \qquad \rightarrow \text{y that maximizes the prior prob. (without looking at the particular patient!)}$$

$$\arg\max_y f_{Y|X}(y/x) \qquad \rightarrow \text{given diagnostic report } x \text{ of patient guess his status of health}$$

max. priorprob.

max. posterior prob. → it is easy to reason out this gives a better guess.

———→

Suppose we want to design a chair which withstands the weight of people who sit on it as well as is not built from too costly or heavy (strong) material! One way is to design it for the heaviest person on earth. But this is too pessimistic and will lead to a chair perhaps too heavy to even move :) On the other hand we want chair to be 'robust' enough to handle heavy people.

One way to put this is to design chair such that it withstands the weight of any 100 random people who sit on it. i.e. design for:

$$M = \max \{X_1, X_2, X_3, \ldots, X_n\}$$

Here $X_1, \ldots, X_n$ represents weights of $n$ people. It is easy to see that they are $\underline{i.i.d.}$ r.v.s (why ??).

Also note that $M$ is a r.v and is infact a function of collection of r.v.s $X_1, \ldots, X_n$!

Let compute the d.f of $M$:

$$F_M(x) = P[M \leq x] = P[\max \{X_1, X_2, \ldots, X_n\} \leq x]$$

$$= P[X_1 \leq x, X_2 \leq x, \ldots, X_n \leq x]$$

$$= P[X_1 \leq x] P[X_2 \leq x] \cdots P[X_n \leq x] \quad \left( \because \text{ they are indyndet.} \right)$$

$$= F_{X_1}(x) F_{X_2}(x) \cdots F_{X_n}(x)$$

$$= (F(x))^n \qquad \left( \text{here } F \text{ is the common, dist. function of all the identically distributed r.v.s } X_1, \ldots X_n \right)$$

$$\Rightarrow f_M(x) = n (F(x))^{n-1} f(x) \qquad \left( \text{here } f \text{ is the common pdf of } X_1, \ldots X_n \right)$$

In words $F(t)$ represent the dist. of body weight among humans!

$\mathrm{III}_b$ one can consider:

$$N = \min\{X_1, X_2, \ldots, X_n\}.$$ $N$ is another function of collections of r.vs.

$$F_N(x) = P[N \le x] = P[\min\{X_1, X_2, \ldots, X_n\} \le x]$$

$$= 1 - P[\min\{X_1, X_2, \ldots, X_n\} > x]$$

$$= 1 - P[X_1 > x, X_2 > x, \ldots, X_n > x]$$

$$= 1 - \left(1 - F(x)\right)^n \qquad (\text{again by } \underline{iid})$$

$$\Rightarrow f_N(x) = n\left(1 - F(x)\right)^{n-1} f(x).$$

Now one can consider the joint dist. of collection of $M, N$ r.vs which are in turn function of collections of r.vs.

$$F_{MN}(x,y) = P[M \le x, N \le y]$$

$$= P[M \le x] - P[M \le x, N > y]$$

$$= \left(F(x)\right)^n - P[y < X_1 \le x, y < X_2 \le x, \ldots, y < X_n \le x]$$

$$= \left(F(x)\right)^n - \left(F(x) - F(y)\right)^n \qquad (\because iid \text{ arguments})$$

In future classes we will look into more eg. of functions of r.vs.

This lecture formalizes the notion of function of multivariate r.v.s. (in other words functions of collections of r.v.s).

Suppose $X_1, X_2, \ldots, X_n$ are r.v.s defined on $\mathbb{P} = (\Omega, \mathcal{F}, P)$. Also, $g: \mathbb{R}^n \to \mathbb{R}$ is given. Consider a new function $Z: \Omega \to \mathbb{R}$ defined as

$$Z(\omega) = g(X_1(\omega), X_2(\omega), \ldots, X_n(\omega)) \quad \forall \omega \in \Omega.$$

The short hand representation of $\downarrow$ is $Z = g(X_1, X_2, \ldots, X_n)$.

Now $Z$ is indeed a function from $\Omega \to \mathbb{R}$, so it is a r.v if:

$$Z^{-1}(B) \in \mathcal{F} \quad \forall B \in \mathcal{B}.$$

i.e.

$$\Rightarrow \{\omega \in \Omega \mid g(X_1(\omega), X_2(\omega), \ldots, X_n(\omega)) \in B\} \in \mathcal{F} \quad \forall B \in \mathcal{B}.$$

Consider the condition $g^{-1}(B) \in \mathcal{B}^n$ $\overset{\forall B \in \mathcal{B}}{\text{ in other words }}$ $g^{-1}(B) = B_1 \times B_2 \times \ldots \times B_n$ each $B_i \in \mathcal{B}$.

It is easy to see that $Z^{-1}(B) = \bigcap_{i=1}^{n} \{\omega \in \Omega \mid X_i(\omega) \in B_i\} \in \mathcal{F}$

since each $X_i$ is a valid r.v!

Hence $g^{-1}(B) \in \mathcal{B}^n$ is the condition on $g$ which makes $Z$ a valid r.v.

Again (not in this class) we can show that if $g$ is a continuous function then $Z = g(X_1, \ldots, X_n)$ is also a valid r.v.

Lets consider an example:

Ⓐ

**eg1** $Z = X + Y$ $\qquad$ (here $g(x,y) = x+y$)

Suppose joint pdf of $X, Y$ is known. Compute $f_Z$.

$$F_Z(z) = P[Z \le z] = P[X+Y \le z]$$

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

Now $f_Z(z) = \dfrac{d F_Z(z)}{dz} = \int_{-\infty}^{\infty} \dfrac{d}{dz} \int_{-\infty}^{z-x} f_{XY}(x,y)\, dy\, dx$

$$= \int_{-\infty}^{\infty} f_{XY}(x, z-x)\, dx$$

Assume now that $X, Y$ are independent rvs.

$$\Rightarrow f_Z(z) = \int_{-\infty}^{\infty} f_X(x)\, f_Y(z-x)\, dx \quad \left\} \begin{array}{l} \text{nothing but} \\ \text{convolution of } f_X, f_y \, dz \end{array} \right.$$

$$= f_X(z) * f_Y(z)$$

Hence, pdf of sum of two rv's is the convolution of the individual pdfs!

Lets now take the special case $f_X(x) = f_Y(x) = \begin{cases} \frac{1}{2} & -1 < x < 1 \\ 0 & \text{otherw} \end{cases}$

In other words we are assuming $X, Y$ are both $\sim U[-1, 1]$
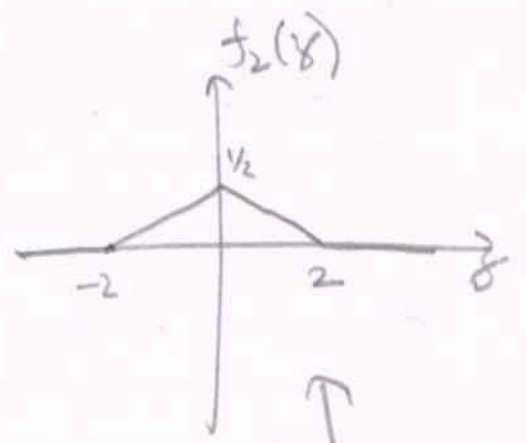(uniform rvs in $[-1, 1]$).

(2)

$$f_Z(z) = f_X(z) * f_Y(z)$$

(X, Y are iid and are uniform r.v. in $(-1, 1]$)

$$= \int_{-\infty}^{\infty} f_X(x) \, f_Y(z-x) \, dx$$
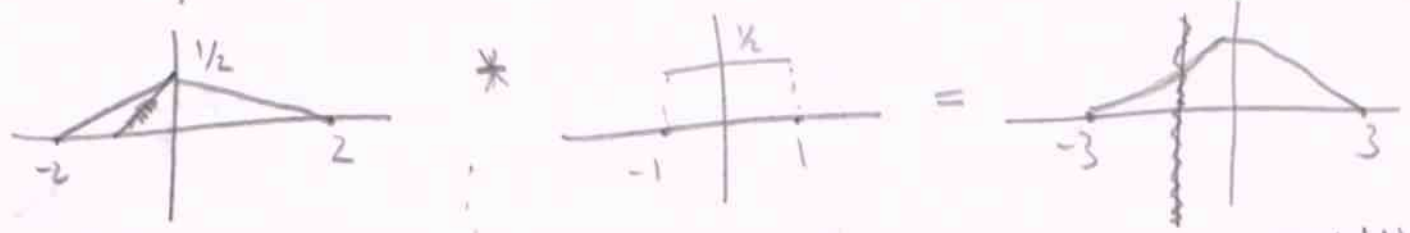
$f_X(x)$ and $f_Y(z-x)$ are non zero iff
$-1 < x < 1$
$-1 < z-x < 1$
i.e. $-1 < x < 1$
$z-1 < x < z+1$

$$= \begin{cases} 0 & z \le -2 \\ \int_{-1}^{z+1} \frac{1}{4} \, dx & -2 < z \le 0 \\ \int_{z-1}^{1} \frac{1}{4} \, dx & 0 < z \le 2 \\ 0 & z > 2 \end{cases}$$

$$= \begin{cases} 0 & z \le -2 \\ \dfrac{z+2}{4} & -2 < z \le 0 \\ \dfrac{2-z}{4} & 0 < z \le 2 \\ 0 & z > 2 \end{cases}$$



So sum of two iid $U[-1,1]$ r.v.s is not $U[-1,1]$ but is flat.
Similarly we can sum three, four, ... $\overset{iid}{}$ r.v's which are $U[-1,1]$:



(infinite sum)

Std. Normal r.v.!

This is an intuition for a special case of Central Limit Theorem.

(Here Sum of infinite iid rvs all $U(-1,1]$ is converging to std. Normal rv.)

$\underrightarrow{\text{eg2}}$ Let $X, Y$ be iid and ~~be dist~~ be Normal rvs.

by above argument : $Z = X + Y$ has the pdf as convolution of pdfs of $X, Y$. It is a well-known result that convolution of any two Gaussian fuctions is a Gaussian fuction.

Using this result we can say $Z$ is again a Normal rv !

(we will see a generic result of this kind later)

(A Gaussian fuction is any fuction of the form:
$$f(x) = a \, e^{-(x-b)^2/c^2}$$
$a, c > 0$. Note that the Normal rv has a pdf as a Gaussian fuc. Hence Normal rv are also known as gaussian rvs!)

$\underrightarrow{\text{eg3}}$ $Z = \dfrac{X}{Y}$ .

$$F_Z(z) = P\{Z \le z\} = P\left[\frac{X}{Y} \le z\right]$$

$$= P\{X \le zY, Y > 0\} + P\{X \ge zY, Y < 0\}$$

$$= \int_0^\infty \int_{-\infty}^{zy} f_{XY}(x, y) \, dx \, dy + \int_{-\infty}^{0} \int_{zy}^{\infty} f_{XY}(x, y) \, dx \, dy$$

$$\Rightarrow f_Z(\gamma) = \frac{dF_Z(\gamma)}{d\gamma} = \int_0^\infty \frac{d}{d\gamma} \int_{-\infty}^{\gamma y} f_{XY}(x,y)\,dx\,dy + \int_{-\infty}^0 \frac{d}{d\gamma} \int_{\gamma y}^0 f_{XY}(x,y)\,dx\,dy$$

$$= \int_0^\infty y\, f_{XY}(\gamma y, y)\,dy + \int_{-\infty}^0 -y\, f_{XY}(\gamma y, y)\,dy$$

$$= \int_{-\infty}^\infty |y|\, f_{XY}(\gamma y, y)\,dy$$

→ Now lets take $X, Y$ as iid and std. Normal rvs.

then
$$f_Z(\gamma) = \int_{-\infty}^\infty |y| \frac{e^{-\frac12 \gamma^2 y^2}}{\sqrt{2\pi}} \frac{1}{\sqrt{2\pi}} e^{-\frac12 y^2}\,dy = \frac{1}{\pi}\int_0^\infty y\, e^{-(\gamma^2+1)\frac{y^2}{2}}\,dy$$

$$= \frac{1}{\pi}\left[\frac{e^{-(\gamma^2+1)t}}{-(\gamma^2+1)}\right]_0^\infty$$

$$= \frac{1}{\pi(1+\gamma^2)}$$

∴ $Z$ is a Cauchy rv !

→ Now, Consider the collection of rv's $Z_1, Z_2, \ldots, Z_n$ each of which are in turn functions of the rvs: $X_1, X_2, \ldots, X_n$.

i.e. Consider
$$Z_1 = g_1(X_1, X_2, \ldots, X_n)$$
$$Z_2 = g_2(X_1, X_2, \ldots, X_n)$$
$$\vdots$$
$$Z_n = g_n(X_1, X_2, \ldots, X_n)$$

we already saw that each $Z_i$ is a rv (from the name initial $\mathbb{P}$). Here the collection of $\{Z_1, Z_2 \ldots Z_n\}$ is indeed a valid multivariate rv. Here we can talk about its dist. fnc. & equivalengtly, the joint distribution of $Z_1, Z_2, \ldots, Z_n$ which are collections of fuctions of $X_1, X_2, \ldots, X_n$ rvs.

$$F_Z(\underline{z}) = F_{Z_1 Z_2 \ldots Z_n}(z_1, z_2, z_3, \ldots, z_n)$$

$$= P[\, Z \leq \underline{z}\,]$$

$$= P[\, Z \in B\,] \qquad \text{where } B = (-\infty, z_1] \times (-\infty, z_2] \times \ldots \times (-\infty, z_n)$$

Now $\quad Z = \underline{g}(x) \qquad$ where $g : \mathbb{R}^n \to \mathbb{R}^n$ &

$$\underline{g}(x_1, x_2, \ldots, x_n) = \Big(g_1(x_1, x_2, \ldots, x_n), \; g_2(x_1, \ldots, x_n), \; \ldots, \; g_n(x_1, \ldots, x_n)\Big)$$

$$\Rightarrow F_Z(\underline{z}) = P[\, \underline{g}(x) \in B\,]$$

$$= P[\, X \in \underline{g}^{-1}(B)\,]$$

$$= \int\int \cdots \int_{\underline{g}^{-1}(B)} f_{X_1 X_2 \ldots X_n}(x_1, x_2, \ldots, x_n)\, dx_1\, dx_2 \ldots dx_n$$

Now suppose ~~that~~ $\underline{g}$ ~~is~~ ~~do the follow~~ invertible:

$$\Rightarrow \exists \underline{h} \ni X = \underline{h}(z) .$$

Also suppose $\underline{g}, \underline{h}$ are continuously differentiable.

Then

$$F_{Z}(\delta) = \int\int \cdots \int_{\underline{g}^{-1}(0)} f_{x_1 x_2 \cdots x_n}(x_1, x_2, \ldots, x_n) \, dx_1 dx_1 \cdots dx_n$$

$$= \int\int \cdots \int_{0} f_{x_1 \cdots x_n}\left( h_1(\gamma_1 \gamma_2 \cdots \gamma_n), h_2(\gamma_1 \cdots \gamma_n), \ldots, h_n(\gamma_1 \cdots \gamma_n) \right)$$
$$|J| \, d\gamma_1 d\gamma_2 \cdots d\gamma_n$$

$$\longrightarrow \text{abs value of the Jacobian.}$$

$$\Rightarrow f_Z(\gamma) = f_{x_1 \cdots x_n}\left( h_1(\gamma_1 \cdots \gamma_n), \ldots, h_n(\gamma_1 \cdots \gamma_n) \right) |J|$$

$$|J| = \text{abs of det. of} \begin{bmatrix} \frac{\partial h_1}{\partial \gamma_1} & \cdots & \frac{\partial h_1}{\partial \gamma_n} \\ \vdots & & \\ \frac{\partial h_n}{\partial \gamma_1} & \cdots & \frac{\partial h_n}{\partial \gamma_n} \end{bmatrix} .$$

We will see an explanation in the next lecture.

We will continue the discussion at end previous lecture (now restricting ourselves to 2-d case):

Consider two rvs: $Z_1 = g_1(X_1, X_2)$          $g_1 : \mathbb{R}^2 \to \mathbb{R}$

$\qquad\qquad\qquad Z_2 = g_2(X_1, X_2)$          $g_2 : \mathbb{R}^2 \to \mathbb{R}$

Let $Z$ be the multivariate rv representing $Z_1, Z_2$
$\quad X$ "                                   "            $X_1, X_2$

Also let, $\underline{g} : \mathbb{R}^2 \to \mathbb{R}^2$ be defined as $\underline{g}(x,y) = (g_1(x,y), g_2(x,y))$

It is easy to see that $Z = \underline{g}(x)$.


Now assume:

① $\underline{g}$ is invertible ($\underline{g}$ is bijection). Let $\underline{h} = \underline{g}^{-1}$.

it is easy to see $X = \underline{h}(Z)$. Also let $\underline{h}(z_1, z_2) = (h_1(z_1, z_2), h_2(z_1, z_2))$

② Assume $\underline{g}, \underline{h}$ are continuously differentiable.

③ Assume $Z, X$ are conts. (multivariate) rvs.

We wish to write down the joint pdf of $Z_1, Z_2$ (i.e. pdf of $Z$) in terms of joint-pdf of $X_1, X_2$ (ie, pdf of $X$). To this end:

$$F_{\underline{Z}}(\underline{z}) = P[Z \le \underline{z}]$$
$$= P[Z \in B]$$
where $B = (-\infty, z_1] \times (-\infty, z_2]$.

$$\Rightarrow F_2(\underline{\gamma}) = P[\underline{g}(x) \in B]$$

$$= P[X \in \underline{h}(B)]$$

$$= \iint_{\underline{h}(B)} f_{X_1 X_2}(\nsim_1 y_2) \, dx_1 \, dx_2 \qquad \boxed{I}$$

Now suppose I do change of dummy variables $x_1, x_2$ in the double integral:

$$x_1 = h_1(\gamma_1, \gamma_2)$$
$$x_2 = h_2(\gamma_1, \gamma_2)$$

$$\left( \begin{array}{l} \text{remember that} \\ h_1 = g_1^{-1} \\ h_2 = g_2^{-1} \end{array} \right)$$

Now $(x_1, x_2) \in \underline{h}(B)$  from the integral limits

$$\Rightarrow \left( h_1(\gamma_1, \gamma_2), h_2(\gamma_1, \gamma_2) \right) \in \underline{h}(B)$$

$$\Rightarrow \underline{h}(\gamma_1, \gamma_2) \in \underline{h}(B) \Rightarrow (\gamma_1, \gamma_2) \in B \qquad \boxed{II}$$

Now in order to proceed with change of variables I need to figure out how elementary area in $\gamma_1, \gamma_2$ coordinate looks like! Before doing that any elementary area is $|J| d\gamma_1 d\gamma_2$ (we will shortly show what is $|J|$)

Then:

$$F_2(\underline{\gamma}) = \iint_B \boxed{f_{X_1 X_2}(h_1(\gamma_1, \gamma_2), h_2(\gamma_1, \gamma_2)) |J| \, d\gamma_1 \, d\gamma_2}$$

from $\boxed{II}$  $\longleftarrow B$   $\longrightarrow$ must be $f_{2, z_2}$ !! (why?)  ②

Since dist. of $Z$ can be computed by integrating a function (over relevant 2-d interval), that same function must be the pdf of $Z$!
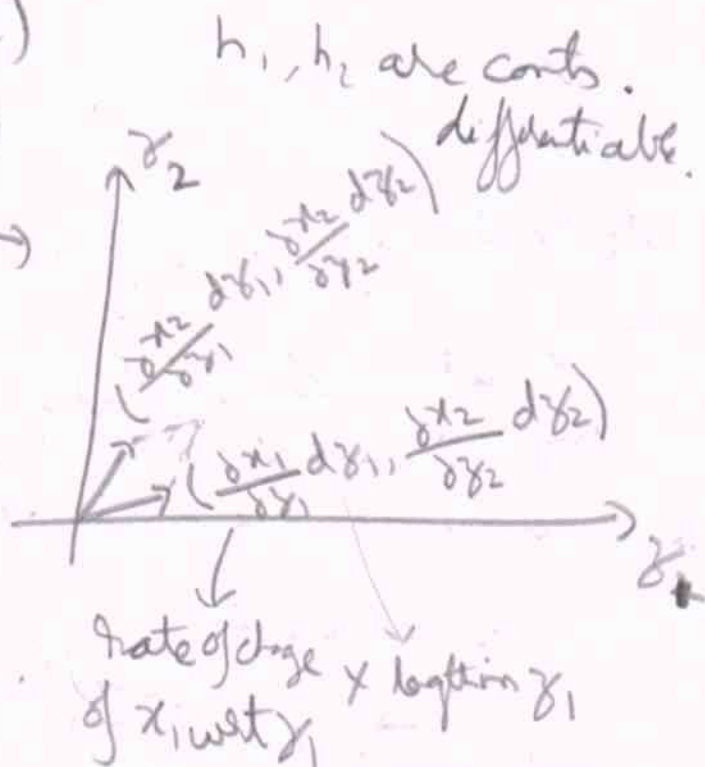
$$\Rightarrow \boxed{f_{Z_1 Z_2}(\gamma_1, \gamma_2) = f_{X_1 X_2}\left(h_1(\gamma_1, \gamma_2), h_2(\gamma_1, \gamma_2)\right) |J|}$$

Hence we are successful in the derivation. Now let us see how $|J|$ can be computed as:

## Change of Variables in Multiple integrals

$$x_1 = h_1(\gamma_1, \gamma_2)$$
$$x_2 = h_2(\gamma_1, \gamma_2)$$

$h_1, h_2$ are conts. differentiable.



Area $= dx_1 dx_2$

rate of change of $x_1$ w.r.t $\gamma_1$ × length in $\gamma_1$

Area in $\gamma_1, \gamma_2$ coordinates is area of parallelogram (for which we know the vectors of sides!)

(3)

Area of llgm is nothing but cross-product of the vectors of base sides:

$$\text{Area vector} = \begin{vmatrix} \hat{i} & \hat{j} & \hat{k} \\ \frac{\partial x_1}{\partial \gamma_1}d\gamma_1 & \frac{\partial x_1}{\partial \gamma_2}d\gamma_2 & 0 \\ \frac{\partial x_2}{\partial \gamma_1}d\gamma_1 & \frac{\partial x_2}{\partial \gamma_2}d\gamma_2 & 0 \end{vmatrix}$$

$$= \left( \frac{\partial x_1}{\partial \gamma_1}d\gamma_1 \frac{\partial x_2}{\partial \gamma_2}d\gamma_2 - \frac{\partial x_1}{\partial \gamma_2}d\gamma_2 \frac{\partial x_2}{\partial \gamma_1}d\gamma_1 \right) \hat{k}$$

$$\text{Area} = \left| \nearrow \right| = \left| \frac{\partial x_1}{\partial \gamma_1} \frac{\partial x_2}{\partial \gamma_2} - \frac{\partial x_2}{\partial \gamma_1} \frac{\partial x_1}{\partial \gamma_2} \right| d\gamma_1 d\gamma_2$$

is called Jacobian and denoted by $|J|$

Note that,

$|J|$ is also abs. of det. of $\longrightarrow$ $\begin{bmatrix} \frac{\partial x_1}{\partial \gamma_1} & \frac{\partial x_1}{\partial \gamma_2} \\ \frac{\partial x_2}{\partial \gamma_1} & \frac{\partial x_2}{\partial \gamma_2} \end{bmatrix}$

Called as Jacobian matrix

In n-dimensional case:

Jacobian matrix $= \begin{bmatrix} \frac{\partial x_1}{\partial \gamma_1} & \frac{\partial x_1}{\partial \gamma_2} & \cdots & \frac{\partial x_1}{\partial \gamma_n} \\ \frac{\partial x_2}{\partial \gamma_1} & \frac{\partial x_2}{\partial \gamma_2} & \cdots & \frac{\partial x_2}{\partial \gamma_n} \\ \frac{\partial x_n}{\partial \gamma_1} & \frac{\partial x_n}{\partial \gamma_2} & \cdots & \frac{\partial x_n}{\partial \gamma_n} \end{bmatrix}$

④

Lets take an eg and work out details:

Eg1     Let    $Z_1 = X + Y$
           $Z_2 = X - Y$    joint pdf of $X, Y$ is given

Compute joint pdf of $Z_1, Z_2$.

We know,

$$f_{Z_1 Z_2}(\gamma_1, \gamma_2) = f_{X_1 X_2}\left(h_1(\gamma_1, \gamma_2), h_2(\gamma_1, \gamma_2)\right) |J|$$

first

i) we need to figure out what are $h_1, h_2$:

i.e. express $X, Y$ in terms of $Z_1, Z_2$:

$$X = \frac{Z_1 + Z_2}{2}$$

$$Y = \frac{Z_1 - Z_2}{2}$$

$$\implies$$

$$h_1(\gamma_1, \gamma_2) = \frac{\gamma_1 + \gamma_2}{2}$$

$$h_2(\gamma_1, \gamma_2) = \frac{\gamma_1 - \gamma_2}{2}$$

Now $|J| = abs \begin{vmatrix} \dfrac{\partial \left(\frac{\gamma_1 + \gamma_2}{2}\right)}{\partial \gamma_1} & \dfrac{\partial \left(\frac{\gamma_1 + \gamma_2}{2}\right)}{\partial \gamma_2} \\[4mm] \dfrac{\partial \left(\frac{\gamma_1 - \gamma_2}{2}\right)}{\partial \gamma_1} & \dfrac{\partial \left(\frac{\gamma_1 - \gamma_2}{2}\right)}{\partial \gamma_2} \end{vmatrix}$

$$= abs \begin{vmatrix} \frac{1}{2} & \frac{1}{2} \\[2mm] \frac{1}{2} & -\frac{1}{2} \end{vmatrix} = \frac{1}{2}$$

$$\Rightarrow f_{Z_1,Z_2}(\gamma_1,\gamma_2) = \frac{1}{2} f_{XY}\left(\frac{\gamma_1+\gamma_2}{2}, \frac{\gamma_1-\gamma_2}{2}\right)$$

Now lets compute marginal $Z_1$:

$$f_{Z_1}(\gamma_1) = \int_{-\infty}^{\infty} \frac{1}{2} f_{XY}\left(\frac{\gamma_1+\gamma_2}{2}, \frac{\gamma_1-\gamma_2}{2}\right) d\gamma_2$$

$$= \int_{-\infty}^{\infty} f_{XY}(t, \gamma_1-t) dt \qquad \left(\text{Put } t = \frac{\gamma_1+\gamma_2}{2}\right)$$

$\hookrightarrow$ This expression is familiar from prev. lecture. This shows we are consistent.

$\Longrightarrow$

$$Z_1 = X/Y$$
$$Z_2 = Y$$

again $\quad X = Z_1 Z_2$
$\qquad Y = Z_2$ $\qquad \left(\begin{array}{l} \text{i.e.} \\ h_1(\gamma_1,\gamma_2) = \gamma_1\gamma_2 \\ h_2(\gamma_1,\gamma_2) = \gamma_2 \end{array}\right)$

$$|J| = abs. \left| \begin{array}{cc} \gamma_2 & \gamma_1 \\ 0 & 1 \end{array} \right| = |\gamma_2|$$

$$\Rightarrow f_{Z_1,Z_2}(\gamma_1,\gamma_2) = |\gamma_2| f_{XY}(\gamma_1\gamma_2, \gamma_2)$$

6

Now again

$$f_{Z_1}(x_1) = \int_{-\infty}^{\infty} |x_2| \, f_{XY}(x_1 x_2, x_2) \, dx_2$$

This expression is also familiar from prev. lecture!

## EXPECTATIONS

Now lets return to the topic of expectations.

Consider $Z = g(x, y)$

we know $Z$ is a r.v.

So we know: $E[Z] = \begin{cases} \int_{-\infty}^{\infty} z \, f_Z(z) \, dz & \text{if } Z \text{ is continuous} \\ \sum_{\forall z_i} z_i \, f_Z(z_i) & \text{if } Z \text{ is discrete r.v.} \end{cases}$

But one can also show:

<u>Theorem</u>: $E[Z] = E[g(x,y)] = \begin{cases} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x,y) \, f_{XY}(x,y) \, dx \, dy \\ \sum_{\forall x_i} \sum_{\forall y_i} g(x_i, y_i) \, f_{XY}(x_i, y_i) \end{cases}$

Recall that we proved a similar theorem for $Z = g(x)$ also.

②

Again like prev. time we will show only for the discrete case:

Proof: $E[z] = \sum_{\forall \gamma_i} \gamma_i f_z(\gamma_i)$

$$= \sum_{\forall \gamma_i} \gamma_i \sum_{\substack{(x_i, y_i): \\ g(x_i, y_i) = \gamma_i}} f_{xy}(x_i, y_i)$$

$$= \sum_{\forall \gamma_i} \sum_{\substack{(x_i, y_i): g(x_i, y_i) = \gamma_i}} g(x_i, y_i) f_{xy}(x_i, y_i)$$

$$= \sum_{\forall (x_i, y_i)} g(x_i, y_i) f_{xy}(x_i, y_i)$$

if $E[z]$ exists then relies sum is abs. convergent so it doesn't matter in which order we take the sum!

In summary, we know how to compute expectation of function of two (of in general 'n') RVs!

In this lecture we will proceed with discussion of expectation in case of collections of r.v.s.

We already showed that:

$$E[g(x,y)] = \begin{cases} \int_{-\infty}^{\infty}\int_{-\infty}^{\infty} g(x,y)\, f_{xy}(x,y)\,dx\,dy & (\text{if } X, Y \text{ are jointly conts.}) \\[4mm] \sum_{\forall x_i}\sum_{\forall y_i} g(x_i, y_i)\, f_{xy}(x_i, y_i) & (\text{if } X, Y \text{ are discrete r.v.s}) \end{cases}$$

(All derivations from now on (unless specified explicitly) take case of X,Y jointly conts. and derive results on expectation using integrals, however the generic results do hold of for discrete r.v.s case also).

Consider $Z = g(x,y) = X + Y$

$$E[Z] = E\{X+Y\} = \int_{-\infty}^{\infty}\int_{-\infty}^{\infty}(x+y)\, f_{xy}(x,y)\,dx\,dy$$

$$= \int_{-\infty}^{\infty}\int_{-\infty}^{\infty} x\, f_{xy}(x,y)\,dy\,dx + \int_{-\infty}^{\infty}\int_{-\infty}^{\infty} y\, f_{xy}(x,y)\,dx\,dy$$

$$= \int_{-\infty}^{\infty} x\,\boxed{\int_{-\infty}^{\infty} f_{xy}(x,y)\,dy}\,dx + \int_{-\infty}^{\infty} y\,\boxed{\int_{-\infty}^{\infty} f_{xy}(x,y)\,dx}\,dy$$
$$\qquad\qquad\qquad\downarrow \qquad\qquad\qquad\qquad\qquad\downarrow$$

$$= \int_{-\infty}^{\infty} x\, f_x(x)\,dx + \int_{-\infty}^{\infty} y\, f_y(y)\,dy$$

$$= E\{X\} + E\{Y\}$$

In general, $E\{X_1 + \cdots + X_n\} = \sum_{i=1}^{n} E\{X_i\}$.

①

i.e. ~~from~~ Expectation of sum of rvs = sum of expectation of rvs.

(Note that we did NOT assume these rvs are independent)  (I)

~~Also Also~~

at $i^{th}$ trial among

∴ independent and identical

eg1  Let $X_i$ = indicator of success in $n$ ~~Berno~~ Bernoulli trials.

i.e. each $X_i$ is a (independent) Bernoulli rv with $P\{X_i=1\} = p$
$$P\{X_i=0\} = 1-p$$

Now $E\{X_i\} = 1 \cdot P\{X_i=1\} + 0 \cdot P\{X_i=0\} = p$

Consider the rv $X = X_1 + X_2 + \dots + X_n$. In words, $X$ is no. successes in $n$ iid Bernoulli trials. Of course $X$ follows a binomial distribution with parameters $(n, p)$.

Let compute $E\{X\}$ using (I):

$$E\{X\} = \sum_{i=1}^{n} E\{X_i\} = \sum_{i=1}^{n} p = np \quad \rightarrow \quad \text{we know this is } E\{X\} \text{ of binomial rv}$$

eg2  Let $X_i$ = indicator of changeover at $i^{th}$ interline between two consecutive tosses in $n$ independent coin tosses of the same coin.
(gap)

We have already seen that $P\{X_i=1\} = 2p(1-p)$

$$\Rightarrow E\{X_i\} = P\{X_i=1\} = 2p(1-p)$$

Now consider $X = X_1 + X_2 + \dots + X_{n-1}$. In words, $X$ is the number of changeovers in $n$ tosses!

$$\Rightarrow E\{X\} = \sum_{i=1}^{n} E\{X_i\} = 2(n-1)p(1-p).$$

(Note that here $X_i$ are not independent Bernoulli rvs, so $X$ is not binomial distributed. However the expectation matches to that of a binomial rv!)

②

Now nothing particular abt $g(x,y) = X+Y$, this linearity prop. of $E$ is ~~indeed~~ followed from the linearity prop. of integrals and summations. So in general one has:

⑫ Consider $g(x,y,z) = \sum\limits_{i=1}^{\ell} a_i \, f_i(x,y,z) + \sum\limits_{i=1}^{m} b_i \, g_i(x,z)$

$$+ \sum\limits_{i=1}^{n} c_i \, h_i(x) + d$$

linear combination of functions of $x,y,z$.

It is easy to see that :

$$E[g(x,y,z)] = \sum\limits_{i=1}^{\ell} a_i \, E[f_i(x,y,z)] + \sum\limits_{i=1}^{m} b_i \, E[g_i(x,z)]$$

$$+ \sum\limits_{i=1}^{n} c_i \, E[h_i(x)] + d,$$

to compute we will need joint dist of $x,y,z$

to compute we need joint.dist of $x,z$

we need only dist. of $x$.

~~This~~ (Also this can be further generalized to functions over $n$ rvs)

This is the linearity property of Expectation.

One can show another property of Expectation:

Suppose $X,Y$ are independent rvs. Then:

$$E[\underbrace{g(x)\,h(y)}_{f(x,y)}] = E[g(x)]\,E[h(y)] \qquad \text{Ⅱ}$$

Proof LHS $= \int\limits_{-\infty}^{\infty}\int\limits_{-\infty}^{\infty} g(x)h(y)\, f_{xy}(x,y)\,dx\,dy = \int\limits_{-\infty}^{\infty}\int\limits_{-\infty}^{\infty} g(x)h(y)\, f_x(x)\,f_y(y)\,dx\,dy$

$$= \left(\int\limits_{-\infty}^{\infty} g(x)\,f_x(x)\,dx\right)\left(\int\limits_{-\infty}^{\infty} h(y)\,f_y(y)\,dy\right) = \text{RHS}.$$

③

LHS involves double integral/summation whereas RHS involves two single integrals/summations. So it is useful observation.

Also, in general, we can show if $X, Y$ are independent then $g(X), h(Y)$ are also independent (Provided $g(X), h(Y)$ are well defined r.v.s!)

Proof: $g(X), h(Y)$ are independent r.v.s

$\Longleftrightarrow \Big[g(X) \in B_1\Big], \Big\{h(Y) \in B_2\Big]$ are independent events

$\forall B_1, B_2 \in \mathcal{B}$

$\Longleftrightarrow \Big[X \in g^{-1}(B_1)\Big], \Big[Y \in \mathbf{3} h^{-1}(B_2)\Big]$  ''

but $g^{-1}(B_1)$ and $h^{-1}(B_2)$ are none label sets!

which is true since $X, Y$ are themselves independent r.v.s!

## Moments of Collections of r.v.s

While of discussing r.v's we defined moments, absolute moments etc.

Now we of can extend these definitions:

$$M_{m,n} = E[x^m y^n] \longrightarrow m, n^{th} \text{ moment of } X, Y$$

(this is some function of $x, y$ hence we can) compute its expectation!

eg $M_{10} = E\{x\} = M_x$, $M_{01} = E\{Y\} = M_y$, $M_{11} = E\{XY\} \equiv M_{xy}$.

so on....

Similarly, one can extend the concept of central moments:

$$\sigma_{m,n} = E\left[(X-E\{X\})^m (Y-E\{Y\})^n\right] \to m,n^{th} \text{ central moment of } X, Y.$$

eg. $\sigma_{10} = 0 = \sigma_{01}$, $\sigma_{20} \doteq var(X) = \sigma_X^2$, $\sigma_{02} = var(Y) = \sigma_Y^2$,

$$\sigma_{11} = E\left[(X-E\{X\})(Y-E\{Y\})\right] \equiv Cov(X,Y)$$

$\searrow$ symbol
new

$\sigma_{11}$ is called as covariance of $X, Y$. (ofcourse $Cov(X,Y) = Cov(Y,X)$).

$Cov(X,Y)$ has some connection with the notion of "how correlated two r.v.s $X, Y$ are". Lets explore this connection now:

① Suppose $X, Y$ are independent. Then we can show $Cov(X,Y) = 0$

Proof: $Cov(X,Y) = E\left[(X-E\{X\})(Y-E\{Y\})\right]$

$$= E\left[(XY + E\{X\}E\{Y\} - XE\{Y\} - YE\{X\})\right] \quad \text{\textdollar}$$

$$= E\{XY\} - E\{X\}E\{X\} \qquad \to \text{linearity prop. of } E$$

$$= E\{X\}E\{Y\} - E\{X\}E\{Y\} \qquad \to \because X, Y \text{ are independent}$$
$$\qquad\qquad\qquad\qquad\qquad\qquad (by \text{ ②})$$

$$= 0$$

So $X, Y$ are independent $\implies Cov(X,Y) = 0$

However $Cov(X,Y) = 0 \not\implies X, Y$ are independent.

Here is the counter eg:
Consider $Y = X^2$ and $X$ is such that $E\{X\} = E\{X^3\} = 0$.
Note that $Y, X$ are surely dependent (not independent!)

⑤

However, for this eg: $\text{Cov}(X,Y) = \text{Cov}(X,X^2)$
$$= E\{X^3\} - E\{X\}E\{X^2\} = 0$$

So $\text{Cov}(X,Y) = 0 \not\Rightarrow X, Y$ are dependent.

(Note that an eg. of $X$ such that $E\{X\} = E\{X^3\} = 0$ is the
std. Normal r.v. In fact in assignment you showed that all
odd moments of a std. Normal r.v are zero.} Also you showed
that $Y = X^2$ has Chi-square distribution if $X$ is std. Normal)

$\longrightarrow$ For $X, Y$ Normal r.v ~~such that~~ $\text{Cov}(X,Y) = 0$ it ~~turns out that~~ turns out that
indeed $X, Y$ are independent! So Normal r.v's are an exception!
and the converse holds!!

We say $X, Y$ are uncorrelated if $\text{Cov}(X,Y) = 0$ .
(uncorrelated is here in some sense weaker cond. than independence)
In fact we can ~~q~~ quantify the "correlation" between
two r.v's using what is known as the correlation coefficient defined:
$$\text{as}$$

$$\rho_{xy} = \frac{\text{Cov}(X,Y)}{\sigma_X \, \sigma_Y} .$$

~~are to 0~~ of course $\rho_{xy} = 0 \Rightarrow \text{Cov}(X,Y) = 0 \Rightarrow X, Y$ are uncorrelated

Also one can show that:

$$|\rho_{xy}| \le 1 \quad \& \quad \rho_{xy} = \pm 1 \overset{\text{implies}}{\underset{\text{~~implies~~}}{\longrightarrow}} \text{ "perfect" correlation in the sense that } X, Y \text{ are linearly dependent!}$$

(6)

**Proof:** TST $|S_{xy}| \leq 1$

i.e. TST $(S_{xy})^2 \leq 1$

i.e. TST $(Cov(x,y))^2 \leq \sigma_x^2 \sigma_y^2$

i.e. TST $\left(E\{(X-E\{x\})(Y-E\{Y\})\}\right)^2 \leq E\{(X-E\{x\})^2\} E\{(Y-E\{Y\})^2\}$

(Lets put $X' = X - E\{x\}$, $Y' = Y - E\{Y\}$)

i.e. TST $\left(E\{x'y'\}\right)^2 \leq E\{x'^2\} E\{y'^2\}$ ⨷

⟶ This is known as Cauchy-Schwartz inequality

This inequality also appeare in linear algebra (vector spaces) and is a fundamental inequality. Infact this being satisfied in ⨷ form motivates the study of of vector spaces of r.v.s !!

⟶ Here is some intuition:

Suppose there exists some vector space in which inner product is given by $E\{x'y'\}$ i.e. $\langle V_1, V_2 \rangle = E[x'y']$

It is easy to see, $\langle V_1, V_1 \rangle = E[x'^2]$

$\langle V_2, V_2 \rangle = E\{y'^2\}$

but we know that $\left(\langle V_1, V_2 \rangle\right)^2 \leq \langle V_1, V_1 \rangle \langle V_2, V_2 \rangle$

$\left(\|V_1\| \|V_2\| \cos\theta\right)^2 \leq \|V_1\|^2 \|V_2\|^2$

⟶ $|\cos\theta| \leq 1$ which is true of course.

So ⨷ is extension of Cauchy-Schwartz inequality in case of Euclidean vectors!

⑦

Proof of ③ is simple (# typical wherever Cauchy-Swartz appears!)

Proof    Consider $E\left[(a\vec{x}'+y')^2\right]$. we know it is $\geq 0$.

$$\Rightarrow a^2 E\{\vec{x}'^2\} + 2a\, E\{x'y'\} + E\{y'^2\} \geq 0$$

$$\Leftrightarrow \left(E[x'y']\right)^2 \leq E\{x'^2\} E\{y'^2\}$$

(discriminant $\leq 0$)

Also note that ———— strict equality appears if and only if
$a x' + y' = 0$   ie.  X, Y are linearly dependent!
(with prob. 1)

This proves the overall claim that $|S_{xy}| \leq 1$ &

$$S_{xy} = \begin{cases} 0 & \text{is case while } X, Y \text{ are uncorrelated} \\ \pm 1 & \text{is case of "highest" correlation} \\ & \text{i.e. } X, Y \text{ are linearly related!} \end{cases}$$

Now go back to the example of swine flu. We want to know which of symptoms $X_1, X_2 \ldots, X_n$ is most important symptom that characterizes Y (presence of swine flu or not).

→(One answer)   Compute $|S_{x_1 y}|, |S_{x_2 y}|, \ldots, |S_{x_n y}|$

whichever symptom has max ↗ we can say it has highest "correlation" with disease and we can hence declare it to be the most important symptom for the disease!

Now lets compute $\text{var}(Z)$ where $Z = X+Y$, in terms of var & cov. of $X, Y$:

$$\text{var}(Z) = \text{var}(X+Y) = E\{(X+Y - E\{X+Y\})^2\}$$
$$= E[((X - E\{X\}) + (Y - E\{Y\}))^2]$$
$$= E\{(X - E\{X\})^2\} + E\{(Y - E\{Y\})^2\} + 2E\{(X - E\{X\})(Y - E\{Y\})\}$$
$$= \text{var}(X) + \text{var}(Y) + 2\,\text{Cov}(X,Y). \qquad \boxed{IV}$$

~~Note that var (x+y) = var(x) + var(y) + 2 cov(x,y) can be written~~
~~in the following (weird) ways:~~

~~var (x + y)~~

$\longrightarrow$

Apart ~~fr~~ from this sometimes "vectorial" versions of mean & variance
are defined. Here's the motivation:

$\longrightarrow$ X is an $n$-dimensional multivariate r.v.

Suppose we want to find $E[a^T X]$. $a^T X = \sum_{i=1}^{n} a_i X_i$.

$$E\{a^T X\} = E\{\sum_{i=1}^{n} a_i X_i\} = \sum_{i=1}^{n} a_i E\{X_i\} \rightarrow \text{by linearity prop. of } E.$$

$$= a^T E\{X\}$$

new notation $E\{X\} = \begin{bmatrix} E\{X_1\} \\ \vdots \\ E\{X_n\} \end{bmatrix}$

This $E\{X\}$ (exp. of multivariate r.v) is nothing but the vector of expectations
of the individual r.v.s. This is also sometimes called as the 'mean vector' of
the multivariate r.v $X$.

~~###~~ Now suppose we wish to find $\text{var}(a^T X)$:

$$\text{var}(a^T X) = E[(a^T X - E\{a^T X\})^2] = E[(a^T X - a^T E\{X\})^2]$$
$$= E\{a^T(X - E\{X\})\, a^T(X - E\{X\})\}$$

since transpose of a number is the number itself we get

$$\text{var}(a^T x) = E\left[a^T \underbrace{(x - E\{x\})(x - E\{x\})^T}_{n \times n \text{ matrix}} a\right]$$

by linearity prop. of $E$ one can show $= a^T \Sigma a$ where

$\Sigma$ is called as the covariance matrix whose entries

are given by $\Sigma_{ij} = \text{Cov}(X_i, X_j)$

$ij^{th}$ element of the covariance matrix.

⟶ for a 2-d case we can go through the steps easily:

$$\text{var}(a_1 X_1 + a_2 X_2) = E\left[(a_1 X_1 + a_2 X_2 - E\{a_1 X_1 + a_2 X_2\})^2\right]$$

$$= a_1^2 \text{var}(X_1) + a_2^2 \text{var}(X_2) + 2 a_1 a_2 \text{Cov}(X_1, X_2)$$

(by repeated appl. of linearity prop. of $E$, similar to (IV))

$$= \begin{bmatrix} a_1 & a_2 \end{bmatrix} \begin{bmatrix} \text{var}(X_1) & \text{Cov}(X_1, X_2) \\ \text{Cov}(X_1, X_2) & \text{var}(X_2) \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$$

$$= a^T \Sigma a$$

⟶ Hence sometimes instead of talking about moments and central moments of collections of rvs, people talk abt mean vector and covariance matrix of the corresponding multivariate rv.

(10)

# CONDITIONAL EXPECTATION

Suppose $X, Y$ are two random variables. Now in all cases

(i) $X, Y$ are discrete (ii) $X, Y$ are jointly conts (iii) one of them is conts. other is discrete,

we defined $f_{X/Y}(x/y) \longrightarrow$ either conditional pmf of conditional pdf

given $Y=y$ | if $X$ is discrete | $X$ is conts.

Now the r.v. for which $\uparrow$ is the pmf or pdf is denoted by:

$$Z = X/Y=y$$

We already now that $Z$ exactly takes those values which $X$ takes and its pmf/pdf is given by $f_{X/Y}(x/y)$.

Since $Z$ is a random variable we can talk about its expectation:

$$E[Z] = E[X/Y=y] = \begin{cases} \displaystyle\int_{-\infty}^{\infty} x \, f_{X/Y}(x/y) \, dx & \text{if } X \text{ is conts r.v.} \\[2mm] \displaystyle\sum_{\forall x_i} x_i \, f_{X/Y}(x_i/y) & \text{if } X \text{ is discrete r.v.} \end{cases}$$

This is called as conditional expectation of $X$ given that $Y=y$.

$\longrightarrow$ Now we can further extend this concept:

$X, Y$ are r.v.s & say $Z = g(x,y)$

we can talk abt. $f_{X/Z}(x/z)$ i.e. $f_{X/g(x,y)}(x/z)$

and in turn talk abt. $E[X/g(x,y)=z]$ and so on $\ldots$

(11)