

Lecture 1: Random Experiments and Probability

1. INTRO:

There are lot of things around us about which we cannot make any precise prediction. For example the weather, the outcome of a coin toss, the roll of a die, etc. Such processes are often referred to as "random experiments". So outcomes of random experiments are uncertain. However in most of these situations we know what are the possible outcomes. Each such outcome is called an elementary event. The set of all elementary events is called the sample space often denoted as Ω .

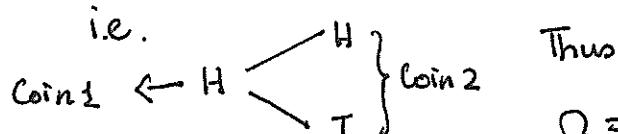
Example 1: If we toss a coin it is not sure in which way it will turn up when it hits the ground. One may say apart from head or tail one may argue that on a field for example a coin may get stuck and remain vertical. Thus in order to build a theory of probability we need make more idealistic assumption. So we talk of a fair coin, which can either be head or tail, when it falls to the ground. Thus

$$\Omega = \{H, T\}, : H \equiv \text{Head} \quad T \equiv \text{Tail}.$$

Example 2: If we roll a fair die, then the sample space is given as,

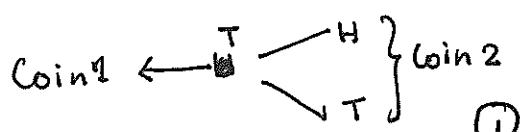
$$\Omega = \{1, 2, 3, 4, 5, 6\}$$

Example 3: Now let us toss two coins simultaneously. What is Ω ? Note that for each outcome of the first coin, there are two outcomes of the second coin, i.e.

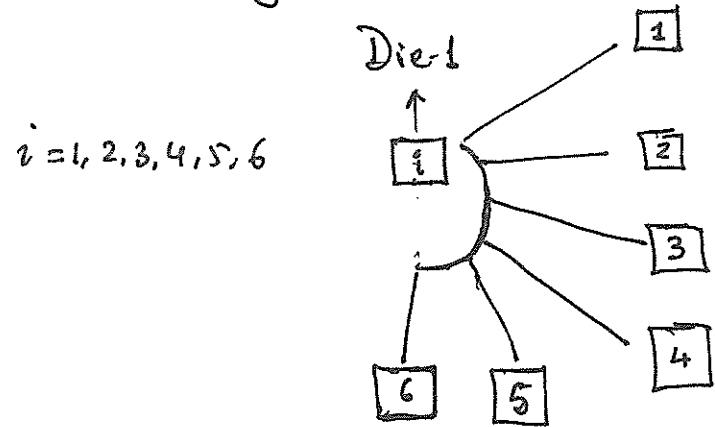


Thus

$$\Omega = \{HH, HT, TH, TT\}.$$



Example 4: If we role a fair die, then ^{roll} another fair die
 \Rightarrow then there are 36 possible outcomes, since for each of the six outcomes of the first die, there corresponds six outcomes of the second die. This is given as



$$\Omega = \left\{ (1, 1), (1, 2), (1, 3), (1, 4), (1, 5), (1, 6), (2, 1), (2, 2), (2, 3), (2, 4), (2, 5), (2, 6), (3, 1), (3, 2), (3, 3), (3, 4), (3, 5), (3, 6), (4, 1), (4, 2), (4, 3), (4, 4), (4, 5), (4, 6), (5, 1), (5, 2), (5, 3), (5, 4), (5, 5), (5, 6), (6, 1), (6, 2), (6, 3), (6, 4), (6, 5), (6, 6) \right\}$$

□

All our examples above may make us feel that for any random experiments the number of all possible outcomes is finite. But these need not be the case. For example consider the number of car accidents that may happen in Kanpur in the next 6 months. In fact theoretically any non-negative integer would do; i.e.

$$\Omega = \{0, 1, 2, 3, 4, \dots\}$$

As we will see later that there can be a sample space with uncountably infinite ~~numbers~~ sample points. We shall however in this chapter consider only sample spaces with finite number of sample points.

2. Events and Probability.

To be precise, for a finite sample space Ω , we can define any subset of Ω to be an event. Thus the power set of Ω , i.e. 2^Ω is called the set of all events associated with random experiment with sample space Ω . Remember the outcome of a trial of a random experiment is some $w \in \Omega$. If $w \in A \subset \Omega$, then we say that the event A has occurred.

For example let us consider the rolling of a fair die and let us say A be the event that an odd number appears. Then A is said to happen if the face of the die shows 1, 3 or 5. Thus $A = \{1, 3, 5\}$, which is a subset of $\Omega = \{1, 2, 3, 4, 5, 6\}$.

We say that two events A and B are mutually exclusive if $A \cap B = \emptyset$. We say events A_1, A_2, \dots, A_k are mutually exclusive if $A_i \cap A_j = \emptyset$, $i, j = 1, \dots, k$ & $i \neq j$. \emptyset or the empty-set denotes the null-event, so very often referred to as the impossible event.

Let us consider a random experiment which has N outcomes, i.e. $\Omega = \{\omega_1, \omega_2, \dots, \omega_N\}$, i.e. $\text{Card}(\Omega) = N$. Note $\{\omega_1\}, \{\omega_2\}, \dots, \{\omega_N\}$ forms a set of mutually exclusive events. For example if we toss a "fair" coin, then either "head" or "tail" appears. Both cannot appear at the same time. However we also "feel" that "head" or "tail" must have an equal chance of appearing materializing. This is what one says as "Head & Tail are equally likely events". In fact each $\omega_i \in \Omega$, is called an elementary event. So we are now ready to give the "Classical" definition of probability.

Classical events: Definition of Probability

Let a random experiment result in N , mutually exclusive, equally likely and exhaustive, elementary events. Let A be an event such that $\text{card}(A) = n_A$. Then probability of the occurrence of the event A , is the number $P(A)$ given as

$$P(A) = \frac{n_A}{N}$$

* Exhaustive means that apart from these N outcomes there are no other outcomes

Certain facts are immediately clear.

i) $0 \leq P(A) \leq 1$

ii) $P(\Omega) = \frac{N}{N} = 1,$

iii) $P(\emptyset) = \frac{0}{N} = 0$

iv) If $\text{card}(A) = n_A$ and $\text{card}(B) = n_B$, then

$$P(A \cup B) = \frac{\text{Card}(A \cup B)}{N} = \frac{n_A + n_B - n_{A \cap B}}{N}$$

$$\therefore P(A \cup B) = \frac{n_A}{N} + \frac{n_B}{N} - \frac{n_{A \cap B}}{N}$$

$$= P(A) + P(B) - P(A \cap B)$$

v) If $A \cap B = \emptyset \Rightarrow P(A \cup B) = P(A) + P(B).$

vi) A^c is ^{denotes} called the complement of the event A ; (ie. NOT A)
In fact $\text{card}(A^c) = N - n_A$. Thus

$$P(A^c) = \frac{N - n_A}{N} = 1 - \frac{n_A}{N}$$

$$\therefore P(A^c) = 1 - \frac{n_A}{N}$$

$$\therefore P(A^c) = 1 - P(A)$$

Another way to look at the last inequality is the following

$$A \cup A^c = \Omega,$$

Since $A \cap A^c = \emptyset, \Rightarrow P(\Omega) = P(A) + P(A^c)$
 $\Rightarrow P(A^c) = 1 - P(A).$

Drawback: One might note that the word "equally likely" must refer to the fact that each elementary event has an equal "probability" of occurrence, or that they are equiprobable. Thus the definition might appear circular. In the next chapter we shall see how Kolmogorov, the great Russian mathematician provided a solid mathematical foundation to probability.

The term equiprobable means: $P(\{\omega_i\}) = \frac{1}{N}$, for $i=1, \dots, N$.

Thus when you toss a coin: $P(\text{Head}) = 0.5, P(\text{Tail}) = 0.5$

Thus when you roll a die: $P(i) = \frac{1}{6}, i=1, 2, \dots, 6$.

In fact if you repeatedly toss a coin you will see after a large number of tosses, the number of heads and tails are almost equal. Thus using these kinds of experimentation one may have a feel that the classical definition is not bad at all.

We will soon see its power, but let us just begin with a story. Chevalier de Mere, pose who was a French gambler posed this problem to his mathematician friend Blaise Pascal

Which is more likelier?

Rolling at least one six in four throws of a die

OR

Getting a double six in 24 throws of a pair of die

The gambler reasoned as follows:

He felt that the average number of successful roll is same in both cases:

Case-I: Probability of getting six in a single throw is $\frac{1}{6}$.

So the expected value in 4 throws is $4 \times \frac{1}{6} = \frac{2}{3}$

Case-II: Probability of getting a double six when a double die is rolled is $\frac{1}{36}$. So the average in 24 throws is $24 \times \frac{1}{36} = \frac{2}{3}$.

But 'Chevalier de Mere' observed that he lost frequently when he used the second gamble. This problem excited Blaise Pascal, and he wrote about this to his friend Pierre de Fermat about this issue, and between them they formed the mathematical edifice of probability theory. We shall solve 'Chevalier de Mere's problem in Lecture 3.

Let us now use the classical definition to solve an interesting problem. Let there be N people in this a party. What is probability that at least two of them have the same birthday. Let us ignore leap years. This problem is solved by first computing the complement probability of the complementary event, i.e. no two people have their birthdays on the same day of the year. Let A^c denote that event

Now $\text{card}(\Omega) = 365^N$ in this case.

$$P(A^c) = \frac{365 \times 364 \times 363 \times \dots \times (365-(N-1))}{365^N}$$

The first person can have his birthday on any of the 365 days and second one on the remaining 364 days, the third on the remaining 363 days and so on. Thus our event A has

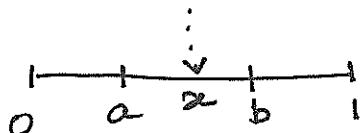
the probability

$$P(A) = 1 - \left(\frac{365 \times 364 \times 363 \times \dots \times (365-N+1)}{365^N} \right)$$

If $N = 25$, then one can show that $P(A^c) \approx 0.40$, showing that $P(A) \approx 0.60$, thus the probability of the birthdays of two people matching increases beyond 0.5 if $N \geq 25$. So if N is large such an event is equally likely.

3. Infinite Sample Spaces & Bertrand's Paradox

Consider the following problem. Let us be given the line $[0, 1]$



We throw a stone or a pin on the line. We would like to know if our pin falls in the interval (a, b) ,

$a < b$, & $0 < a < b < 1$, so our sample space and what is the probability that it will fall there.

It is intuitively clear that $\Omega = [0, 1]$. One way of calculating the probability of the event $A = \{w : a < w < b\}$, is to consider the lengths as a kind of "cardinality" of these intervals.

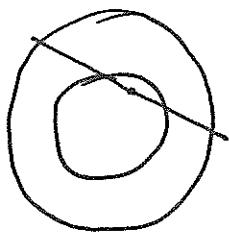
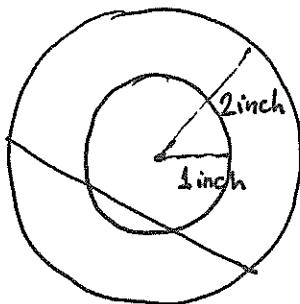
Thus we may write

$$P(A) = \frac{\text{length}(a, b)}{\text{length}[0, 1]} = b - a.$$

However observe that if this is the way we compute the probability in the above case then if $\tilde{A} = \{w : w = a\}$, then $P(\tilde{A}) = 0$. Thus it is impossible to throw a pin at an exact point on $[0, 1]$.

However such approaches have its own perils while we are handling the inf case of infinite sample spaces. This is aptly demonstrated by the Bertrand's paradox, where in the same event is shown to have different probabilities.

Bertrand's Paradox :

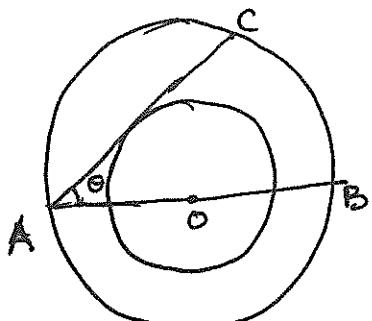


Consider a circular circle of radius 2 inch and a concentric circle of radius 1 inch. A chord is drawn at random in the circle of radius 2 inch. What is the probability of the chord cutting through the circle of radius one inch?

Solution: Approach 1: The chord of course is unique if it is not a diameter, in the sense that it can be identified through its mid-point. The chord to cut through the smaller concentric circle the mid point of the chord must lie in the circular disc formed by the smaller circle. Let us denote by A as the event that the chord cuts through the smaller circle. Here the sample space is the circular disc of radius 2inch. We shall get the probability of A as a ratio of areas

$$\text{P(A)} = \frac{\text{Area of Smaller circle}}{\text{Area of bigger circle}} = \frac{\pi}{\pi 4} = \frac{1}{4}$$

Solution: Approach 2



Fix a diameter and draw AB of the bigger circle and observe that any chord, starting at A will vary from 0 to π while while it varies from 0 to $\pm \frac{\pi}{6}$ to cut through the smaller circle. So, we have

$$\text{P(A)} = \frac{\frac{2\pi}{6}}{\pi} = \frac{1}{3}.$$

Note that for any chord AC we actually draw a horizontal diameter AB and argue as above.

You can see the contradiction!!

Lecture 2: Kolmogorov's Axioms for Probability

We have seen that though for finite sample spaces the idea of defining probability through counting works pretty well, we run into rough weather the moment we try to extend such an idea to an infinite sample space setting.

It was the great mathematicians, Emile Borel, Henri Lebesgue, Felix Hausdorff and Cantelli who realized that at a much deeper level the notion of probability is intimately linked with notion of measure of a set and thus notion of probability needed a fresh mathematicalization, which finally was completed by Kolmogorov. Borel and his co-workers gained a tremendous insight into the nature of probability by considering questions about infinite coin tosses and thus each instance can be represented as a point in the interval $[0, 1]$. This finally led to the strong law of large numbers which we will mention later in the course. Roughly what happens is the following. When we make an arbitrarily large number of trials the proportion of heads is nearly $\frac{1}{2}$ and remains there forever with probability one.

Borel and his co-workers understood the fact that every subset of an infinite sample space may not be an event. It was believed that events must have the following property

i) Ω , \emptyset must be events (sure event & null event must be there)

ii) If A is an event then so is A^c .

iii) If $A_1, A_2, \dots, A_n, \dots$ be a countable sequence of events then so is $\bigcup_{i=1}^{\infty} A_i$.

The class of such subsets of Ω is often denoted by \mathcal{F} and \mathcal{F} is often referred to as the σ -algebra or σ -fields of events.

One might be worried that we spoke of a countably infinite sequence of events but did not say anything about finite number events. But it is simple to show that if A_1, \dots, A_k are a finite number of events then $\bigcap_{i=1}^k A_i$ is also an event, i.e. a member of \mathcal{F} . Set $A_{n+1} = \emptyset$, $A_{n+2} = \emptyset, \dots$. Since by i) \emptyset is an event, we have

$$\bigcup_{i=1}^k A_i = A_1 \cup A_2 \cup \dots \cup A_k \cup \emptyset \cup \emptyset \cup \dots \cup \emptyset \dots$$

is also an event.

Further if A_1, \dots, A_n, \dots is a countable sequence of events then $\bigcap_{i=1}^{\infty} A_i$ is also an event (try this out as an assignment).

Kolmogorov assumed or rather viewed that given the pair (Ω, \mathcal{F}) , we can define a function $P: \mathcal{F} \rightarrow [0, 1] \subset \mathbb{R}$ (A set function actually) which satisfies the following rules

- i) For any $A \in \mathcal{F}$, $0 \leq P(A) \leq 1$ (This is clear from defn of P)
 - ii) $P(\Omega) = 1$
 - iii) If $\{A_n\}_{n=1}^{\infty}$ is a sequence of mutually exclusive events (i.e. $A_i \cap A_j = \emptyset$, if $i \neq j$), then
- $$P\left(\bigcup_{n=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i)$$

The third rule or axiom The above rules are often called Kolmogorov's axioms of probability. The third rule is a central one and provides us with some important facts about the function P which is called the probability function or probability measure.

Our first result which is a consequence of the axiom is as follows,

$$\boxed{P(\emptyset) = 0}$$

One might wonder why this was not mentioned in the axiom. In fact $P(\emptyset) = 0$ is a consequence of the third axiom, since if we set $A_i = \emptyset$, $\forall i \in \mathbb{N}$, (Here \mathbb{N} denotes the set of natural numbers), then $\{\emptyset\}_{n=1}^{\infty}$ is a mutually disjoint sequence of sets, and iii) gives us

$$P(\emptyset) = \sum_{n=1}^{\infty} P(\emptyset)$$

$$P(\emptyset) = P(\emptyset) + P(\emptyset) + \dots + P(\emptyset) + \dots$$

This equation is only possible if $P(\emptyset) = 0$ which is what we already know from the classical approach, but deriving this fact using Kolmogorov's axioms does not put any restriction on the sample space Ω .

Our next result is the following, which is again a consequence of the third rule and the fact that $P(\emptyset) = 0$. Let A_1, A_2, \dots, A_k be k mutually disjoint events, then

$$P\left(\bigcup_{n=1}^k A_n\right) = \sum_{n=1}^k P(A_n).$$

The strategy is the same; set $\bigcup_{n=k+1}^{\infty} A_n = \emptyset$, $A_{n+2} = \emptyset, \dots$

Thus

$$\bigcup_{n=1}^k A_n = A_1 \cup A_2 \cup A_3 \dots \cup A_k \cup \emptyset \cup \emptyset \dots \cup \emptyset$$

Thus using (iii) we have

$$\begin{aligned} P\left(\bigcup_{n=1}^k A_n\right) &= \sum_{n=1}^k P(A_n) + \sum_{n=k+1}^{\infty} P(\emptyset) \\ &= \sum_{n=1}^k P(A_n) \quad (\because P(\emptyset) = 0). \end{aligned}$$

(3)

Thus we have that if A_1, \dots, A_k are mutually disjoint events, then

$$\boxed{P\left(\bigcup_{n=1}^k A_n\right) = \sum_{n=1}^k P(A_n).} \longrightarrow \textcircled{\#}$$

Now as $A \cup A^c = \Omega$, using $\textcircled{\#}$ we have

$$\begin{aligned} P(\Omega) &= P(A \cup A^c) \\ &= P(A) + P(A^c) \text{ by } \textcircled{\#} \end{aligned}$$

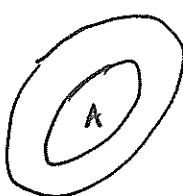
$$\therefore P(A) + P(A^c) = 1 \quad (\text{since } P(\Omega) = 1, \text{ from second axiom})$$

Hence irrespective of the nature of Ω we have

$$\boxed{P(A^c) = 1 - P(A)}$$

We will now derive some more results about the probability function based on the Kolmogorov's axioms.

- Let $A \subseteq B$, then $P(A) \leq P(B)$.



In this case we have

$$B = A \cup (B \setminus A) \quad [\text{show that } B \setminus A \in \mathcal{F}]$$

$$\therefore P(B) = P(A) + P(B \setminus A), \quad (\text{Using } \textcircled{\#})$$

$$\therefore P(B) - P(A) = P(B \setminus A)$$

But by i) $P(B \setminus A) \geq 0 \Rightarrow P(B) \geq P(A)$.

- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$, for any $A, B \in \mathcal{F}$.

$$\text{Now } A \cup B = (A \setminus B) \cup (B \setminus A) \cup (A \cap B)$$

We observe that above three sets on the left are mutually exclusive.

Hence using \oplus we have

$$P(A \cup B) = P(A \setminus B) + P(B \setminus A) + P(A \cap B) \longrightarrow \textcircled{*}$$

Again

$$A = E(A \setminus B) \cup (A \cap B).$$

But as $(A \setminus B)$ and $(A \cap B)$ are disjoint, i.e. mutually exclusive we have using \oplus

$$P(A) = P(A \setminus B) + P(A \cap B)$$

$$\therefore P(A \setminus B) = P(A) - P(A \cap B)$$

Further

$$B = E(B \setminus A) \cup (A \cap B)$$

and by similar arguments as \oplus before

$$P(B) = P(B \setminus A) + P(A \cap B)$$

$$\therefore P(B \setminus A) = P(B) - P(A \cap B).$$

Putting these expressions back in $\textcircled{*}$ we have

$$\begin{aligned} P(A \cup B) &= P(A) + P(B) - 2P(A \cap B) + P(A \cap B) \\ \therefore P(A \cup B) &= P(A) + P(B) - P(A \cap B) \end{aligned} \longrightarrow \textcircled{++}$$

We can immediately conclude that

$$P(A \cup B) \leq P(A) + P(B) \quad \text{since } P(A \cap B) \geq 0$$

We can extend this using mathematical induction to the following result. If A_1, \dots, A_k are elements of E

then

$$P\left(\bigcup_{n=1}^k A_n\right) \leq \sum_{n=1}^k P(A_n)$$

This is called Boole's Inequality.

The equality \leftrightarrow can be easily generalized for more than two events and is given as:-

$$\begin{aligned} P\left(\bigcup_{k=1}^n A_k\right) &= \sum_{k=1}^n P(A_k) - \sum_{k_1 < k_2} P(A_{k_1} \cap A_{k_2}) \\ &\quad + \sum_{\substack{k_1 < k_2 < k_3}} P(A_{k_1} \cap A_{k_2} \cap A_{k_3}) \\ &\quad - \sum_{\substack{k_1 < k_2 < k_3 < k_4}} P(A_{k_1} \cap A_{k_2} \cap A_{k_3} \cap A_{k_4}) + \dots \\ &\quad \dots + (-1)^{n+1} P\left(\bigcap_{k=1}^n A_k\right) \end{aligned}$$

This is sometimes referred to the as the "Principal of Inclusion and Exclusion".

Since the Kolmogorov's axiom involves, infinite sequences of events let us look into some properties of such sequences.

- Consider the set of non-decreasing events, i.e. $A_n \subseteq A_{n+1}$. The set $A = \bigcup_{n=1}^{\infty} A_n$ is then called the limit of $\{A_n\}$
- Similarly consider the sequence of non-increasing events, i.e. $A_{n+1} \subseteq A_n$. Then $A = \bigcap_{n=1}^{\infty} A_n$ is call the limit of $\{A_n\}$

Let us look at the following result, which is important enough to be mentioned as a theorem

Theorem 2.1: Consider a sequence $\{A_n\}$ of non-decreasing event.

Then $\lim_{n \rightarrow \infty} P(A_n) = P\left(\bigcup_{n=1}^{\infty} A_n\right)$.

Proof: Let us set $A = \bigcup_{n=1}^{\infty} A_n$. We can write A as follows

$$A = A_n + \bigcup_{j=n}^{\infty} (A_{j+1} \setminus A_j)$$

Again using the third axiom of Kolmogorov, (Think why??) we have

$$P(A) = P(A_n) + \sum_{j=n}^{\infty} P(A_{j+1} \setminus A_j)$$

Note that

$$0 \leq \sum_{j=1}^{\infty} P(A_{j+1} \setminus A_j) \leq 1$$

Thus $\sum_{j=1}^{\infty} P(A_{j+1} \setminus A_j)$ is a convergent series of non-negative numbers. That Thus

$$\lim_{n \rightarrow \infty} \sum_{j=n}^{\infty} P(A_{j+1} \setminus A_j) = 0 \quad \# (\text{see details})$$

Hence $n \rightarrow \infty$

$$P(A) = \lim_{n \rightarrow \infty} P(A_n) + \lim_{n \rightarrow \infty} \sum_{j=n}^{\infty} P(A_{j+1} \setminus A_j)$$

$$\boxed{\lim_{n \rightarrow \infty} P(A_n) = P(A)}$$

— x —

let $S = \sum_{i=1}^{\infty} a_i$, $S_n = \sum_{i=1}^n a_i$, Then for any $\epsilon > 0$, we have

~~#~~ $n_0 \in \mathbb{N}$, such that for all $n \geq n_0$

$$|S_n - S| < \epsilon$$

$$\Rightarrow \left| \cancel{a_{n+1} + a_{n+2} + \dots} \right| < \epsilon$$

$$\Rightarrow \left| \sum_{i=n+1}^{\infty} a_i \right| < \epsilon$$

$$\Rightarrow \lim_{n \rightarrow \infty} \sum_{i=n+1}^{\infty} a_i = 0.$$

(7)

Lecture 3: Conditional Probability, Independence & Bayes Theorem

Sec 1: Conditional Probability

In this lecture we shall introduce the reader to a very important notion of probability theory called "Conditional Probability". Mathematically stated it means the following: Given a probability space (Ω, \mathcal{F}, P) , and events $A, B \in \mathcal{F}$, we ask the question as to what is the probability of the event A given the fact or conditioned on the fact that the event B has already occurred?

We ask such questions pretty often asked in real-life scenarios. For example given that a train has departed from the originating station with a ten minute delay, and we may be interested to know the probability that it will reach its' destination on time.

The probability of an event A , given that the event B has occurred is called the conditional probability of A given that B , and is denoted by $P(A|B)$.

Definition 3.1 : Conditional Probability: Given (Ω, \mathcal{F}, P) to be a probability space, and $P(B) > 0$; where $B \in \mathcal{F}$. Then the conditional probability $P(A|B)$, for any $A \in \mathcal{F}$ is given as

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Let us view this definition from the classical angle. Let Ω be a sample space, with N possible outcomes. The event B has occurred means ~~no~~ number that some event B has occurred. Now what is the probability that $A \in B$ has occurred. Now what is the probability that $A \in B$ will occur. So now our total number of possibilities have shrunk to only those which represent the event B . So among those events possibilities we need to what is proportion that is common to ~~the~~ the event A , i.e. the sample points which represent the occurrence $A \cap B$.

$$\therefore P(A|B) = \frac{n_{A \cap B}}{n_B} = \frac{\frac{n_{A \cap B}}{N}}{\frac{n_B}{N}} = \frac{P(A \cap B)}{P(B)}, \text{ where } P(B) > 0.$$

Example 3.1 Let us consider an urn having 5 red & 5 black balls. Find the probability that a red ball appears on the 2nd draw given that the first draw resulted in a black ball. The first ball is not replaced back in the urn.

Let B be the event that the first ball drawn was black. Let A be the event that second ball drawn is red.

$$\therefore P(A|B) = \frac{5}{9}$$

How did we come to this value. We have already drawn a black ball. So before making the second draw we have only 9 balls left so that is exactly our total possible outcomes and we have 5 balls are red, which gives the above value
Observe that from the definition of $P(B|A) \rightarrow P(A|B)$

we have

$$P(A \cap B) = P(A|B) P(B)$$

So in Example 3.1 we have

$$P(A \cap B) = P(A|B) P(B)$$

$$\text{Here } P(B) = \frac{5}{10} = \frac{1}{2} \quad \therefore P(A \cap B) = \frac{5}{9} \times \frac{1}{2} = \frac{5}{18}$$

Of course $A \cap B$ is the event that a black ball appears in the first draw and a red in the second

Given any $B \in \mathcal{F}$ with $P(B) > 0$, we define the conditional probability function given B as $P(\cdot|B) : \mathcal{F} \rightarrow [0, 1]^R$, i.e.

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \forall A \in \mathcal{F}.$$

Now let us see if the axioms of Kolmogorov are satisfied.
For any $A \in \mathcal{F}$, it is obvious that $P(A|B) \geq 0$. Further as $A \cap B \subset B$, $\Rightarrow P(A \cap B) \leq P(B)$. Thus $P(A|B) \leq 1$.

Now if $A = \Omega$ we have

$$P(\Omega|B) = \frac{P(\Omega \cap B)}{P(B)} = \frac{P(B)}{P(B)} = 1.$$

Let $\{A_n\}$ be a sequence of mutually disjoint events, i.e.

then

$$\begin{aligned} P\left(\bigcup_{n=1}^{\infty} A_n | B\right) &= \frac{P\left(\left(\bigcup_{n=1}^{\infty} A_n\right) \cap B\right)}{P(B)} \\ &= \frac{P\left(\bigcup_{n=1}^{\infty}(A_n \cap B)\right)}{P(B)} \end{aligned}$$

As $\{\cdot \cap B\}$, forms a disjoint sequence of events. Hence

$$\begin{aligned} P\left(\bigcup_{n=1}^{\infty}(A_n \cap B)\right) &= \sum_{n=1}^{\infty} P(A_n \cap B) \\ \therefore P\left(\bigcup_{n=1}^{\infty} A_n | B\right) &= \sum_{n=1}^{\infty} \frac{P(A_n \cap B)}{P(B)} \\ &= \sum_{n=1}^{\infty} P(A_n | B). \end{aligned}$$

Thus all the three laws of Kolmogorov holds. Thus we have shown that for any given event B , with $P(B) > 0$, $P(\cdot | B)$ is a probability measure or probability function.

We shall now state an important result called the "Theorem of Total Probabilities" which will lead to do computations with the famous Bayes Theorem, which we will discuss in the next section.

Theorem of Total Probability

Theorem 3.1.1 : Let (Ω, \mathcal{F}, P) be a probability space. Let $\{B_1, \dots, B_n\}$ be a collection of n mutually disjoint collection of events such that $\Omega = \bigcup_{i=1}^n B_i$. Further assume that for each $i \in \{1, 2, \dots, n\}$, $P(B_i) > 0$.

Then for any $A \in \mathcal{F}$

$$P(A) = \sum_{i=1}^n P(B_i) P(A | B_i)$$

Proof: Note that since $\{B_1, \dots, B_n\}$ forms a mutually disjoint partition of Ω , then A either occurs jointly with B_1 or B_2 or \dots or B_n , i.e.

$$A = \bigcup_{i=1}^n (A \cap B_i)$$

while we know that $(A \cap B_1), \dots, (A \cap B_n)$ are mutually disjoint. Then

$$P(A) = P\left(\bigcup_{i=1}^n (A \cap B_i)\right)$$

$$= \sum_{i=1}^n P(A \cap B_i) \quad \left[\because A \cap B_i, i=1, \dots, n \text{ are mutually exclusive} \right]$$

By using exclusion-inclusion principle.

$$\therefore P(A) = \sum_{i=1}^n P(B_i) P(A|B_i)$$

This theorem remains true if $n = \infty$.
Think how.

Section 2: Bayes' Theorem

Thomas Bayes was a priest, an English priest. ~~The~~ The theorem that now bears his name, was published after his notes were edited by Robert Price. Thomas Bayes had not published his result. The idea of Bayes later led to the field of Bayesian statistics which has huge applications in modern medicine and machine learning, for example.

Bayes Theorem Let us consider the probability space (Ω, \mathcal{F}, P) .

Let $\{B_1, \dots, B_n\}$ be a finite sequence of events such that $P(B_i) > 0, \forall i = 1, \dots, n$ & $\Omega = \bigcup_{i=1}^n B_i$. Let $A \in \mathcal{F}$ such that $P(A) > 0$, then

$$P(B_i|A) = \frac{P(B_i \cap A)}{P(A)}$$

$$P(B_i|A) = \frac{P(B_i) P(A|B_i)}{\sum_{i=1}^n P(B_i) P(A|B_i)}$$

(4)

This theorem also holds for $n = \infty$. Now we shall provide the simple proof.

$$P(B_i | A) = \frac{P(B_i \cap A)}{P(A)}$$

$$P(B_i \cap A) = P(A) P(B_i) P(A|B_i)$$

while by the theorem of total probability

$$P(A) = \sum_{i=1}^n P(B_i) P(A|B_i).$$

A more simple version of the Bayes Theorem can be provided with two events $A \& B \in F$, with $P(A) > 0, P(B) > 0$.

Note that if $B \in F \Rightarrow B^c \in F$ and we know that

$$\Omega = B \cup B^c$$

$$\begin{aligned} P(A) &= P(A \cap B) + P(A \cap B^c) \\ &= P(B) P(A|B) + P(B^c) P(A|B^c) \end{aligned}$$

$$\therefore P(B|A) = \frac{P(B) P(A|B)}{P(B) P(A|B) + P(B^c) P(A|B^c)}$$

What the Bayes theorem does is as follows. It can help us answering the following questions. For example in an urn there are n_w white balls and n_b black balls. Suppose balls are drawn at random without replacement. Suppose we draw the balls twice and suppose that the second ball is white: Then we can ask the question as to what is the probability that the first ball drawn was white or the first ball drawn is black. Bayes theorem helps us in answering such questions.

If B be the event that the first ball drawn is black then B^c is naturally the event that the first ball drawn is white. Let W be the event that the second ball drawn is white. So

$$P(B|W) = \frac{P(B) P(W|B)}{P(W)} \quad (\text{By Bayes Theorem}).$$

$$\text{Now } P(B) = \frac{n_B}{n_B + n_{B^c}}$$

$$P(W|B) = \frac{n_W}{n_B + n_W - 1} \quad \frac{n_W}{n_B + n_W - 1}$$

while

$$\begin{aligned}
 P(W) &= P(W \cap B) + P(W \cap B^c) \\
 &= P(B) \cdot P(W|B) + P(B^c) \cdot P(W|B^c) \\
 &= \frac{n_B}{n_B + n_W} \cdot \frac{n_W}{n_B + n_W - 1} + \frac{n_W}{n_B + n_W} \cdot \frac{n_W - 1}{n_B + n_W - 1} \\
 &= \frac{n_B n_W + n_W(n_W - 1)}{(n_B + n_W)(n_B + n_W - 1)} \\
 &= \frac{n_B n_W + n_W^2 - n_W}{(n_B + n_W)(n_B + n_W - 1)} \\
 \therefore P(B|W) &= \frac{\frac{n_B}{n_B + n_W} \cdot \frac{n_W}{n_B + n_W - 1}}{\frac{n_B n_W + n_W^2 - n_W}{(n_B + n_W)(n_B + n_W - 1)}} \\
 &= \frac{n_B n_W}{n_B n_W + n_W^2 - n_W}
 \end{aligned}$$

Since $n_W \in \mathbb{N}$, i.e. $n_W \geq 1 \Rightarrow n_W^2 \geq n_W \Rightarrow n_W^2 - n_W \geq 0$

$$\Rightarrow n_B n_W + n_W^2 - n_W \geq n_B n_W$$

$\Rightarrow P(B|W) \leq 1$, and thus fulfilling the basic requirement of the probability function.

$$P(B|W) = \frac{n_B n_W}{n_B n_W + n_W^2 - n_W}$$

We urge the reader to compute $P(B^c|A) = ??$

Independence of Events

I do not even remember which year but I remember reading a book called "The Enigmas of Chance", by the famous mathematician Mark Kac. where he discussed an important idea called "Independence of Events" i.e. two events A and B, are called independent if and only if

$$P(A|B) = P(A), \text{ if } P(B) > 0$$

$$P(B|A) = P(B), \text{ if } P(A) > 0$$

i.e.

$$P(A \cap B) = P(A) \cdot P(B)$$

→ (the key fact)

This idea as I read in the book was due to Hugo Steinhaus, who developed this notion, when he was in hiding during the second world war. Steinhaus was a famous Polish mathematician, who made major contributions to an area of mathematics.

Let us see how this can be applied. Consider the simple situation where you toss a coin twice. What is the probability that head (H) appears in the second toss for the first time

Let A be the event that head appears in the second toss for the first time; and B be the event that tail appears. Thus we are trying to find the probability the configuration HT. We know that $P(HT) = \frac{1}{4}$ from direct computation. But

$$HT = B \cap A$$

$$\therefore P(HT) = P(B \cap A) = \frac{1}{4}$$

$$\therefore P(B \cap A) = \frac{1}{2} \cdot \frac{1}{2} = P(B) \cdot P(A)$$

Thus B & A are independent events and which appears to be intuitively true.

Let us now introduce the notion of independence of three events $A, B, C \in \mathcal{F}$. The definition of independence now means the following: Three events A, B & C are independent events if and only if

$$P(A \cap B) = P(A) \cdot P(B)$$

$$P(B \cap C) = P(B) \cdot P(C)$$

$$P(C \cap A) = P(C) \cdot P(A)$$

$$P(A \cap B \cap C) = P(A) \cdot P(B) \cdot P(C).$$

It is important that pairwise independence of events do not mean they are independent. The following example from the book [2] by Moods, Graybill and Boes, ~~shows~~ establishes what we have just mentioned above.

Example: Pairwise independence does not mean independence.

Let

A_1 : Event that odd face appears in the first die

A_2 : Event that odd face appears in the second die.

A_3 : The sum of the two numbers in a random throw of two die is odd.

$$P(A_1) = \frac{3}{6} = \frac{1}{2} \quad \& \quad P(A_2) = \frac{3}{6} = \frac{1}{2}$$

$$P(A_1 \cap A_2) = \frac{9}{36} = \frac{1}{4} = \frac{1}{2} \cdot \frac{1}{2} = P(A_1) P(A_2).$$

We urge the reader to show that

$$P(A_1 \cap A_3) = P(A_1) P(A_3)$$

$$P(A_2 \cap A_3) = P(A_2) P(A_3)$$

$$P(A_3) = \frac{18}{36} = \frac{1}{2} \left(A_3 = \{(1, 2), (2, 1), (1, 4), (4, 1), (1, 6), (6, 1), (2, 3), (3, 2), (5, 6), (6, 5), (3, 6), (6, 3), (4, 3), (3, 4), (5, 4), (4, 5), (2, 5), (5, 2)\} \right)$$

$$P(A_1 \cap A_3) = \frac{9}{36} = \frac{1}{4} \quad (\text{Finding the cases in } A_3 \text{ with odd first term})$$

$$= \frac{1}{2} \cdot \frac{1}{2} = P(A_1) P(A_3)$$

We leave $P(A_2 \cap A_3) = P(A_2) P(A_3)$ verification to the reader.

It is symmetric.

$$P(A_1 \cap A_2 \cap A_3) = 0.$$

$$\text{But } P(A_1) \cap P(A_2) \cap P(A_3) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{8}$$

$$\therefore P(A_1 \cap A_2 \cap A_3) \neq P(A_1) P(A_2) P(A_3)$$

At the end of this section let us solve the "Chevalier de Mere" problem and also revisit the birthday problem.

For the first case in the "Chevalier de Mere" problem one is looking at the appearance of atleast one six in the throw of four dies dice or four throws of a die.

Let A_i be the event that six does not appear on the i -th throw

$$\text{then } P(A_i) = \frac{5}{6}, \quad i=1, \dots, 4$$

Thus if A is the event the six does not appear in any of the four throws, then

$$A = \bigcap_{i=1}^4 A_i$$

$$\therefore P(A) = P\left(\bigcap_{i=1}^4 A_i\right) = P(A_1)P(A_2)P(A_3)P(A_4)$$

o [as A_1, A_2, A_3, A_4 are independent events]

$$= \left(\frac{5}{6}\right)^4$$

~~If B is the probability~~ Hence A^c is the event at least one six appears we have

$$P(A^c) = 1 - \left(\frac{5}{6}\right)^4$$

His second case was a double six appearing on 24 throws of a pair of die.

Let B_i be the event that a double six does not occur in the i -th throw, then $P(B_i) = \left(\frac{35}{36}\right)$, $i=1, \dots, 24$.

Let B be the event that a double six never occurs, then

$$B = \bigcap_{i=1}^{24} B_i$$

Hence the independence of B_i 's showed that

$$P(B) = \left(\frac{35}{36}\right)^{24}$$

Our required event is B^c , and hence

$$P(B^c) = 1 - \left(\frac{35}{36}\right)^{24}$$

which is bigger $P(B^c)$ or $P(A^c)$? They don't appear to be equal in the face of it.

From the internet for example from mathworld.wolfram.com. we gather that

$$P(A^c) = 1 - \left(\frac{5}{6}\right)^4 \approx 0.5177$$

$$P(B^c) = 1 - \left(\frac{35}{36}\right)^4 \approx 0.4914$$

and thus $P(A^c) > P(B^c)$. and that is why Chevalier de Mere was ^{mostly} losing when he chose the second bet.

Let us revisit the "Birthday problem" in terms of conditional probability. Instead of asking the question, as what is the probability of at least two people having the same birthday we ask the question, ^{that} in a group of people what is the probability that none of the members have the same birthday. We provide the following approach from [S1].

Suppose we choose two people from the group. What is the prop probability that they have different birthdays?

Call that event B_2 , then

$$P(B_2) = 1 - \frac{1}{365}. \quad [\text{not considered leap year.}]$$

What happens ~~we~~ if we have three people. What is $P(B_3)$? In fact we can write

$$B_3 = (A_3 \cap B_2)$$

if A_3 is the event that the third person has a birthdate which does not match with any one of the other two.

$$P(B_3) = P(A_3 \cap B_2) = P(A_3 | B_2) P(B_2)$$

$$\text{Now } P(A_3 | B_2) = \frac{363}{365} = 1 - \frac{2}{365}$$

$$\therefore P(B_3) = \left(1 - \frac{2}{365}\right) \left(1 - \frac{1}{365}\right) \approx 0.9918$$

Thus in general we can write when we are looking at a group of n -people we can write

$$B_n = A_n \cap B_{n-1}$$

$$\therefore P(B_n) = P(A_n | B_{n-1}) P(B_{n-1})$$

$$= \left(1 - \frac{(n-1)}{365}\right) P(B_{n-1})$$

$$= \left(1 - \frac{n-1}{365}\right) \left(1 - \frac{n-2}{365}\right) P(B_{n-2})$$

$$\vdots \quad \vdots$$

$$= \left(1 - \frac{n-1}{365}\right) \dots \left(1 - \frac{2}{365}\right) P(B_2)$$

$$= \left(1 - \frac{n-1}{365}\right) \dots \left(1 - \frac{2}{365}\right) \left(1 - \frac{1}{365}\right)$$

It was shown for example in [S1], that $P(B_{22}) \approx 0.5243$, while $P(B_{23}) \approx 0.4927$. So if we have 23 or more people in the group their probability of having different birthdays goes down below 0.5 and has less chance of occurrence.

[S1]: F.M. Dekking, C. Kraaikamp, H.P. Lopuhaä, & L.E. Meester, A Modern Introduction to Probability and Statistics: Understanding Why and How. Springer 2010

Lecture 4: Random Variables, Their Distribution and their Expectation

Sec 1: RANDOM VARIABLES

We are all accustomed to handling real numbers. Thus it might be a good idea to shift the arena for computing probability from the zone of abstract events to the real lines. This is done through the vehicle of random variables. But before we embark into the definition of a random variable ~~lets~~ ^{let} us get a fact straight: Random Variable is a function.

A random variable $X : \Omega \rightarrow \mathbb{R}$, a function from the sample space Ω , to the real line \mathbb{R} , such that for any real number $a \in \mathbb{R}$, such that the set

$$S_a = \{ \omega \in \Omega : X(\omega) \leq a \} \subset \mathcal{F}$$

i.e. the set S_a is an event.

Consider the random experiment of tossing a coin, i.e. $\Omega = \{H, T\}$
Define random variable

$$X(H) = 1, \quad X(T) = 0.$$

Suppose we consider a σ -algebra made out of open intervals in \mathbb{R} , i.e. do taking their unions, complements, etc. We call such a σ -algebra as the Borel σ -algebra. In fact for a random variable the set

$$S_I = \{ \omega : X(\omega) \in I \},$$

where I is an open interval is also an event. Define for each I

$$P_X(I) = P \{ \omega : X(\omega) \in I \},$$

We leave it to the reader should be able to show that $P_X(I)$ indeed satisfies the Kolmogorov Axioms. Thus if \mathcal{B} is the notation for Borel σ -algebra, then $(\mathbb{R}, \mathcal{B}, P_X)$ is a probability space induced by the random variable X .

Suppose the range of X is finite, i.e. $\text{Range } X = \{x_1, \dots, x_n\}$

Then consider the sets

$$\Omega_i = \{\omega : X(\omega) = x_i\}, i=1, \dots, n$$

Then $\Omega_i \cap \Omega_j = \emptyset, \forall i \neq j$ & $\bigcup_{i=1}^n \Omega_i = \Omega$, i.e. a random variable X with a finite range induce a partition of the sample space Ω . We urge the reader to extend this idea to a random variable (r.v. for short) with a countable range.

Any random variable X whose range is either finite or countable is called discrete random variable. When studying discrete random variable we often ask the question, what is the probability that X takes the value say x_k or just say x .

We denote this by $P(X=x)$, which means

$$P(X=x) = P\{\omega \in \Omega : X(\omega) = x\}$$

Now suppose the range of the r.v. is countable, i.e

$$\text{range } X = \{x_1, \dots, x_n, \dots\} \subseteq \mathbb{R}$$

Any ~~x while~~ $x \in \mathbb{R}$, which is not in the range has zero probability. Thus it is meaningful to ask the question

$$\text{What is } P(X=x_k) = P(\{\omega \in \Omega : X(\omega) = x_k\})$$

$$= P(\Omega_k)$$

$$\text{As } \Omega = \bigcup_{k=1}^{\infty} \Omega_k \Rightarrow P(\Omega) = \sum_{k=1}^{\infty} P(\Omega_k)$$

$$\therefore \boxed{\sum_{k=1}^{\infty} P(X=x_k) = 1}$$

i.e. the total probability is one.

Now for any $x \in \mathbb{R}$, even for a discrete random variable X , it is meaningful to ask the question, what is the probability that X takes values less than x . Thus

$$P(X \leq x) = \sum_{x_k \leq x} P(X = x_k)$$

So more formally we write

$$f_X(x) = P(X = x) \quad (\text{p.m.f for short})$$

which is called the probability mass function of a discrete random variable, and

$$F_X(x) = P(X \leq x) = \sum_{x_k \leq x} f_X(x_k)$$

is called the cumulative distribution function (c.d.f), or distribution function of the random variable.

Consider the following simple example, of throwing a fair die. Define a random variable $X: \Omega \rightarrow \mathbb{R}$ as follows;

$$X(\omega) = \begin{cases} 1, & \text{if } \omega \text{ is odd} \\ 0, & \text{if } \omega \text{ is even} \end{cases}$$

The range of X is given as $\text{range } X = \{0, 1\}$. The probability mass function thus given as

$$f_X(x) = \begin{cases} \frac{1}{2} & P(X=1) = \frac{1}{2} \\ & P(X=0) = \frac{1}{2} \end{cases}$$

Note: $P(X=1) = P(\{\omega : X(\omega)=1\}) = P(\{1, 2, 3\}) = \frac{1}{2}$, same approach for calculating $P(X=0)$.

The distribution function is given as

$$F_X(x) = \begin{cases} 0 & \text{if } x < 0 \\ \frac{1}{2} & \text{if } 0 \leq x < 1 \\ 1 & \text{if } x \geq 1 \end{cases}$$

Observe the way we do the segmenting of the variable, x . The custom

is to have strict inequality on the right-hand side which allows us to graphically represent as a step function, the cdf of a discrete r.v. X . In our example of the die, the graph of $F(x)$ is given F_X is given as

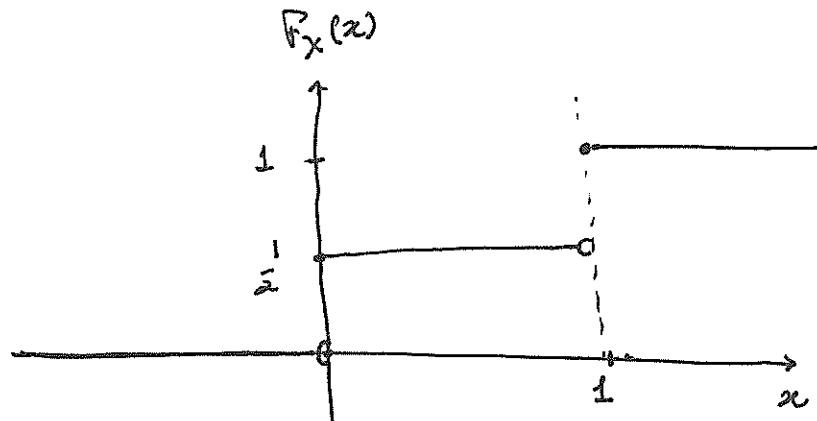


Fig 4.1 c.d.f of X in the previous example.

The c.d.f. or distribution function has the following properties.

- i) F_X is a non-negative function, and naturally $F_X(x) \leq 1$
- ii) If $x < y$, we have $F_X(x) \leq F_X(y)$.

This is what we get just by observing the figures above figure of the cdf. It is not a difficult task to prove the above two from the very definition of a c.d.f.

Another important result shows us how to compute a p.m.f. value of the c.d.f. of a discrete r.v. is given.

Let us arrange the values of a r.v. X with finite range in the ascending order

$$x_1 < x_2 < x_3 < \dots < x_n$$

and $f(x_i) = F(x_i)$, and

$$f(x_i) = F_X(x_i) - F_X(x_{i-1}), \quad i=1, 2, \dots, n$$

This can be proved from the definition of F_X as

$$F_X(x_i) = F_X(x_{i-1}) + f(x_i),$$

(A)

Further observe that F_X is right continuous, i.e.

$$\lim_{x \rightarrow a^+} F(x) = F(a)$$

It is however not left continuous. Note that we are writing all these properties by looking at the graph of the cdf in Fig 4.1. Also observe that

$$\lim_{x \rightarrow -\infty} F_X(x) = 0$$

while

$$\lim_{x \rightarrow +\infty} F_X(x) = 1$$

You might wonder why the above two properties uses limits. It is done so to consider the countably infinite random variables which has a countably infinite range. Let us now summarize the properties of the distribution function F_X of a discrete random variable X , that we have learnt.

- i) $F_X \geq 0 \text{ & } F_X \leq 1$
- ii) $\lim_{x \rightarrow a^+} F_X(x) = F(a)$
- iii) If $x < y \Rightarrow F_X(x) \leq F_X(y)$
- iv) $\lim_{x \rightarrow -\infty} F_X(x) = 0 \text{ & } \lim_{x \rightarrow +\infty} F_X(x) = 1$

Every random variable need not have a finite or countable range. A random variable may have a range which is uncountably infinite. A r.v. X is called continuous if 'its' range is uncountably infinite.

Suppose we want to measure height of all students here at IIT Kanpur. Then we cannot say that there are only a finite set of point height measures possible. In fact we can reasonably say that height of students are in the interval $[4, 7]$, i.e. between 4 ft and 7 ft. So the height of a student is a random variable.

$$\Omega = \text{Students} \xrightarrow{\text{Height} = X} \mathbb{R}$$

So if X is the random variable denoting the height of a student, then

$$4 \leq X(\omega) \leq 7$$

and thus ~~X~~ X has an uncountably infinite range. Now which question is more meaningful,

$$P(X(\omega) = 5 \text{ ft}) \quad \text{or} \quad P(4.8 \leq X(\omega) \leq 5.2).$$

So lets first consider the fact $P(4.8 \leq X(\omega) \leq 5.2)$, since whenever we make a measurement, random errors and observational errors creep in and it is impossible to say that an interval is of exactly this particular length. In fact if we measure the same line segment five times, each time we are likely to get very slightly different results. Keeping our intuition of the classical approach we may write

$$P(4.8 \leq x \leq 5.2) = \frac{\text{length of } [4.8, 5.2]}{\text{length of } [4, 7]} = \frac{1}{3}.$$

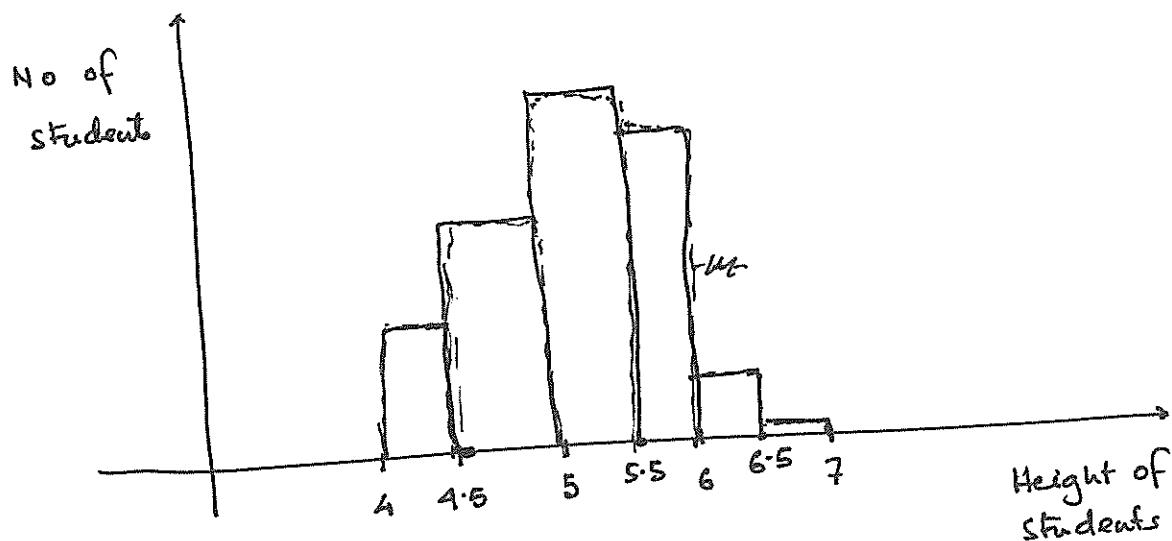
$$\text{However } P(x = 5) = \frac{\text{length of } \{5\}}{\text{length of } [4, 7]} = 0$$

(6) (6)

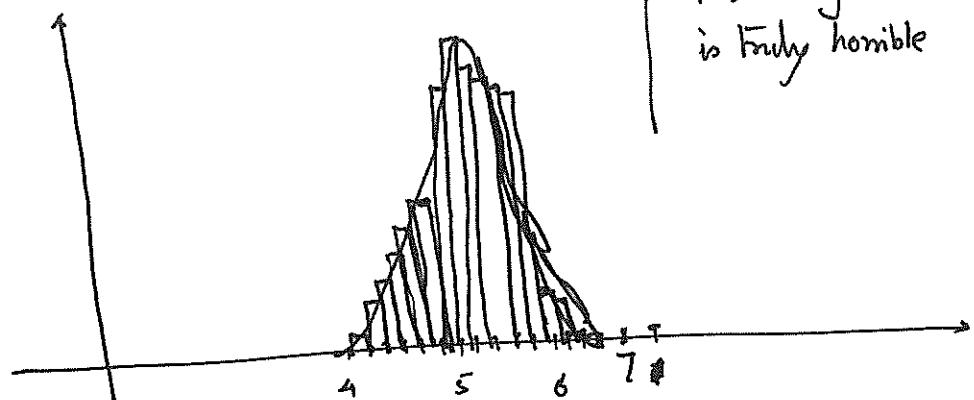
* We are aware of the perils of such an approach. Recall Bertrand's paradox.

But is that the only way we deal with the case ~~case~~ when range of X is uncountably infinite. There could be many situations where our intuition from classical probability may fail. A better path is to pass through pictorial representation of data.

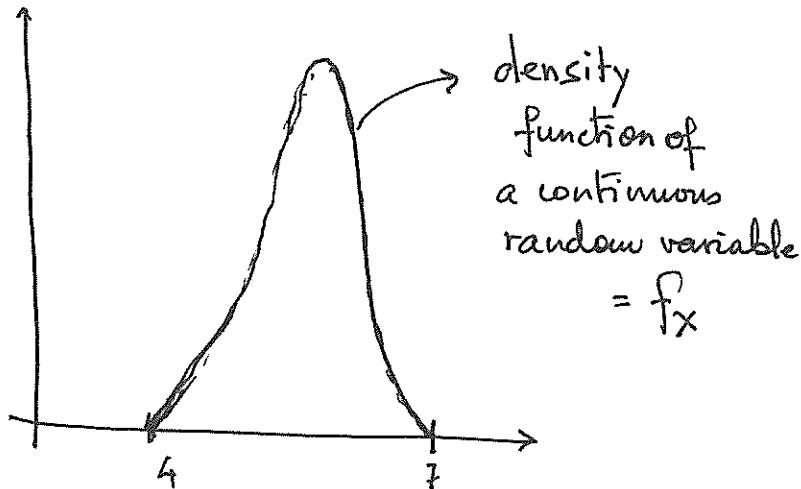
Consider again the heights of students at IIT Kanpur. Let us divide the height range $[4, 7]$ into intervals of length 0.5 along the horizontal axis, and in the vertical axis the number of ~~per~~ students with height in ~~that~~ those intervals. The resulting diagram is called a histogram which we depict below.



If we make more finer partitions of the interval, then we have a diagram as



So as we make the partitioning intervals smaller and smaller we can represent the histogram by a smooth curve which we will call the density function, since each rectangle in some sense represents the density of the people with heights in that interval.



Note that the values of the density function of a continuous random variable denoted as f_x does not represent the probabilities, as at individual points the probability is zero while as we see the density function is not.

From a formal perspective, the probability density function of a continuous random variable X , is a function f_x which has the following properties

$$\text{i)} \quad f_x(x) \geq 0, \quad \forall x$$

$$\text{ii)} \quad \int_{-\infty}^{\infty} f_x(x) dx = 1$$

$$\text{iii)} \quad P(a \leq x \leq b) = \int_a^b f_x(x) dx.$$

Another approach to view a continuous random variable X is to assume that it has a distribution function F_X which is continuous. Indeed we have

$$F_X(x) = P(X \leq x) = \int_{-\infty}^x f_x(z) dz.$$

Using (iii). Remember that all other properties of the distribution function remains same, as we have seen for the discrete variable.

Thus for a continuous random variable, X , we have

$$\begin{aligned} P(a \leq X \leq b) &= \int_a^b f_X(x) dx \\ &= \int_{-\infty}^b f_X(x) dx - \int_{-\infty}^a f_X(x) dx \\ &= F_X(b) - F_X(a). \end{aligned}$$

Thus

$$P(a \leq X \leq b) = F_X(b) - F_X(a)$$

One can also show that

$$\frac{dF_X}{dx} = f_X.$$

One might have a sense that somehow, all these we have discussed for the continuous case is largely intuitive. However there is a result which we will not prove, but will every thing clear.

We have already discussed that given a probability space (Ω, \mathcal{F}, P) and X be a random variable associated with it. Then using X we can shift our working space to $(\mathbb{R}, \mathcal{B})$, where \mathcal{B} is the Borel σ -algebra in \mathbb{R} . The following Theorem shows how a unique probability measure \hat{P} can be defined for the space $(\mathbb{R}, \mathcal{B})$. We present the result as given in Robert Ash*

Theorem 4.1 : Let f_X be a non-negative function on \mathbb{R} , with

$$\int_{-\infty}^{\infty} f_X(x) dx = 1. \text{ Then there exists a unique probability measure } \hat{P}$$

such that for any Borel subset of $B \in \mathcal{B}$ we have

$$\hat{P}(B) = \int_B f_X(x) dx$$

is a probability measure on \mathbb{R} .

The reader can show that $\hat{P}(B)$ satisfies the Kolmogorov Axioms.

* Robert B. Ash : Basic Probability Theory, Dover 2008,

As an example consider, the following c.d.f or distribution function.

$$F_X(x) = \begin{cases} 1 - e^{-x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

It is a continuous function. It is non-decreasing also and
 $\lim_{x \rightarrow -\infty} F_X(x) = 0 \Rightarrow \lim_{x \rightarrow +\infty} F_X(x) = 1$. Hence F_X represents a continuous distribution function. Hence X is a continuous random variable.

The probability density function associated with X is

$$f_X(x) = \begin{cases} e^{-x}, & \text{when } x \geq 0 \\ 0, & \text{when } x < 0 \end{cases}$$

(We just differentiated F_X , i.e we set $f_X(x) = \frac{dF_X}{dx}$).

Let us now check that f_X is truly a density function. Of course $f_X \geq 0, \forall x \in \mathbb{R}$ &

$$\begin{aligned} \int_{-\infty}^{\infty} f_X(x) dx &= \int_0^{\infty} e^{-x} dx \\ &= \left[-e^{-x} \right]_0^{\infty} = 1. \end{aligned}$$

There are ~~several~~ other random variables, which are of the mixed type, i.e. it has both discrete part and continuous.

For our purposes, we shall stick to discrete and continuous variables and in the next chapter we shall look at some special (means largely applicable) distribution of discrete random variables, followed by a chapter on distributions for continuous random variables. Among the continuous variable distributions the normal distribution will play a crucial role in statistics.

[*Mathematicians by the way term the normal distribution as the Gaussian distribution, after the Gauß, the originator of the idea. More on that later.]

Section 2: Expectation of a Random Variable

If we take a simplistic approach, then expectation means average. We are also aware of the idea of weighted average.

Let us consider n observations (could be height or weight...blah blah)

$$x_1, x_2, \dots, x_n.$$

To each of these observation assign weights w_1, w_2, \dots, w_n respectively & $w_i \geq 0$, for all $i = 1, \dots, n$. Then weighted mean means is defined as

$$\bar{x} = \frac{w_1 x_1 + w_2 x_2 + \dots + w_n x_n}{w_1 + w_2 + \dots + w_n}$$

If we set $\bar{w}_i = \frac{w_i}{\sum_{i=1}^n w_i}$, where $\sum_i w_i = w_1 + \dots + w_n$,

then

$$\bar{x} = \bar{w}_1 x_1 + \bar{w}_2 x_2 + \dots + \bar{w}_n x_n$$

Now for each i , $0 \leq \bar{w}_i \leq 1$ & $\sum_{i=1}^n \bar{w}_i = 1$, and thus the weights \bar{w}_i can be viewed as probability associated with occurrence of the ~~the i -th observation~~ of the i -th observation turning out to be x_i . Thus taking a cue from the weighted mean the expectation of a discrete random variable x is given as

$$E(x) = \sum_{i=1}^{\infty} x_i P(x=x_i) = \sum_{i=1}^{\infty} x_i f_X(x_i) \quad \begin{cases} \text{provided} \\ \text{the series} \\ \text{is converges} \end{cases}$$

The definition of expectation for a continuous random variable is quite analogous and is given as

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

provided the integral is finite.

Given a random variable $X: \Omega \rightarrow \mathbb{R}$, & $g: \text{Range } X \rightarrow \mathbb{R}$, the function $Y = g \circ X$, is a function from $\Omega \rightarrow \mathbb{R}$. Is this also a random variable? (The value of Y is computed as $Y(\omega) = g(X(\omega))$)

Now consider the set $\{\omega: Y(\omega) \in I\}$, for any interval I , we need to show that $\{\omega: Y(\omega) \in I\} \in \mathcal{F}$, if $(\Omega, \mathcal{F}, \mathbb{P})$ is the probability space we are working on.

Observe that $\{\omega: Y(\omega) \in I\} = Y^{-1}(I)$, i.e. to prove that Y is a random variable we have to show that $Y^{-1}(I) \in \mathcal{F}$ for any interval I in \mathbb{R} . Observe that

$$\begin{aligned} Y^{-1}(I) &= (g \circ X)^{-1}(I) \\ &= X^{-1} \circ g^{-1}(I) \end{aligned}$$

Now $g^{-1}(I) \subset \text{Range } X$. ~~Thus for any~~ ^{Any} subset of $\text{Range } X$ is a Borel set and hence $X^{-1}(g^{-1}(I)) \in \mathcal{F}$ as X is a random variable.

So for $Y = g(X)$, we have, when X is discrete

$$E(Y) = E[g(X)] = \sum_{i=1}^{\infty} g(x_i) f_X(x_i) \quad (\text{provided the series converges})$$

and when X is continuous

$$E(Y) = E[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx, \quad (\text{provided the integral is finite}).$$

E.g. If $g(x) = X + c$, where X is constant then

$$\begin{aligned} E[X+c] &= \int_{-\infty}^{\infty} (x+c) f_X(x) dx = \int_{-\infty}^{\infty} x f_X(x) dx + c \int_{-\infty}^{\infty} f_X(x) dx \\ &= E(X) + c \\ &\because \int_{-\infty}^{\infty} f_X(x) dx. \end{aligned}$$

You can check out some properties in the assignments. Consider again a the random variable X with the distribution

$$f_X(x) = \begin{cases} e^{-x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

$$\begin{aligned} \text{Then } E(X) &= \int_{-\infty}^{\infty} x f_X(x) dx \\ &= \int_0^{\infty} x e^{-x} dx \\ &= -xe^{-x} \Big|_0^{\infty} - \int_0^{\infty} \frac{dx}{dx} (-e^{-x}) dx \\ &= 0 + \int_0^{\infty} e^{-x} dx = 1 \end{aligned}$$

$$\text{i.e. } E(X) = 1.$$

However for every distribution one need not have a finite expectation. Consider the discrete distribution, where the random variable has countable values, i.e. X has values $\{0, 1, 2, 3, 4, \dots\}$.

Let us have

$$P_X(x) = \frac{1}{2^x}$$

Consider $g(x) = 2^x$. then

$$E(g(x)) = \sum_{x=0}^{\infty} 2^x \frac{1}{2^x} = \sum_{x=0}^{\infty} 1 + 1 + \dots + 1 + \dots$$

Thus The series $1 + 1 + \dots + 1 + \dots$ is not convergent and hence $E(g(x))$ does not exist. This is often called the St. Petersburg Paradox.

In statistics, $E(x)$ is traditionally given the symbol μ .
 Another important measure in statistics is that of variance.
 Variance is a measure of dispersion. It measures the average square deviation of the random variable values from the mean. Thus for a random variable X

$$\sigma_x^2 = \text{Variance of } X = \text{Var}(X) = E(x - \mu)^2$$

$$\begin{aligned} \text{Thus } E(x - \mu)^2 &= E(x^2 - 2x\mu + \mu^2) \\ &= E(x^2) - 2\mu E(x) + \mu^2 \\ &= E(x^2) - 2\mu^2 + \mu^2 \quad (\text{check this}) \\ &= E(x^2) - 2E(x)^2 + E(x)^2 \\ &= E(x^2) - E(x)^2 \end{aligned}$$

$$\therefore \boxed{\text{Variance of } X = \text{Var } X = E(x^2) - E(x)^2}$$

Let us finish this chapter by computing the variance of the random variable X , whose pdf is

$$f_X(x) = \begin{cases} e^{-x}, & x \geq 0 \\ 0, & x < 0. \end{cases}$$

$$\begin{aligned} \text{Thus } \text{Var}(X) &= E(x^2) - E(x)^2 \\ &= E(x^2) - 1 \\ &= \int_0^\infty x^2 e^{-x} dx - 1 \\ &= \left[-x^2 e^{-x} \right]_0^\infty - 2 \int_0^\infty x e^{-x} dx - 1 \\ &= [0 + 2E(x)] - 1 \\ &= 2 - 1 = 1. \end{aligned}$$

Section 3: Chebychev's Theorem

We begin with a general form of Chebychev's Theorem:

Theorem 4.2 Let g be a non-negative function of a random variable X , defined on a probability space (Ω, \mathcal{F}, P) . Then for any $k > 0$ we have

$$\text{Pf} \quad \boxed{P(g(X) \geq k) \leq \frac{E[g(X)]}{k}}$$

Proof: Observe that if X is a continuous random variable,

$$\begin{aligned} E(g(X)) &= \int_{-\infty}^{\infty} g(x) f_X(x) dx \\ &= \int_{\{x: g(x) \geq k\}} g(x) f_X(x) dx + \int_{\{x: g(x) < k\}} g(x) f_X(x) dx \\ &\geq \int_{\{x: g(x) \geq k\}} g(x) f_X(x) dx \quad \left[\text{Since } \int_{\{x: g(x) < k\}} g(x) f_X(x) dx \geq 0 \right] \\ &\geq k \int_{\{x: g(x) \geq k\}} f_X(x) dx = k P[g(\overset{X}{\cancel{x}}) \geq k] \end{aligned}$$

$$\therefore P[g(\overset{X}{\cancel{x}}) \geq k] \leq \frac{E[g(X)]}{k} \quad \square$$

Further if $g(x) = (x - \mu)^2$ we have $P[(x - \mu)^2 \geq k] \leq \frac{\text{Var}(x)}{k}$

Chebychev's inequality says that

$$P(|x-\mu| \geq r\sigma_x) \leq \frac{1}{r^2}, \quad \text{where } \sigma_x^2 = \text{Var}(X) \neq 0$$

~~In~~ Observe that

$$\begin{aligned} P(|x-\mu| \geq r\sigma_x) &= P((x-\mu)^2 \geq r^2\sigma_x^2) \\ &\leq \frac{\text{Var}(x)}{r^2\sigma_x^2} \quad (\text{By Theorem 4.2}) \\ \Rightarrow P(|x-\mu| \geq r\sigma_x) &\leq \frac{\cancel{x-\mu}^2}{r^2\sigma_x^2} \\ \Rightarrow P(|x-\mu| \geq r\sigma_x) &\leq \boxed{\frac{1}{r^2}} \end{aligned}$$

This is Chebychev's inequality. What it says is the following

$$\begin{aligned} -P(|x-\mu| \geq r\sigma_x) &\geq -\frac{1}{r^2} \\ \Rightarrow 1 - P(|x-\mu| \geq r\sigma_x) &\geq 1 - \frac{1}{r^2} \\ \text{c.e.} \quad \boxed{P(|x-\mu| < r\sigma_x) \geq 1 - \frac{1}{r^2}} &\rightarrow \textcircled{O} \end{aligned}$$

In fact this says that

$$\boxed{P(\mu - r\sigma_x < x < \mu + r\sigma_x) \geq 1 - \frac{1}{r^2}}$$

This says that the probability, that the value of the random variable X falls in the open interval $(\mu - r\sigma_x, \mu + r\sigma_x)$ increases as r increases. We will make use of this idea later in our study.

Section 4: Moments and Moment Generating Function

The concept of moments play an important role in the practical applications. The ~~n~~-r-th moment is simply the expectation of the random variable function $g(x) = x^r$ of the random variable X , i.e.

$$\mu'_r = E(X^r)$$

If $r=1$, then $\mu'_1 = E(X) = \mu_X$, the mean or expectation.

There is also a notion of central moment, i.e. moment centered at a value μ_X , i.e.

$$\mu_r = E[(x - \mu_X)^r]$$

In this case $\mu_1 = 0$, while $\mu_2 = \sigma_X^2$.

But how do we compute moments. This is done through the device of moment generating function, which we now discuss.

The moment generating function, or mgf for short is given as ^{of a r.v. X}

$$m_X(t) = E[e^{tX}]$$

For the continuous case we have

$$m_X(t) = E[e^{tX}] = \int_{-\infty}^{\infty} e^{tx} f_X(x) dx$$

of course for $m_X(t)$ to be finite this integral must exist for ~~every~~ the given t . In fact it is much better that $m_X(t)$ is finite at least over an open interval, ~~for~~ i.e. $\exists h > 0$, s.t. $\forall t \in (-h, h)$, the function $m_X(t)$ is finite.

For the discrete case of course we have

$$m_X(t) = E(e^{tx}) = \sum_x e^{tx} f_X(x).$$

Now if g_1 & g_2 are functions of the random variable X , then one can easily show that

$$E[g_1(x) + g_2(x)] = E[g_1(x)] + E[g_2(x)]$$

Can this be generalized to the countable case. Let

$$Y = \sum_{i=1}^{\infty} g_i(x)$$

and assume that $Y(w)$ is finite for all w and $E(Y)$ exists, and if $E[g_i(x)]$ exists for each i , then

$$E[Y] = \sum_{i=1}^{\infty} E(g_i(x)) \quad [\text{of course under some conditions}]$$

provided $\sum_{i=1}^{\infty} E[g_i(x)]$ is also convergent.

$$\begin{aligned} E[e^{tx}] &= E\left[1 + \frac{tx}{1!} + \frac{t^2 x^2}{2!} + \frac{t^3 x^3}{3!} + \dots + \frac{t^n x^n}{n!} + \dots\right] \\ &= 1 + \mu'_1 t + \mu'_2 \frac{t^2}{2!} + \dots + \mu'_n \frac{t^n}{n!} + \dots \end{aligned}$$

In fact the coefficient of t^n is μ'_n , for $n=1, 2, \dots$

If the expectation $E[e^{tx}]$ is finite then one can show that the above series also converges. Suppose we can integrate under the differential sign. Suppose we can differentiate under the integral sign. Suppose we differentiate the mom mgf, r times, then we get

$$\frac{d^r}{dt^r} m(t) = \int_{-\infty}^{\infty} x^r e^{tx} f_X(x) dx$$

$$\Rightarrow \boxed{\frac{d^r}{dt^r} m(0) = \int_{-\infty}^{\infty} x^r f_X(x) dx = \mu'_r} \quad (18)$$

Lecture -5

Chapters: Binomial & Poission Distributions

In this chapter we seek to study some specialized probability distributions, which has a lot of applicability. In this chapter in particular we shall focus on Binomial and Poission distribution which are important distributions for discrete random variables.

For a discrete random variable X , its probability distribution is simply an algebraic expression of its probability mass function. Of course special distributions are for particular type of random variables, describing special situations.

Section 1: Binomial Distributions

Suppose we make repeated coin tossing. Of course a fair coin is used. Suppose we toss a coin 100 times and ask the question that what is the probability that head appears 25 times? The Binomial Distribution answers such questions.

Binomial Distribution deals with repeated trials and each trial has ~~less~~ only two possible outcomes, which mark as success or failure. Each such trial is often called a Bernoulli trial. For example when we toss a coin for example we can consider the appearance of head (H) as success and appearance of tail as failure. If we denote success as S and failure as F, then any repeated random trial with ~~n~~ trials looks like may look like the following cases

SSFF...FSSS....S

So there could be k successes out of n trials. We ask ourselves what is the probability of such an event if each success has a fixed probability p . Note when n trials are carried out k successes can appear in $\binom{n}{k}$ ways.

Thus if X denotes the random variable of denoting the number of success, then

$$P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$$

In fact S denote.

So the Binomial p.m.f is usually written as

$$f_X(x) = \binom{n}{x} p^x q^{n-x}, \text{ where } q=1-p$$

Let us first check that f_X is indeed a probability mass function.

i) $f_X(x) \geq 0$, is of course clear as $p \geq 0$

$$\begin{aligned} \text{ii)} \quad \sum_{x=0}^n f_X(x) &= \sum_{x=1}^n \binom{n}{x} p^x q^{n-x} \\ &= (p+q)^n = 1. \quad (\text{Last step used Binomial Theorem}) \end{aligned}$$

So why it is called a Binomial distribution is clear.

Before providing any example of a numerical problem let us compute two important measures associated with a random variable, namely expectation and variance.

Thus if X is a binomial random variable, then we will compute $E(X)$ & $\text{Var}(X)$. There are two approaches to it. We shall choose the approach using the moment generating function. Now the mgf of the binomial

random variable is

$$\begin{aligned} m_X(t) &= \sum_{x=0}^n e^{tx} \binom{n}{x} p^x q^{n-x} \\ &= (pe^t + q)^n \end{aligned}$$

Thus

$$m_x(t) = (pe^t + q)^n$$

$$\therefore m'_x(t) = n(p e^t + q)^{n-1} p e^t$$

$$\therefore m'_x(0) = n(p+q)^{n-1} p$$

$$m'_x(0) = np \quad (\text{as } p+q=1)$$

$$\text{Hence } \mu_x = E(x) = m'_x(0) = np.$$

Now to compute the variance we compute $m''_x(0)$. So

$$m''_x(t) = n(n-1)(p e^t + q)^{n-2} p e^t p e^t \\ + n(p e^t + q)^{n-1} p e^t$$

$$\therefore m''_x(0) = n(n-1)(p+q)^{n-2} \cancel{(p e^0)} p^2 \\ + n(p+q)^{n-1} p.$$

$$\therefore m''_x(0) = n(n-1)p^2 + np \\ = (n^2 - n)p^2 + np \\ = n^2p^2 - np^2 + np \\ = n^2p^2 + np - np^2 \\ = n^2p^2 + np(1-p) \\ = n^2p^2 + npq$$

$$\therefore E(x^2) = n^2p^2 + npq$$

$$\therefore \text{Var}(x) = E(x^2) - [E(x)]^2 = n^2p^2 + npq - n^2p^2$$

$$\therefore \boxed{\text{Var}(x) = npq}$$

The whole idea of using m.g.f might appear counter-intuitive as if something is forced on to us, without motivation and does not seem to have immediate motivation connect with the definition. Of course we can definitely compute $E[X]$ and $\text{Var}(X)$ of a binomial random variable, using brute force, which we demonstrate below.

Brute force computation:

$$\begin{aligned}
 \mu_x = E(x) &= \sum_{x=0}^n x \binom{n}{x} p^x q^{n-x} \\
 &= \sum_{x=0}^n x \frac{n!}{x!(n-x)!} p^x q^{n-x} \\
 &= \sum_{x=0}^n \frac{n!}{(x-1)!(n-x)!} p^x q^{n-x} \\
 \therefore \mu_x &= np \sum_{x=0}^n \frac{(n-1)!}{(x-1)!(n-x)!} p^{x-1} q^{n-x} \\
 &= np \sum_{x=0}^n \binom{n-1}{x-1} p^{x-1} q^{n-x} \\
 &= np \left(p + q \right)^{n-1} \\
 &= np.
 \end{aligned}$$

However computing $E(X^2)$ using brute force technique is not so simple. There is a small trick. This is as follows.

$$\begin{aligned}
 E(X^2) &= E(X^2 - x + x) \\
 &= E(X^2 - x) + E(x) \\
 &= E(X(X-1)) + E(x).
 \end{aligned}$$

$$\begin{aligned}
 E(X(x-1)) &= \sum_{x=0}^n x(x-1) \binom{n}{x} p^x q^{n-x} \\
 &= \sum_{x=0}^n x(x-1) \frac{n!}{x!(n-x)!} p^x q^{n-x} \\
 &= \sum_{x=0}^n x(x-1) \frac{n!}{x(x-1)(x-2)!(n-x)!} p^x q^{n-x} \\
 &= \sum_{x=0}^n \frac{n!}{(x-2)!(n-x)!} p^x q^{n-x} \\
 &= n(n-1) p^2 \sum_{x=0}^n \frac{(n-2)!}{(x-2)!(n-2)!} p^{x-2} q^{n-x} \\
 &= n(n-1) p^2 (p+q)^{n-2} \\
 &= n(n-1) p^2
 \end{aligned}$$

$$E(x^2) = n(n-1)p^2 + np^0 -$$

Hence $\text{Var}(X)$ is given as

$$\begin{aligned}
 \text{Var}(X) &= E(x^2) - [E(x)]^2 \\
 &= n(n-1)p^2 + np^0 - n^2 p^2 \\
 &= n^2 p^2 - np^2 + np - n^2 p^2 \\
 &= np(1-p) = npq.
 \end{aligned}$$

The square root of the variance of X is called the standard deviation of X . $\sigma_X = \sqrt{\text{Var}(X)}$. Let us now look into a numerical example, where we can see the application of binomial distribution.

Example 5.1: A jumbo-jet, say Boeing-747-400 has four engines.

they operate independently. Let each engine has a probability p of failure and for a successful flight at least two engines should work. What is the probability that the flight will be successful.

Solution: The probability of engine not failing is p

Let X be the number of engines working during the flight
We have to find the probability that $X \geq 2$, i.e

$$P(X \geq 2) = 1 - P(X < 2)$$

Of course X follows a binomial distribution with $n=4$ &
probability of success p . Thus

$$\begin{aligned} P(X \geq 2) &= 1 - (P(X=0) + P(X=1)) \\ &= 1 - \left(\binom{4}{0} p^0 q^{4-0} + \binom{4}{1} p^1 q^{4-1} \right) \\ &= 1 - [q^4 + 4pq^3] \end{aligned}$$

$$\boxed{P(X \geq 2) = 1 - q^4 - 4pq^3}$$

□

Observe that if X is a binomial random variable, then
its distribution depends on two parameters, the number of
trials n and the probability of success p . So often
one writes

$$\boxed{X \sim B(n, p)}$$

which means that X is a binomial random variable
with parameters n & p .

Section 2: Poisson Distribution

Suppose you are editing an essay with 5,000 words and you want to know if there what is the probability, that there are five spelling mistakes, given that the probability of committing the mistake is 0.005. If we consider that find occurrence of a spelling mistake is a success, then modelling through the Binomial distribution shows that

$$f_x(s) = \binom{5000}{5} (0.005)^5 (0.995)^{4995}$$

Of course this is a cumbersome computation. The question is can we make it simpler. We are in a situation where n is large and p small while the product np is moderate. Our question is such cases, what form would the Binomial pmf take. Trying to figure out this leads us to the Poisson distribution.

Let us assume the $np = \lambda > 0$, then setting $p = \frac{\lambda}{n}$, we can rewrite the Binomial distribution as

$$\begin{aligned} f_x(x) &= \binom{n}{x} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x} \\ &= \frac{n!}{x!(n-x)!} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x} \left(1 - \frac{\lambda}{n}\right)^{-x} \\ &= \frac{n(n-1)(n-2)\dots(n-x+1)(n-x)!}{x!(n-x)!} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x} \end{aligned}$$

Thus

$$f_x(x) = \frac{1}{x!} \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \dots \left(1 - \frac{(x-1)}{n}\right) \lambda^x \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x}$$

Now we need to see what happens to the right hand side when $n \rightarrow \infty$. We of course know the celebrated formula

$$\lim_{n \rightarrow \infty} \left(1 + \frac{z}{n}\right)^n = e^z$$

Thus

$$\lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^n = e^{-\lambda}$$

Further $\left(1 - \frac{\lambda}{n}\right)^{-x} \rightarrow 1$ as $n \rightarrow \infty$. Hence the rhs becomes (as $n \rightarrow \infty$)

$$\boxed{\frac{\lambda^x e^{-\lambda}}{x!}}$$

This gives us a new kind of probability distribution, which we call a Poisson distribution with parameter λ .

$$\boxed{f_X(x) = \frac{\lambda^x e^{-\lambda}}{x!}}, \quad x = 0, 1, 2, 3, \dots$$

Note that though x can take any non-negative value, in practice it is small, but even if it is large we can always use Stirlings approximation to $x!$. Let us now compute the mean and variance of the Poisson random variable. In fact we write $X \sim \text{Poisson}(\lambda)$, i.e. X follows a Poisson distribution with parameter λ . Let us use the m.g.f technique and leave the brute force technique to the reader.

$$\begin{aligned} m_X(t) &= \sum_{x=0}^{\infty} e^{tx} \frac{\lambda^x e^{-\lambda}}{x!} \\ &= e^{-\lambda} \sum_{x=0}^{\infty} \frac{(\lambda e^t)^x}{x!} \\ &= e^{-\lambda} e^{\lambda e^t} \end{aligned}$$

Now

$$\begin{aligned} m'_X(t) &= e^{-\lambda} \lambda e^{\lambda e^t} e^t = \lambda e^{-\lambda} e^t e^{\cancel{\lambda e^t}} \\ \boxed{m'_X(0) = \mu_X = \lambda} \end{aligned}$$

$$m''_x(t) = \lambda e^{-\lambda} [e^t e^{\lambda e^t} + \lambda e^t e^{\lambda e^t} e^t]$$

$$m''_x(0) = \lambda e^{-\lambda} [e^\lambda + \lambda e^\lambda]$$

$$\boxed{m''_x(0) = \lambda + \lambda^2}$$

Thus

$$\begin{aligned} \text{Var}(x) &= E(x^2) - (\mu_x)^2 \\ &= m''_x(0) - (\mu_x)^2 \\ &= \lambda + \lambda^2 - \lambda^2 \\ &= \lambda \end{aligned}$$

$$\boxed{\text{Var}(x) = \lambda}$$

I find this amazing, both mean and variance have the same value, λ the parameter of the Poisson distribution.

In his small book called "Facts from Figures", M. J. Moroney almost magically brings up the Poisson Distribution, for a discrete random variable X , with countable range, i.e.

$$\begin{aligned} \text{Range } X &= \{0, 1, 2, 3, \dots\} \\ &= \mathbb{N} \cup \{0\}. \end{aligned}$$

Observe that,

$$+ e^{-\lambda} e^\lambda = e^{-\lambda} \left(1 + \frac{\lambda}{1!} + \frac{\lambda^2}{2!} + \dots\right)$$

$$1 = e^{-\lambda} e^\lambda = e^{-\lambda} \left(1 + \frac{\lambda}{1!} + \frac{\lambda^2}{2!} + \dots + \frac{\lambda^n}{n!} + \dots\right)$$

$$1 = e^{-\lambda} + \frac{\lambda e^{-\lambda}}{1!} + \frac{\lambda^2 e^{-\lambda}}{2!} + \dots + \frac{\lambda^n e^{-\lambda}}{n!} + \dots$$

Observe that each term on the rhs is a term of the Poisson

distribution for different values of x . In fact to show that how Poission distribution works effectively Moroney tested it on ~~for~~ a data collected by the great statistician R.A. Fisher. This a data on the death of a cavalryman getting killed by a horsekick in the course of a year. The data has been collected over twenty years from ten army corps (of Britain I guess) thus has 200 readings

Army Corp	Years	1	2	3	4	5	6	7	8	9	10	11, 12 ... 20	
		no of death											
1													
2													
3													
.													
10													

So the data is represent in a 10×20 matrix, which has 200 entries. It can be summarized as follows

No of deaths	Frequency of occurrence of such deaths [No cells in the above matrix with the given no]
0	109
1	65
2	22
3	3
4	1
5	0
6	0

$$\text{Total death} = (0 \times 109) + (1 \times 65) + (2 \times 22) + (3 \times 3) + (4 \times 1) + (5 \times 0) + (6 \times 0)$$

= 122 (in twenty years) among 200 observations

So average death is = $\frac{122}{200} = 0.61$
per year per corps

So $\lambda = 0.61$ in this case. Once we fix it we have

$e^{-\lambda} \approx 0.543$. Let us assume that number of deaths is a random variable following Poisson with the assumption that $\lambda = 0.61$ is fixed.

No of deaths	Poisson Prob.	Poisson freq. = 200 x Poisson Prob	Actual
0	0.543	109	109
1	0.331	66.3	65
2	0.101	20.2	22
3	0.021	4.1	3
4	0.003	0.6	1

[Table is taken from page 98 of "Facts and Figures" by M.J. Moroney, Pelican, 1951]

So you see this data can be indeed very well modelled very well by Poisson distributions.

"Thus statistics works"

—x—



Lecture 6 : Distributions of Continuous Random Variables

In this chapter we will be more focused with collecting and collating information, rather than delving into conceptual issues. Let us note that we shall list in this chapter some important distributions associated with continuous random variables. We are not going to make detailed computation of mean and variances or mgf of each and every distribution, that we list here. But we will of course do so for all some of the most well known ones, like the exponential distribution, Gamma distribution and above all the jewel of statistics, the normal distribution.

Let's start with a simple one first, the uniform distribution. This is a distribution whose density function remains constant over an interval

$$f_x(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$

It is clear that f_x is a density function. Observe that

$$f_x(x) \geq 0 \quad \forall x \text{ and}$$

$$\int_{-\infty}^{\infty} f_x(x) dx = \int_a^b \frac{1}{b-a} dx$$

$$\text{Further } E(x) = \int_a^b x \frac{1}{b-a} dx = \frac{1}{b-a} \left(\frac{b^2}{2} - \frac{a^2}{2} \right) = \frac{b+a}{2}$$

while $\text{Var}(x) = \frac{(b-a)^2}{12}$ (Compute it yourself)

The mgf for the uniform distribution is given as

$$m_X(t) = E[e^{tx}] = \frac{e^{bt} - e^{at}}{(b-a)t}$$

Before we discuss the normal distribution, let us also mention few important distributions which are used in statistics.

Section 1: The Gamma Distribution and Exponential distribution

The Gamma distribution depends on the Gamma function

The Gamma Function is given as

$$\Gamma(z) = \int_0^\infty e^{-x} x^{z-1} dx, z > 0$$

If $z \in \mathbb{N}$ then

$$\Gamma(nz) = (z-1)!$$

Further if $z = \frac{1}{2}$ then it is well-known that $\Gamma(\frac{1}{2}) = \sqrt{\pi}$

The Gamma function, is an improper integral. For more details see any good book on real analysis or advanced calculus. In fact

$$\boxed{\Gamma(z+1) = z\Gamma(z)}$$

Just for fun let us see why $\Gamma(\frac{1}{2}) = \sqrt{\pi}$. Note that

$$\Gamma(\frac{1}{2}) = \int_0^\infty \frac{e^{-x}}{\sqrt{x}} dx. \quad \text{Set } x = t^2, \text{ i.e. } dx = 2t dt$$

$$\therefore \Gamma(\frac{1}{2}) = \int_0^\infty \frac{e^{-t^2}}{t} 2t dt = 2 \int_0^\infty e^{-t^2} dt$$

There is a famous formula in calculus, i.e

$$\int_0^\infty e^{-t^2} dt = \frac{\sqrt{\pi}}{2}.$$

Thus $\Gamma(\frac{1}{2}) = \sqrt{\pi}$.

A random variable X is said to follow the Gamma distribution
if it has the density

$$f_X(x) = \frac{\lambda}{\Gamma(r)} (\lambda x)^{r-1} e^{-\lambda x}, \quad \begin{cases} x > 0 \\ x \neq 0 \end{cases} \quad (\lambda > 0, r > 0)$$

Here r is the parameter. We say that $X \sim \text{Gamma}(r; \lambda)$.
What happens if $r=1$. Then

$$f_X(x) = \frac{\lambda}{\Gamma(1)} (\lambda x)^0 e^{-\lambda x}, \quad x \geq 0, \lambda > 0$$

$$\therefore f_X(x) = \lambda e^{-\lambda x} \quad \rightarrow \textcircled{*}, \quad x \geq 0$$

Note that $\Gamma(1) = 1$, since $0! = 1$ (By convention)

The expression in $\textcircled{*}$ is also a p.d.f which is called exponential distribution. Note for the exponential pdf $f_X(0)$ can be defined as λ ~~0~~, since $\lim_{x \rightarrow 0} f_X(x) = \lambda$. ~~This~~

The exponential distribution has the following amazing property called "Memoryless". Let us see why it is called so.

Let $X \sim \text{exp}(\lambda)$, i.e. X is a random variable with exponential distribution with parameter λ . Then

$$P[X > a+b | X > a] = P[X > b]. \quad a > 0, b > 0$$

It does not keep the information $X > a$, used in the conditioning.

$$\begin{aligned} P[X > a+b | X > a] &= \frac{P[X > a+b \cap X > a]}{P[X > a]} \\ &= \frac{P[X > a+b]}{P[X > a]} \end{aligned}$$

$$\begin{aligned}
 P[X > a+b] &= 1 - P[X \leq a+b] \\
 &= 1 - \int_0^{a+b} \lambda e^{-\lambda x} dx \\
 &= 1 - \left[\lambda \int_0^{a+b} e^{-\lambda x} dx \right] \\
 &= 1 - \left[\lambda \left[\frac{e^{-\lambda x}}{-\lambda} \right] \right]_0^{a+b} \\
 &= 1 - \cancel{\lambda} \left[-e^{-\lambda x} \right]_0^{a+b} \\
 &= 1 - \left[-e^{\lambda(a+b)} + 1 \right] \\
 &= 1 + e^{-\lambda(a+b)} - 1 \\
 &= e^{-\lambda(a+b)}
 \end{aligned}$$

Similarly $P[X > a] = e^{-\lambda a}$. Thus

$$P[X > a+b | X > a] = \frac{e^{-\lambda(a+b)}}{e^{-\lambda a}} = e^{-\lambda b} = P[X > b].$$

The idea of memoryless-ness is used in reliability theory. For example let X denote the random variable which represents the lifetime of a component in a device counted for example in hours. Then

$P[X > a+b | X > a]$ seeks to find the probability, that given that the component has functioned for more than a hours, what is the probability that it will work for more than additional b hours, i.e. it will work for $a+b$ hours. Memory less-ness tells us that this probability is same as the probability that the component has worked more than b hours.

For seeing the application of exponential distribution see for example the book "Probability and Statistics; with reliability, queuing and computer science applications" by Kishore Trivedi (Wiley 2002).

We ask the reader to compute the mean and variance of an exponential distribution. The answers are

$$E[x] = \lambda \text{ and } Var(x) = \frac{1}{\lambda^2}$$

For the Gamma distribution, we shall compute the mean and variance using mgf technique. Of course we have

$$m_x(t) = E[e^{tx}]$$

$$= \int_0^\infty e^{tx} \frac{\lambda^r}{\Gamma(r)} e^{-\lambda x} x^{r-1} dx$$

$$= \int_0^\infty \frac{\lambda^r}{\Gamma(r)} e^{-(\lambda-t)x} x^{r-1} dx$$

$$= \frac{\lambda^r}{(\lambda-t)^r} \int_0^\infty \frac{1}{\Gamma(r)} e^{-(\lambda-t)x} x^{r-1} dx$$

$$= \lambda^r \int_0^\infty \frac{1}{\Gamma(r)} x^{r-1} e^{-(\lambda-t)x} dy$$

$$\text{Set } (\lambda-t)x = y \Rightarrow dy = (\lambda-t)dx$$

$$\therefore m_x(t) = \frac{\lambda^r}{(\lambda-t)^r} \int_0^\infty \frac{(\lambda-t)^r}{\Gamma(r)} x^{r-1} e^{-(\lambda-t)x} dx \quad (\text{Assume } \lambda > t)$$

$$= \frac{\lambda^r}{(\lambda-t)^r} \frac{1}{\Gamma(r)} \int_0^\infty y^{r-1} e^{-y} dy$$

$$= \frac{\lambda^r}{(\lambda-t)^r} \frac{1}{\Gamma(r)} \Gamma(r) = \frac{\lambda^r}{(\lambda-t)^r}$$

$$\therefore m'_x(t) = \lambda^r (\lambda-t)^{-r} r^{r-1} \Rightarrow \boxed{E(x) = m'_x(0) = \frac{r}{\lambda}}$$

To compute $E[X^2]$, we need to compute $m''_X(0)$. Observe that

$$m''_X(t) = r(r+1) \lambda^r (\lambda-t)^{-r-2}$$

Hence

$$\begin{aligned} m''_X(0) &= \frac{r(r+1)}{\lambda^2} \\ \therefore \text{Var}(x) &= \frac{r(r+1)}{\lambda^2} - \frac{r^2}{\lambda^2} = \frac{r}{\lambda^2} \end{aligned}$$

The Gamma distribution plays a key role in queueing theory.

Section 2: The Beta Distribution

The ~~notion~~ notion of Beta distribution depends on the idea of ~~a~~ Beta-functions an important class of functions in mathematical analysis, which is given as

$$B(a, b) = \int_0^1 x^{a-1} (1-x)^{b-1} dx, \quad a>0, b>0.$$

To begin with we will first explore the relation between Beta and Gamma function. The ~~realti~~ relation is

$$B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

(See any book on advanced calculus for a proof. You need double integrals for proving this). The Beta distribution has the p.d.f

$$f_X(x) = \frac{1}{B(a,b)} x^{a-1} (1-x)^{b-1}; \quad 0 < x < 1$$

with $a>0$ & $b>0$.

It is of course simple to show that f_X is a p.d.f

We say $X \sim \text{Beta}(a, b)$. If $a=1$ & $b=1$, then the Beta distribution is just the uniform distribution. The distribution function or cdf can be in fact very compactly represented

$$F_X(x) = \begin{cases} 0, & x \leq 0 \\ \int_0^x \frac{1}{B(a,b)} z^{a-1} (1-z)^{b-1} dz, & x \in (0,1) \\ 1, & x \geq 1. \end{cases}$$

Let us compute mean and variance of the beta distribution. We present here the approach given in the book "Introduction to the Theory of Statistics" by Mood, Graybill and Boes, McGraw-Hill, 1974 Third Edn. We first compute

$$\begin{aligned} E[X^k] &= \frac{1}{B(a,b)} \int_0^1 x^k x^{a-1} (1-x)^{b-1} dx \\ &= \frac{1}{B(a,b)} \int_0^1 x^{k+a-1} (1-x)^{b-1} dx \\ &= \frac{B(k+a, b)}{B(a, b)} = \frac{\Gamma(k+a+b)}{\cancel{\Gamma(k+a)\Gamma(b)}} \cdot \frac{\cancel{\Gamma(a+b)}}{\Gamma(a)\Gamma(b)} \\ &= \frac{\Gamma(k+a)\Gamma(b)}{\Gamma(k+a+b)} \cdot \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \\ \therefore E[X^k] &= \frac{\Gamma(k+a)\Gamma(a+b)}{\Gamma(k+a+b)\Gamma(a)} \end{aligned}$$

$$\begin{aligned} \therefore \text{For } k=1 \quad E[X] &= \frac{\Gamma(a+1)\Gamma(a+b)}{\Gamma(a)\Gamma(a+b+1)} \\ &= \frac{a\Gamma(a)\Gamma(a+b)}{\Gamma(a)\cdot (a+b)\Gamma(a+b)} \\ &= \frac{a}{a+b} \end{aligned}$$

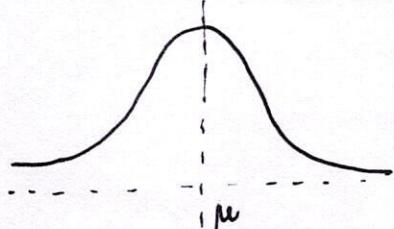
Now we will set $k=2$

$$\begin{aligned}
 E[x^2] &= \frac{\Gamma(a+2)\Gamma(a+b)}{\Gamma(2+a+b)\Gamma(a)} \\
 &= \frac{(a+\frac{1}{2})\Gamma(a+1)\Gamma(a+b)}{(a+\frac{b}{2}+1)\Gamma(a+b+1)\Gamma(a)} \\
 &= \frac{a(a+1)\Gamma(a)\Gamma(a+b)}{(a+b)(a+b+1)\Gamma(a+b)\Gamma(a)} \\
 &= \frac{a(a+1)}{(a+b)(a+b+1)} \\
 \text{Var}(x) &= \frac{a(a+1)}{(a+b)(a+b+1)} - \left(\frac{a}{a+b}\right)^2 \\
 &= \frac{ab}{(a+b)^2(a+b+1)}
 \end{aligned}$$

So we now see the advantage of computing $E[x^n]$

Section 3: Normal distribution

Normal distribution is the poster-boy of statistics. Even many members of the general public have heard about it and will immediately recognize the bell-shaped curve, which you see on the left. On a more formal level we say that a continuous random variable X has a normal distribution if it has a p.d.f.



$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad x \in \mathbb{R}$$

Here μ and σ are the parameters, and we say or rather that X follows a normal distribution with mean μ and parameter σ^2 and symbolically $X \sim N(\mu, \sigma^2)$.

It will turn out that μ and σ^2 are the ~~not~~ mean and variance of a normal random variable. We shall call a continuous random variable, normal random variable, if it follows the normal distribution. Mind you, ~~mathematicians~~ mathematicians do not like the term normal distribution. They want to use the term Gaussian distribution, after Karl, Federich Gauss, the great mathematician who introduced the idea ~~with~~ of the normal curve while tabulating astronomical data. Gauss was trying to measure the radius of the moon. As he repeated his measurement he got a different value (even if the difference was very less). When he plotted the values of his measurement he found that they seem to fall on a bell-shaped curve which the statisticians called as ~~a~~ normal curve. The normal curve is also called an error curve as errors in measurements lie on the bell-shaped curve. We will say more on this later.

First let us show that f_x is a p.d.f. The fact that $f_x(x) \geq 0$ is obvious. Now observe that

$$\int_{-\infty}^{\infty} f_x(x) dx = \int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$$

$$\text{Let } \frac{x-\mu}{\sigma} = t \Rightarrow dx = \sigma dt. \text{ Hence}$$

$$\begin{aligned} \int_{-\infty}^{\infty} f_x(x) dx &= \int_{-\infty}^{\infty} e^{-\frac{t^2}{2}} \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} e^{-t^2} dt \\ &= \frac{\sqrt{\pi}}{\sqrt{\pi}} = 1. \quad (\because \int_{-\infty}^{\infty} e^{-t^2} dt = 1, \text{ famous result in analysis}) \end{aligned}$$

Of course that shows that f_x is indeed a p.d.f.

Suppose you do not know ~~to~~ the fact that $\int_{-\infty}^{\infty} e^{-t^2} dt = \sqrt{\pi}$.

Then in that case we ^{can} show that $\int_{-\infty}^{\infty} f_x(x) dx = 1$; using the techniques double integral. This attempt is made in Mood, Graybill Boes, [though I believe they tacitly use the above formula].

Theorem 6.1 : If $X \sim N(\mu, \sigma^2)$, Then

$$E[X] = \mu, \text{Var}[X] = \sigma^2$$

$$\text{and } \mu t + \frac{1}{2} \sigma^2 t^2$$

$$m_X(t) = e^{t\mu + \frac{1}{2}\sigma^2 t^2}$$

Proof: You can of course directly compute the mean and the variance directly. However we will compute the $m_X(t)$ and through that compute μ & σ^2 through that.

$$\begin{aligned} m_X(t) &= E[e^{tx}] \\ &= E[e^{tx - t\mu + t\mu}] \\ &= e^{t\mu} E[e^{t(x-\mu)}] \\ &= e^{t\mu} \int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx \\ &= e^{t\mu} \int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{t(x-\mu)} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx \\ &= e^{t\mu} \int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} \end{aligned}$$

$$\begin{aligned} \text{Now } (x-\mu)^2 - 2\sigma^2 t^2 (x-\mu) &= (x-\mu)^2 - 2\sigma^2 t (x-\mu) + \sigma^4 t^2 - \sigma^4 t^2 \\ &= (x-\mu - \sigma^2 t)^2 - \sigma^4 t^2 \\ &\quad - \left[\frac{(x-\mu - \sigma^2 t)^2}{2\sigma^2} \right] dx \end{aligned}$$

$$\begin{aligned} \therefore m_X(t) &= e^{t\mu} e^{\frac{\sigma^2 t^2}{2}} \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{(x-\mu - \sigma^2 t)^2}{2\sigma^2}} dx \\ &= e^{t\mu + \frac{\sigma^2 t^2}{2}} \\ &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{(x-(\mu+\sigma^2 t))^2}{2\sigma^2}} dx = 1 \quad (\text{Think why!!}). \end{aligned}$$

$$\therefore m_X(t) = e^{t\mu + \frac{1}{2}\sigma^2 t^2}$$

$$m'_X(t) = (\mu + \sigma^2 t) e^{t\mu + \frac{1}{2}\sigma^2 t^2}$$

$$\boxed{E(x) = m'_X(0) = \mu}$$

$$m''_X(t) = \sigma^2 e^{t\mu + \frac{1}{2}\sigma^2 t^2} + (\mu + \sigma t)(\mu + \sigma t) e^{t\mu + \frac{1}{2}\sigma^2 t^2}$$

$$m''_X(0) = \sigma^2 + \mu^2 = E(x^2)$$

$$\therefore \text{Var}(x) = E(x^2) - [E(x)]^2$$

$$= \sigma^2 + \mu^2 - \mu^2$$

$$\boxed{\text{Var}(x) = \sigma^2}$$

If $\mu=0$ and $\sigma^2=1$, we say x follows a standard normal distribution, and its pdf is

$$f_X(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{1}{2}x^2}$$

$$\boxed{f_X(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}}$$

In fact if we consider $X \sim N(\mu, \sigma^2)$, then the

random variable $Z = \frac{X-\mu}{\sigma} \sim N(0,1).$

The distribution function of a standard normal variable has a symbol $\Phi(x)$, ie

$$\Phi(x) = P(X \leq x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{1}{2}x^2} dx.$$

Now consider any $X \sim N(\mu, \sigma^2)$

then

$$P(x \leq x) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx$$

Now observe

$$P(x \leq x) = P\left(\frac{x-\mu}{\sigma} \leq \frac{x-\mu}{\sigma}\right) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx$$

Now $\frac{x-\mu}{\sigma} \sim N(0, 1)$. Thus if $X \sim N(\mu, \sigma^2)$

Then

$$F_x(x) = \Phi\left(\frac{x-\mu}{\sigma}\right)$$

$$\boxed{F_x(x) = \Phi\left(\frac{x-\mu}{\sigma}\right)}$$

The values of the function Φ is tabulated and these tables are available.

$$\therefore P(a \leq x \leq b) = F_x(b) - F_x(a).$$

$$= \Phi\left(\frac{b-\mu}{\sigma}\right) - \Phi\left(\frac{a-\mu}{\sigma}\right)$$

$$\therefore \boxed{P(a \leq x \leq b) = \Phi\left(\frac{b-\mu}{\sigma}\right) - \Phi\left(\frac{a-\mu}{\sigma}\right)}$$

For a more historical aside, let us note that De Moivre also had independently discovered the normal curve as a distribution of errors in measurement. Now observe that

$$\begin{aligned}\Phi(x) &= P(X \leq x) = 1 - P(X > x) \\ &= 1 - \int_{x}^{+\infty} f_X(x) dx\end{aligned}$$

Now using Fig B we have

$$P(X > x) = P(X \leq -x)$$

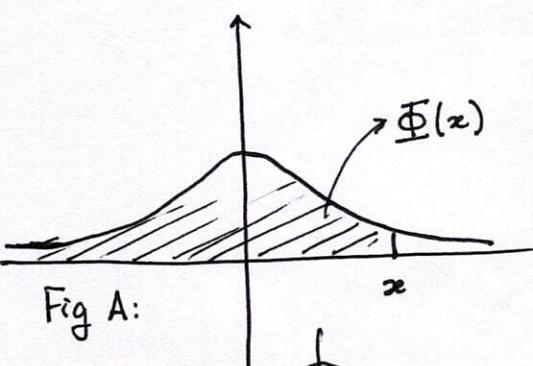


Fig A:

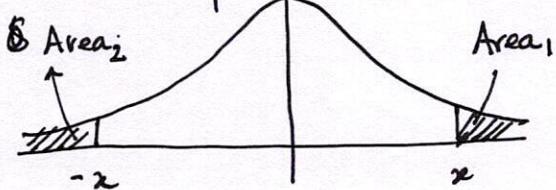


Fig B: Area 1 = Area 2 by symmetry.

Thus

$$\bar{\Phi}(x) = 1 - P(X \leq -x) \quad \left[\text{Always keep in mind that } P(X=x)=0 \right]$$
$$= 1 - \bar{\Phi}(-x)$$

$$\bar{\Phi}(x) = 1 - \bar{\Phi}(-x)$$

The binomial distribution which we know as a discrete distribution actually links up beautifully with normal distribution when the number of trials is very large. ~~and thus links up beautifully with the normal distib~~

Thus the binomial distribution is a link between the discrete and the continuous. The following result asserts this

Theorem 8.1: Let $X \sim B(n, p)$, then

$$Z = \frac{X - np}{\sqrt{np(1-p)}}$$

approaches the standard normal distribution with mean zero and variance 1 as $n \rightarrow \infty$.

Proof: How do we detect whether a given r.v. tends to some other distribution. Here we use the following approach.
(A sketch only)
We show that as $n \rightarrow \infty$, the mgf of the binomial random variable will tend to the mgf of standard normal. We know from an assignment problem, that

$$\bullet \quad m_{\frac{X+a}{b}}(t) = e^{\frac{at}{b}} m_X\left(\frac{t}{b}\right) \quad \begin{array}{l} \text{See practice problems set} \\ \text{No 2. Problem} \end{array}$$

Since $X \sim B(n, p)$, we have

$$m_X(t) = (pe^t + (1-p))^n, \quad 1-p = q \quad (\text{in the notes})$$

$$\therefore \text{For } Z \text{ we have } m_Z(t) = e^{-\frac{np}{\sigma^2} t} \left[pe^{t/\sigma} + (1-p) \right]^n$$

$$\text{where } \sigma = \sqrt{np(1-p)}.$$

Let us set $np = \mu$. Then

$$m_Z(t) = e^{-\frac{\mu t}{\sigma}} \left[1 + p(e^{\frac{t}{\sigma}} - 1) \right]^n$$

In order to proceed to the limit we need to delink $m_Z(t)$ since we will consider $n \rightarrow \infty$.

$$\begin{aligned} &= \ln(m_Z(t)) \\ \therefore \ln(m_Z(t)) &= -\frac{\mu t}{\sigma} + n \ln \left[1 + p \left(e^{\frac{t}{\sigma}} - 1 \right) \right] \\ &\approx -\frac{\mu t}{\sigma} + n \ln \left[1 + p \left(\frac{t}{\sigma} + \frac{t^2}{2\sigma^2} + \frac{1}{3!} \frac{t^3}{\sigma^3} + \dots \right) \right] \end{aligned}$$

Remember the series $\ln(1+x) = x - \frac{1}{2}x^2 + \frac{1}{3}x^3 - \frac{1}{4}x^4 + \dots$ converges $|x| < 1$. Assuming that such conditions hold we have (Think why)

$$\begin{aligned} \ln(m_Z(t)) &= -\frac{\mu t}{\sigma} + np \left[\frac{t}{\sigma} + \frac{1}{2!} \frac{t^2}{\sigma^2} + \frac{1}{3!} \frac{t^3}{\sigma^3} + \dots \right] \\ &= \cancel{-\frac{\mu t}{\sigma}} + -\frac{np^2}{2!} \left[\frac{t}{\sigma} + \frac{1}{2!} \frac{t^2}{\sigma^2} + \frac{1}{3!} \frac{t^3}{\sigma^3} + \dots \right]^2 \\ &\quad + \frac{np^3}{3!} \left[\frac{t}{\sigma} + \frac{1}{2!} \frac{t^2}{\sigma^2} + \frac{1}{3!} \frac{t^3}{\sigma^3} + \dots \right]^3 - \dots \end{aligned}$$

Now by collecting collecting the powers of t we obtain

$$\begin{aligned} \ln(m_Z) &= \left(-\frac{\mu}{\sigma} + \frac{np}{\sigma} \right) t + \left(\frac{np}{2\sigma^2} - \frac{np^2}{2\sigma^2} \right) t^2 \\ &\quad + \left(\frac{np}{6\sigma^3} - \frac{np^2}{2\sigma^3} + \frac{np^3}{3\sigma^3} \right) t^3 + \dots \\ &= \frac{1}{\sigma^2} \left(\frac{np - np^2}{2} \right) t^2 + \left(\frac{np}{6\sigma^3} - \frac{np^2}{2\sigma^3} + \frac{np^3}{3\sigma^3} \right) t^3 + \dots \\ &= \frac{1}{\sigma^2} \cdot \frac{\sigma^2}{2} t^2 + \frac{n}{\sigma^3} \left[\frac{p - 3p^2 + 2p^3}{6} \right] t^3 \end{aligned}$$

$$\begin{aligned}
 \therefore \ln(m_Z(t)) &= \frac{1}{2}t^2 + \frac{n}{\sigma^3} \left[\frac{p - 3p^2 + 2p^3}{6} \right] t^3 + \dots \\
 &= \frac{1}{2}t^2 + \frac{n}{(\sqrt{np(1-p)})^3} \left[\frac{p - 3p^2 + 2p^3}{6} \right] t^3 + \dots \\
 &= \frac{1}{2}t^2 + \frac{n}{n^{3/2} p^{3/2} (1-p)^{3/2}} \left[\frac{p - 3p^2 + 2p^3}{6} \right] t^3 + \dots \\
 \text{So } \ln m_Z(t) &= \frac{1}{2}t^2 + \frac{1}{\sqrt{n} (p^{3/2} (1-p)^{3/2})} \left[\frac{p - 3p^2 + 2p^3}{6} \right] t^3 + \dots
 \end{aligned}$$

We can now safely argue, that for $r > 2$ the coefficient of t^r will go to zero, as $n \rightarrow \infty$. Hence

$$\ln m_Z(t) \rightarrow \frac{1}{2}t^2 \text{ as } n \rightarrow \infty$$

as $\frac{1}{2}t^2$ is the mgf of the standard normal random variable.

What we Theorem 8.1, is a special case of the central limit theorem we will learn later on in the course.

[The sketch of the proof given here is from the book:

John. E. Freund's Mathematical Statistics by Miller & Miller]

Pearson : 2014.]

—x—

Suppose X and Y are two random variables on the same probability space. Can we give any meaning to the following question?

Do these two variables taken jointly can have distribution?

Are these two variables independent of each other or are they related? Means we ask the question

When is X and Y stochastically independent?

and

If X and Y are related then how do we measure it?

We seek to answer these questions here.

Section 1: Joint distribution of a random vector

In this chapter our aim is to study a random vector and its distribution. A random vector X is the following mapping

$$X: \Omega \rightarrow \mathbb{R}^k$$

i.e. if $\omega \in \Omega$, then $X(\omega) = (x_1, x_2, \dots, x_k)$

Note that each x_1, \dots, x_k depend on $\omega \in \Omega$, thus can be viewed as a realization of a random variable. So we can write

$$X(\omega) = (X_1(\omega), \dots, X_k(\omega))$$

where $X_i(\omega) = x_i$, $i=1, \dots, k$. In short

$$X = (X_1, \dots, X_k)$$

The expression

$$X \leq x \Leftrightarrow X_1 \leq x_1, \dots, X_k \leq x_k$$

In fact

$$\{x \leq z\} = \{x_1 \leq z_1, \dots, x_k \leq z_k\} = \bigcap_{i=1}^k \{x_i \leq z_i\}$$

The distribution function of this random vector X , is given as

$$F_X(\mathbf{z}) = F_{X_1 \dots X_k}(z_1 \dots z_k) = P(X_1 \leq z_1, \dots, X_k \leq z_k)$$

Instead of focusing on k -random variables we just focus on two random variables X and Y in order to make the discussion simpler at least for the time being. For the sake of curiosity largely we state the properties of the distribution functions of two random variables, X and Y .

- Properties of $F_{X,Y}(x,y)$ (never mind if you forget these immediately after reading)

i) $\lim_{x \rightarrow -\infty} F(x,y) = 0$, for all y & $\lim_{y \rightarrow -\infty} F(x,y) = 0$, for all x .

and

$$\begin{array}{c} \lim_{x \rightarrow \infty} F(x,y) = 1 \\ \lim_{y \rightarrow \infty} F(x,y) = 1 \end{array}$$

ii) $F(x,y)$ is right continuous in each variable

$$\lim_{h \rightarrow 0^+} F(x+h, y) = \lim_{h \rightarrow 0^+} F(x, y+h) = F(x, y).$$

Further one can relate the computing the probability using distribution functions. i.e. let $x_1 < x_2$ & $y_1 < y_2$, such that

$$P[x_1 < X \leq x_2, y_1 < Y \leq y_2]$$

$$= F(x_2, y_2) - F(x_2, y_1) - F(x_1, y_2) + F(x_1, y_1)$$

- Discrete bi-variate r.v.

Let X and Y be discrete random variable. Then its' the joint probability mass function or joint pmf of X and Y , is given by

$$f_{X,Y}(x,y) = P[X=x, Y=y]$$

Of course one must have

- $f_{X,Y}(x,y) \geq 0, \quad \forall (x,y) \in \text{Ran}(X) \times \text{Ran}(Y)$

- $\sum_x \sum_y f(x,y) = 1 = \sum_y \sum_x f(x,y)$

- $\sum_x \sum_y f_{X,Y}(x,y) = 1 = \sum_y \sum_x f_{X,Y}(x,y).$

Example: 6.1: Consider an urn having 3 red, 4 black and 1 green ball. Two balls are drawn at random.

let X be the random variable denoting the number of green/red balls among the two drawn balls and Y denotes the number of black balls. Can we write down the joint probability distribution

X	Y	$X \setminus Y$	0	1	2	$\sum_y f_{X,Y}(x,y)$
		$X \setminus Y$	0	$\frac{4}{28}$	$\frac{6}{28}$	$\frac{10}{28}$
		0	$\frac{3}{28}$	$\frac{12}{28}$	0	$\frac{15}{28}$
		2	$\frac{3}{28}$	0	0	$\frac{3}{28}$
$\sum_x f_{X,Y}(x,y)$			$\frac{6}{28}$	$\frac{16}{28}$	$\frac{6}{28}$	1

- The details of how we fill the table.

Some parts can be filled by the reader

Probability table

$$\sum_x \sum_y f_{X,Y}(x,y)$$

(3)

There are 8 balls and we are choosing just two, ie in $\binom{8}{2}$ ways.

$$\binom{8}{2} = \frac{8!}{2! 6!} = \frac{8! \times 7! \times 6!}{2! 6!} = \frac{56}{2} = 28$$

So there are 28 all possible outcomes. Let us first compute

$$f_{X,Y}(0,0) = 0 \quad (\text{Think why!})$$

$$f(0,1) = \frac{\binom{3}{0} \binom{4}{1} \binom{1}{1}}{28} = \frac{4}{28} \quad [0! = 1 \text{ remember}]$$

$$f(0,2) = \frac{\binom{3}{0} \binom{4}{2} \binom{1}{0}}{28} = \frac{6}{28}$$

$$f(1,0) = \frac{\binom{3}{1} \binom{4}{0} \binom{1}{1}}{28} = \frac{3}{28}$$

$$f(1,1) = \frac{\binom{3}{1} \binom{4}{1} \binom{1}{0}}{28} = \frac{12}{28}$$

$$f(2,0) = \frac{\binom{3}{2} \binom{4}{0} \binom{1}{0}}{28} = \frac{6}{28}$$

$$f(2,1) = 0 \quad & f(2,2) = 0$$

I have skipped writing X, Y every time, so $f(0,1) = f_{X,Y}(0,1)$.
 Looking at the table one can see that all the properties of
 the joint pmf is satisfied. \square

Of course we can now write down the distribution function as follows

$$F_{X,Y}(x,y) = \sum_{s \leq x} \sum_{t \leq y} f(s,t), \quad x \in \mathbb{R} \text{ & } y \in \mathbb{R}.$$

Exercise. Using the above table in the previous page compute

$$F_{X,Y}(2,2).$$

Now we can of course think of independent continuous random variables X, Y . Thus $f_{X,Y}(x,y)$ is a joint probability density function if

$$f_{X,Y}(x,y) \geq 0 \quad \forall (x,y) \in \text{Ran}(X) \times \text{Ran}(Y)$$

and

$$\int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f_{X,Y}(x,y) dy dx = 1$$

which is same as saying

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x,y) dx dy = 1.$$

The distribution function can now be written as

$$F_{X,Y}(x,y) = \int_{-\infty}^y \int_{-\infty}^x f_{X,Y}(z,w) dz dw, \quad x \in \mathbb{R}, y \in \mathbb{R}$$

$$= P(X \leq x, Y \leq y).$$

For your happiness only.
 $f_{X,Y}(x,y) = \frac{\partial^2 F_{X,Y}(x,y)}{\partial x \partial y}$

Large We urge the reader to extend this definition for k random variables.

Example: Continuous case: Consider

$$f(x,y) = K(x+y), \quad 0 < x < 1, 0 < y < 1$$

If $K \geq 0$, then $f(x,y) \geq 0, \forall (x,y) \in [0,1] \times [0,1]$

Now to be a joint pmf density we must have

$$\int_0^1 \int_0^1 K(x+y) dx dy = 1$$

$$\therefore K \int_0^1 \left[\int_0^1 (y+x) dx \right] dy = 1 \Rightarrow K \int_0^1 \left[y + \frac{1}{2} \right] dy = 1$$

$$\Rightarrow K \int_0^1 y dy + K \int_0^1 \frac{1}{2} dy = 1$$

$$\Rightarrow \frac{K}{2} + \frac{K}{2} = 1$$

$$\Rightarrow K = 1.$$

Knowing the joint distribution of (X, Y) can we get the individual distributions of X and Y .

The marginal distribution of X : $f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x,y) dy$

The marginal distribution of Y : $f_{Y|X}(y) = \int_{-\infty}^{\infty} f_{X,Y}(x,y) dx$

We have to show that $f_X(x)$ and $f_Y(y)$ are actually densities.
Of course we have to show that $f_X(x) \geq 0$ & $f_Y(y) \geq 0$, $\forall (x,y) \in \text{Ran}(X) \times \text{Ran}(Y)$.

$$\begin{aligned} \int_{-\infty}^{\infty} f_X(x) dx &= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} f_{X,Y}(x,y) dy \right] dx = \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x,y) dy dx = 1 \end{aligned}$$

$$\begin{aligned} \int_{-\infty}^{\infty} f_Y(y) dy &= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} f_{X,Y}(x,y) dx \right] dy \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x,y) dx dy = 1 \end{aligned}$$

We can also define all the above ideas for the discrete case.

Following on the ideas of conditional distribution probability,
we define, what is called conditional pmf and conditional p.d.f

let X and Y be two continuous random variables. Then

the conditional density of Y given $X=x$ is given as

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

This is defined if
 $f_X(x) > 0$ and
not defined $f_X(x)=0$.

Let X and Y be discrete, then, with $P[X=x_i] > 0$, $\forall x_i \in \text{Ran}(X)$, then

$$f_{Y|X}(y_j | x_i) = P[Y=y_j | X=x_i] = \frac{P[X=x_i, Y=y_j]}{P[X=x_i]}$$

Our question again remains the same. If $f_{Y|X}(y|x)$ a density of (x,y) are jointly continuous random variable. Let us check this fact. If $f_{Y|X}(\cdot|x)$ is a density function of y . The fact that $f_{Y|X}(y|x) \geq 0$, $\forall y$ is clear. Then.

$$\begin{aligned}\int_{-\infty}^{\infty} f_{Y|X}(y|x) dy &= \int_{-\infty}^{\infty} \frac{f_{Y|X}(y|x)}{f_X(x)} \frac{f_{X,Y}(x,y)}{f_X(x)} dy \\ &= \frac{1}{f_X(x)} \int_{-\infty}^{\infty} f_{X,Y}(x,y) dy \\ &= \frac{1}{f_X(x)} \cdot f_X(x) = 1\end{aligned}$$

Hence $f_{Y|X}(\cdot|x)$ is truly a density function. How can we write down the cumulative distribution function in this case. If $f_X(x) > 0$ then then the cdf of Y given $X=x$ is given as

$$F_{Y|X}(y|x) = \int_{-\infty}^y f_{Y|X}(z|x) dz$$

Let us now look into the important concept of independence of the random variables X and Y .

We say X and Y are two given random variables.

If X and Y are discrete

Independence \Leftrightarrow Joint pmf $f_{X,Y}(x,y) = f_X(x) f_Y(y)$

marginal pmf

If X and Y are continuous

Independence \Leftrightarrow Joint density = Product of marginals.
 $f_{X,Y}(x,y) = f_X(x) f_Y(y)$

$$\Rightarrow f_{Y|X}(y|x) = f_Y(y)$$

$$\& f_{X|Y}(x|y) = f_X(x)$$

A collection of k random variables X_1, X_2, \dots, X_k iff

$$f_{X_1, \dots, X_n}^{\text{P}}(x_1, x_2, \dots, x_n) = f_{X_1}(x_1) \cdots f_{X_n}(x_n).$$

Depends on whether you consider pmf or pdf.

E.g. $f_{X,Y}(x,y) = e^{-(x+y)}$; $x \geq 0, y \geq 0$

Is X and Y independent. Let us first compute the marginals

$$\begin{aligned} f_X(x) &= \int_0^\infty e^{-(x+y)} dy = e^{-x} \\ f_Y(y) &= \int_0^\infty e^{-(x+y)} dx = e^{-y} \\ f_{X,Y}(x,y) &= e^{-(x+y)} = e^{-x}e^{-y} \\ &= f_X(x)f_Y(y). \end{aligned}$$

— $\rightarrow X$ —

Expectation, Covariance and Conditional Expectation

look at finding the expectation of a random function of a random vector if $X: \Omega \rightarrow \mathbb{R}^k$ joint distribution is known.

$Y = g(X)$ be a function of random vector X and it is real valued, i.e. $g(X(\omega)) \in \mathbb{R}$, and $g: \text{Ran}(X) \rightarrow \mathbb{R}$.

We shall also study the way to measure the relation between two random variables X and Y , using an idea called covariance and correlation coefficient.

Section 1: Expectation of $g(x)$

Let $X: \Omega \rightarrow \mathbb{R}^k$ be random vector given as

$$X = (X_1, X_2, \dots, X_k)$$

If each X_1, X_2, \dots, X_k be discrete random variables and let $g: \mathbb{R}^k \rightarrow \mathbb{R}$, then expectation of $g(X)$ is given as.

$$E[g(X_1, \dots, X_k)] = \sum g(x_1, x_2, \dots, x_k) f_X(x_1, \dots, x_k)$$

If X_1, \dots, X_k are continuous random variables, then

$$E[g(X_1, \dots, X_k)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} g(x_1, \dots, x_k) f_X(x_1, \dots, x_k) dx_1 \dots dx_k.$$

Theorem 8.1: Let X and Y be two continuous random variable with $f_{XY}(\cdot, \cdot)$, representing their joint density function. Then if $E(X) < \infty$ & $E(Y) < \infty$ and $E(X+Y) < \infty$ then

$$E(X+Y) = E(X) + E(Y).$$

[• Here $g(X, Y) = X+Y$. For two variables the density function is given written as $f_{X,Y}(x, y)$.]

Proof:

$$E(X+Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x+y) f_{X,Y}(x,y) dx dy$$

$$\therefore E(X+Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x f_{X,Y}(x,y) dx dy + \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y f_{X,Y}(x,y) dx dy$$

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

$$= \int_{-\infty}^{\infty} x f_X(x) dx + \int_{-\infty}^{\infty} y f_Y(y) dy \quad \begin{array}{l} \text{Using the defn} \\ \text{of marginal distribution} \end{array}$$

$$\therefore E(X+Y) = E(X) + E(Y) \quad \square$$

Think how will you extend it to more than 2 variables. In fact one can show that

$$E\left[\sum_{i=1}^k x_i\right] = \sum_{i=1}^k E[x_i]. \longrightarrow (\#)$$

We urge the reader to prove the above result in Thm 8.1, for the discrete case.

What about $\text{Var}(X+Y)$, where X and Y are two random variables, with finite expectations.

$$\begin{aligned} \text{Var}(X+Y) &= E[(X+Y - E(X+Y))^2] \\ &= E[(X+Y - E(X)-E(Y))^2] \quad \text{Using Thm 8.1} \\ &= E[(X-E(X)) + (Y-E(Y))^2] \\ &= E[(X-E(X))^2] + E[(Y-E(Y))^2] + 2E[(X-E(X))(Y-E(Y))] \\ &= E[(X-E(X))^2] + \{E[(Y-E(Y))^2]\} \\ &\quad + 2E[(X-E(X))(Y-E(Y))] \end{aligned}$$

(2) Using (#)

$$\therefore \text{Var}(x+y) = \text{Var}(x) + \text{Var}(y) + 2 E[(x - E(x))(y - E(y))]$$

We shall now focus on the term

$$\begin{aligned} & E[(x - E(x))(y - E(y))] \\ &= E[xy - xE(y) - yE(x) + E(x)E(y)] \\ &= E[xy] - E[y]E[x] - E[y]E[x] + E[x]E[y] \\ &= E[xy] - E[x]E[y] \end{aligned}$$

$$\therefore \text{Var}(x+y) = \text{Var}(x) + \text{Var}(y) + 2 [E[xy] - E[x]E[y]]$$

What happens if x and y are independent. What happens to the term $E[xy] - E[x]E[y]$ in such a case.

Let just check out this for the continuous case. We will first compute $E[xy]$, thus

$$\begin{aligned} E[xy] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f_{x,y}(x,y) dx dy \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f_x(x) \cancel{f_y(y)} dx dy \quad \left[\begin{array}{l} f_{xy}(x,y) \\ = f_x(x) f_y(y) \\ \text{as } x \text{ and } y \text{ are independent} \end{array} \right] \\ &= \left(\int_{-\infty}^{\infty} x f_x(x) dx \right) \left(\int_{-\infty}^{\infty} y f_y(y) dy \right) \\ &= E[x] E[y] \end{aligned}$$

$$\therefore \boxed{E[xy] = E[x] E[y]}$$

(3)

Theorem 8.2 Let X and Y are random variables, with finite expectations. Then If X and Y are independent then

$$E[XY] = E[X]E[Y].$$

Thus $E[XY] - E[X]E[Y] = 0$ if X and Y are independent

So we have the following: If X and Y are independent random variables, with finite expectations, then

$$\boxed{\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)}$$

In general we shall call the term

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

as the Covariance of X and Y as it seems to measure the expected ^{joint} variation of X and Y from the mean their mean values. In symbols

$$\text{Cov}(X, Y) = E[(X - E(X))(Y - E(Y))]$$

$$\therefore \text{Cov}(X, Y) = E[XY] - E[X]E[Y]$$

Thus:

$$\rightarrow \text{If } X \text{ and } Y \text{ are independent, then } \text{Cov}(X, Y) = 0$$

The converse is not true. We shall provide examples but before that we state, that for any two random variables X and Y

$$\boxed{\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)}$$

- Note Covariance is symmetric in X and Y . Thus

~~Cov(X, Y)~~
Cov(X, Y) = Cov(Y, X)

Here we give two examples. One for the discrete case and other for the continuous case to show that $E[XY] = E[X]E[Y]$ is a necessary condition for X and Y to be independent, but not sufficient.

Eg: (from Math. Stat. Freund's Math. Stat., by Miller & Miller)

Consider two discrete random variables X and Y , taking the values $-1, 0, +1$.

	x	-1	0	1	$P(Y=y)$
y	-1	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{6}$	$\frac{2}{3}$
	0	0	0	0	0
	1	$\frac{1}{6}$	0	$\frac{1}{6}$	$\frac{1}{3}$
$P(x=x)$		$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	1 → Total prob.

The above probability table gives the joint distribution of X and Y .

$$\text{we have } E(X) = -1 \times \frac{1}{3} + 0 \times \frac{1}{3} + 1 \times \frac{1}{3} = 0$$

$$\& E(Y) = -1 \times \frac{2}{3} + 0 \times 0 + 1 \times \frac{1}{3} = -\frac{1}{3}$$

$$E(XY) = (-1) \cdot (-1) \cdot \frac{1}{6} + (-1) \cdot (0) \cdot \frac{1}{3} + (-1) \cdot (1) \cdot \frac{1}{6} \\ + 1 \cdot (-1) \cdot \frac{1}{6} + (1 \cdot 0) \cdot \frac{1}{3} = 0$$

$$\therefore \boxed{E(XY) - E(X)E(Y) = 0}$$

Now consider $x=-1, y=-1$, then $f_{XY}(x,y) = \frac{1}{6}$

$$f_X(x,y) = f_X(-1,-1) = \frac{1}{6}$$

$$\text{Now } f_X(x) = f_X(-1) = \frac{1}{3}$$

$$f_Y(y) = f_Y(-1) = \frac{2}{3}$$

$$\therefore \boxed{f_X(x)f_Y(y) = \frac{2}{9} \neq \frac{1}{6} = f_{XY}(x,y)} \Rightarrow X \text{ and } Y \text{ are not independent.}$$

e.g2: Let X & Y be two continuous random variables given as

$$X = \sin 2\pi U$$

$$Y = \sin \cos 2\pi U$$

where U is a uniformly distributed random variable in $(0, 1)$; i.e.

$$f_U(u) = \begin{cases} 1 & \text{if } u \in (0, 1) \\ 0 & \text{otherwise.} \end{cases}$$

If X is known then $U = \frac{1}{2\pi} \sin^{-1} X$, and thus Y is known.

$$\text{Now } E[Y] = \int_0^1 \cos 2\pi u du = 0 \quad \& \quad E[X] = \int_0^1 \sin 2\pi u du = 0$$

$$\text{Now } E[XY] = \int_0^1 \sin 2\pi u \cos 2\pi u du = 0$$

Note that,

$$XY = \sin 2\pi U \cos 2\pi U = g(U)$$

$$\therefore E[XY] = E[g(u)] = \int_0^1 g(u) f_u(u) du.$$

Here $E[XY] - E[X]E[Y] = 0$, though X and Y are not independent. \square

Let us call two random variables X and Y un-correlated if $\text{cov}(X, Y) = 0$.

From the above discussion we have learned the following

If X and Y are independent $\Rightarrow X$ and Y are uncorrelated
But the converse is not true

We can normalize the correlation measure, i.e. bring in the value of $\text{cov}(X, Y)$ within the interval. This leads to the introduce the notion of a Pearson's correlation coefficient, denoted by ρ_{xy} and denoted given as.

Thus we can write

$$P_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

Let us now focus on the following result, from Rohatgi [Statistical Inference]
Dover Pub. 2003.

Theorem 8.3: If $E(X^2) < \infty$ & $E(Y^2) < \infty$, then $\text{Cov}(X, Y)$
~~exists~~ exists and X and Y have finite mean, then-

$$\bullet [\text{Cov}(X, Y)]^2 \leq \sigma_X^2 \sigma_Y^2$$

$$\text{i.e. } [\text{Cov}(X, Y)]^2 \leq \text{Var}(X) \text{Var}(Y). \rightarrow \textcircled{#}$$

Proof: Set $E(X) = \mu_1$ & $E(Y) = \mu_2$.

$$\therefore (x - \mu_1)(y - \mu_2) \leq \frac{(x - \mu_1)^2 + (y - \mu_2)^2}{2}$$

\therefore In terms of random variables, then

$$(x - \mu_1)(y - \mu_2) \leq \frac{(x - \mu_1)^2 + (y - \mu_2)^2}{2}.$$

$$\therefore E[(x - \mu_1)(y - \mu_2)] \leq \frac{1}{2} \text{Var}(X) + \frac{1}{2} \text{Var}(Y)$$

$$\Rightarrow E[(x - \mu_1)(y - \mu_2)] \leq \frac{1}{2} [\text{Var}(X) + \text{Var}(Y)].$$

$$\therefore \text{Cov}(X, Y) \leq \frac{1}{2} [\sigma_X^2 + \sigma_Y^2], \quad \begin{matrix} \text{where } \sigma_X^2 = \text{Var}(X) \\ \sigma_Y^2 = \text{Var}(Y) \end{matrix}$$

Thus if $\sigma_X^2 < \infty$ and $\sigma_Y^2 < \infty$, then $\text{Cov}(X, Y)$ is bounded above, and is also finite.

Let $a, b \in \mathbb{R}$, then

$$E[a(x - \mu_1) + b(y - \mu_2)]^2 = a^2 \sigma_X^2 + 2ab \text{Cov}(X, Y) + b^2 \sigma_Y^2 \rightarrow \textcircled{#}$$

Let us assume that either $\sigma_X^2 = 0$ or $\sigma_Y^2 = 0$, then $\textcircled{#}$ holds automatically

Now let us assume that $\sigma_x^2 > 0$ & $\sigma_y^2 > 0$. The equation (6) holds for any $a \in b \in \mathbb{R}$. Thus set

$$a = -\frac{\text{Cov}(x, y)}{\sigma_x^2}$$

$$\text{Now from } \textcircled{6} \rightarrow \frac{(\text{Cov}(x, y))^2}{\sigma_x^4} \cdot \sigma_x^2 + b^2 \sigma_y^2 - 2 \cancel{b} \frac{\text{Cov}(x, y)^2}{\sigma_x^2} \geq 0 \rightarrow \textcircled{6}$$

The above expression is non-negative since

$$E[a(x-\mu_1) + b(x-\mu_2)]^2 \geq 0$$

Put $b=1$ in $\textcircled{6}$, and then we get

$$\frac{(\text{Cov}(x, y))^2}{\sigma_x^2} + \sigma_y^2 \geq 2 \frac{\text{Cov}(x, y)^2}{\sigma_x^2}$$

$$\Rightarrow \frac{\text{Cov}(x, y)^2}{\sigma_x^2} \leq \sigma_y^2$$

$$\Rightarrow \boxed{(\text{Cov}(x, y))^2 \leq \sigma_x^2 \sigma_y^2}$$

This completes the proof. \square

Theorem 8.4: $-1 \leq \rho_{xy} \leq 1$, if $\sigma_x^2 > 0$, $\sigma_y^2 > 0$ and finite.

Proof: From Theorem 8.3, $(\text{Cov}(x, y))^2 \leq \sigma_x^2 \sigma_y^2$

$$\Rightarrow \left(\frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} \right)^2 \leq 1$$

$$\Rightarrow \rho_{xy}^2 \leq 1$$

$$\Rightarrow -1 \leq \rho_{xy} \leq 1. \quad \square$$

(8)

Section 8.2: Bivariate Normal Distribution

Can a normal distribution be defined in a joint manner for random variables X and Y ? The answer turns out to be yes and before we proceed let us say little more about covariance. Sometimes one uses the symbol σ_{xy} to denote $\text{Cov}(x, y)$. In this formalism, we have $\text{Var}(x) = \text{Cov}(x, x) = \sigma_{xx} \sigma_{xx}$. Thus given any two random variables X and Y we call the following matrix Σ , as the variance given as

$$\Sigma = \begin{bmatrix} \sigma_{xx} & \sigma_{xy} \\ \sigma_{yx} & \sigma_{yy} \end{bmatrix}$$

as the variance-covariance matrix and often an helpful tool to represent variance. Let $a, b \in \mathbb{R}$, then

$$\begin{aligned} \text{Var}(ax + by) &= a^2 \text{Var}(x) + b^2 \text{Var}(y) + 2ab \text{Cov}(x, y) \\ &= a^2 \sigma_{xx} + b^2 \sigma_{yy} + 2ab \sigma_{xy} \\ &= \langle w, \Sigma w \rangle \end{aligned}$$

where $w = \begin{bmatrix} a \\ b \end{bmatrix}$. Note that as $\text{Var}(ax + by) \geq 0$, we have $\langle w, \Sigma w \rangle \geq 0$, $\forall w \in \mathbb{R}^2$. Thus Σ is a positive semidefinite matrix. Another important idea is that of the moment generating function of more than one random variable. This is defined as follows

$$m_{X,Y}(t_1, t_2) = E[e^{t_1 X + t_2 Y}]$$

Note that if X and Y are independent, then

$$\begin{aligned} m_{X,Y}(t_1, t_2) &= E[e^{t_1 X} e^{t_2 Y}] \\ &= E[e^{t_1 X}] E[e^{t_2 Y}] \\ &\stackrel{\text{---}}{=} m_X(t_1) m_Y(t_2) \\ &= m_X(t_1) m_Y(t_2). \end{aligned}$$

Let us now see how do we compute the expectations in this case.

We have

$$E[X] = \left. \frac{\partial m_{x,y}(t_1, t_2)}{\partial t_1} \right|_{(t_1, t_2) = (0,0)}$$

$$E[Y] = \left. \frac{\partial m_{x,y}(t_1, t_2)}{\partial t_2} \right|_{(t_1, t_2) = (0,0)}$$

$$E[X,Y] = \left. \frac{\partial^2 m_{x,y}(t_1, t_2)}{\partial t_1 \partial t_2} \right|_{(t_1, t_2) = (0,0)}$$

I am not providing the details here which I believe one can be checked easily by the reader. Further the above expressions can be written down by even thinking in an intuitive way. Till now the most important learning of the previous section is the following:

X and Y are independent $\Rightarrow X$ and Y are uncorrelated

However the converse is not true. Now let us ask the following question

Under what situation we can have

X and Y are independent \Leftrightarrow X and Y are uncorrelated
if and only if

The answer is as follows: If X and Y jointly follow the bi-variate bivariate normal distribution then the above assertion is true. So we shall start with the bi-variate bivariate normal density.

The bivariate normal distribution has a joint pdf given by

$$f_{X,Y}(x,y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left\{ \frac{(x-\mu_x)^2}{\sigma_x^2} - 2\rho \left(\frac{x-\mu_x}{\sigma_x} \right) \left(\frac{y-\mu_y}{\sigma_y} \right) + \frac{(y-\mu_y)^2}{\sigma_y^2} \right\}}$$

where $-\infty < x < +\infty$, $-\infty < y < +\infty$, and $\sigma_x > 0$, $\sigma_y > 0$. $\rho \in [-1, +1]$ are finite constants and further μ_x & μ_y are also finite constant real-valued constants.

We shall first show that it is a pdf.

- $f_{X,Y}(x,y) \geq 0$, $\forall (x,y) \in \mathbb{R}^2$ is clear.
- We will now have to show that

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x,y) dy dx = 1$$

To begin with let us substitute

$$u = \frac{x-\mu_x}{\sigma_x} \quad \text{and} \quad v = \frac{y-\mu_y}{\sigma_y}$$

Hence

$$\begin{aligned} & \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-\frac{1}{2} \frac{1}{1-\rho^2} (u^2 - 2\rho uv + v^2)} du dv \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} [(u-\rho v)^2 + (1-\rho^2)v^2]} du dv \end{aligned}$$

(check it yourself)

- Again let us substitute keeping v fixed.

$$w = \frac{u-\rho v}{\sqrt{1-\rho^2}}, \quad \therefore dw = \frac{du}{\sqrt{1-\rho^2}}$$

$$\therefore \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x,y) dy dx = \left(\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-w^2/2} dw \right) \left(\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-v^2/2} dv \right) = 1$$

(Think why??)

(Think about the standard normal distribution).

Thus $f_{X,Y}$ is a pdf.

Let us now compute the mgf of the bivariate normal distribution

$$m_{x,y}(t_1, t_2) = E[e^{t_1 x + t_2 y}] \\ = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{t_1 x + t_2 y} f_{x,y}(x, y) dy dx$$

Set $u = \frac{x - \mu_x}{\sigma_x}$ & $v = \frac{y - \mu_y}{\sigma_y}$, hence

$$m_{x,y}(t_1, t_2) = e^{t_1 \mu_x + t_2 \mu_y} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{t_1 \sigma_x u + t_2 \sigma_y v - \frac{1}{2} \frac{1}{1 - \rho_{xy}^2} (u^2 - 2\rho_{xy}uv + v^2)} dv du \\ \therefore m_{x,y}(t_1, t_2) = e^{t_1 \mu_x + t_2 \mu_y} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{1}{2\sqrt{1 - \rho_{xy}^2}} e^{-\frac{1}{2(1 - \rho_{xy}^2)} [u^2 - 2\rho_{xy}uv + v^2 - 2(1 - \rho_{xy}^2)t_1^2 - 2(1 - \rho_{xy}^2)t_2^2]} dv du$$

Let $w = \frac{u - \rho_{xy}v - (1 - \rho_{xy}^2)t_1 \sigma_x}{\sqrt{(1 - \rho_{xy}^2)}}$

$$z = v - \rho_{xy} t_1 \sigma_x - t_2 \sigma_y$$

$$\therefore m_{x,y}(t_1, t_2) = e^{t_1 \mu_x + t_2 \mu_y} e^{\left[\frac{1}{2} [t_1^2 \sigma_x^2 + 2\rho_{xy} t_1 t_2 \sigma_x \sigma_y + t_2^2 \sigma_y^2] \right]} \\ = e^{t_1 \mu_x + t_2 \mu_y} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{1}{2\pi} e^{-\frac{w^2}{2} - \frac{z^2}{2}} dw dz \\ = e^{t_1 \mu_x + t_2 \mu_y + \frac{1}{2} (t_1^2 \sigma_x^2 + 2\rho_{xy} t_1 t_2 \sigma_x \sigma_y + t_2^2 \sigma_y^2)}$$

$$\therefore m_{x,y}(t_1, t_2) = e^{t_1 \mu_x + t_2 \mu_y + \frac{1}{2} (t_1^2 \sigma_x^2 + 2\rho_{xy} t_1 t_2 \sigma_x \sigma_y + t_2^2 \sigma_y^2)}$$

[Note that $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{1}{2\pi} e^{-\frac{w^2}{2} - \frac{z^2}{2}} dw dz = 1$]

Think why

$$\text{Set } h = t_1 \mu_x + t_2 \mu_y + \frac{1}{2} (t_1^2 \sigma_x^2 + 2\rho_{xy} t_1 t_2 \sigma_x \sigma_y + t_2^2 \sigma_y^2)$$

(12)

$$\frac{\partial m_{x,y}(t_1, t_2)}{\partial t_1} = (\mu_x + t_1 \sigma_x^2 + \rho_{xy} t_2 \sigma_x \sigma_y) e^h$$

$$\frac{\partial m_{x,y}(t_1, t_2)}{\partial t_2} = (\mu_y + t_2 \sigma_y^2 + \rho_{xy} t_1 \sigma_x \sigma_y) e^h$$

$$\frac{\partial^2 m_{x,y}(t_1, t_2)}{\partial t_1 \partial t_2} = P_{xy} \sigma_x \sigma_y e^h + \left(\frac{(\mu_x + t_2 \sigma_y^2 + \rho_{xy} t_1 \sigma_x \sigma_y)^2}{(\mu_y + t_2 \sigma_y^2 + \rho_{xy} t_1 \sigma_x \sigma_y)} \right) \\ (\mu_x + t_1 \sigma_x^2 + \rho_{xy} t_2 \sigma_x \sigma_y)$$

$$\therefore \frac{\partial m_{x,y}(t_1, t_2)}{\partial t_1} \Big|_{(t_1, t_2) = (0,0)} = \mu_x = E(x)$$

$$\therefore \frac{\partial m_{x,y}(t_1, t_2)}{\partial t_2} \Big|_{(t_1, t_2) = (0,0)} = \mu_y = E(y)$$

$$\frac{\partial^2 m_{x,y}(t_1, t_2)}{\partial t_1 \partial t_2} \Big|_{(t_1, t_2) = (0,0)} = \cancel{P_{xy} \sigma_x \sigma_y} + \mu_x \mu_y$$

$$\therefore \text{Cov}(x, y) = \sigma_{xy} = E[xy] - E[x]E[y] \\ = P_{xy} \sigma_x \sigma_y + \mu_x \mu_y - \mu_x \mu_y$$

$$\therefore \text{Cov}(x, y) = P_{xy} \sigma_x \sigma_y$$

\therefore Correlation coefficient = P_{xy} .

We leave it to the reader to check that $\text{Var}(x) = \sigma_x^2 + \text{Var}(y) = \sigma_y^2$

Now we come to the main result of this section.

Theorem 8.5: If X and Y are two random variables which jointly follow the bivariate normal distribution, Then X and Y are independent $\Leftrightarrow X$ and Y are uncorrelated.

Proof: If X and Y are independent, then X and Y are uncorrelated is a known fact.

For the converse, let X and Y be uncorrelated

$$\Rightarrow \text{Cov}(X, Y) = 0$$

$$\Rightarrow \rho_{XY} = 0 = -\frac{1}{2} \left[\left(\frac{x-\mu_X}{\sigma_X} \right)^2 + \left(\frac{y-\mu_Y}{\sigma_Y} \right)^2 \right]$$

Thus

$$\begin{aligned} f_{X,Y}(x, y) &= \frac{1}{2\pi\sigma_X\sigma_Y} e^{-\frac{1}{2} \left[\left(\frac{x-\mu_X}{\sigma_X} \right)^2 + \left(\frac{y-\mu_Y}{\sigma_Y} \right)^2 \right]} \\ &= \left[\frac{1}{\sqrt{2\pi}\sigma_X} e^{-\frac{1}{2} \left(\frac{x-\mu_X}{\sigma_X} \right)^2} \right] \left[\frac{1}{\sqrt{2\pi}\sigma_Y} e^{-\frac{1}{2} \left(\frac{y-\mu_Y}{\sigma_Y} \right)^2} \right] \\ &= f_X(x) f_Y(y), \end{aligned}$$

$$\text{where } X \sim N(\mu_X, \sigma_X^2) \quad \text{and} \quad Y \sim N(\mu_Y, \sigma_Y^2).$$

Of course one can easily check that when $\rho_{XY} = 0$, the marginal distributions are ~~$f_{X,Y}$~~ $f_X(x)$ and $f_Y(y)$. This shows that X and Y are independent random variable \square .

— x —

Section 8.3: Conditional Expectation

Conditional Expectation is a natural outcome from the fact that we have defined the notion of a conditional density. Let us define it formally.

Conditional Expectation of the r.v.'s X and Y

Let (X, Y) be a two-dimensional random vector. Let $Z = g(X, Y)$ be a function of the two random variables. Then the conditional expectation of $g(X, Y)$, given $X = x$ is given as

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

This formula holds good if X and Y are jointly continuous. If X and Y are jointly discrete we have

$$E[g(X, Y) | X = x] = \sum_j g(x, y_j) f_{Y|X}(y_j|x).$$

E.g: Let $f_{X,Y}$ be the joint density of two continuous random variable given as

$$f_{X,Y}(x, y) = \begin{cases} x+y; & \text{if } x \geq 0 \text{ and } y \geq 0 \text{ and } x \in (0,1) \\ & \quad \text{and } y \in (0,1) \\ 0 & \text{otherwise.} \end{cases}$$

How shall we now proceed to compute $E[Y | X = x]$ (i.e. $g(x, Y) = Y$)

First we compute

$$f_{Y|X}\left(\frac{y}{x}\right) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{x+y}{f_X(x)}, \quad x, y \in (0,1).$$

$$f_X(x) = \int_0^1 (x+y) dy = x + \frac{1}{2} \quad \therefore f_{Y|X} = \boxed{\frac{x+y}{x+\frac{1}{2}} = f_{Y|X}(y|x)}$$

(15)

Here $g(x, Y) = Y$
In fact $E[g(X, Y) | X = x]$
is a function of X

Thus

$$\begin{aligned} E[Y|X=x] &= \int_0^1 y f_{Y|X}(y|x) dy \\ &= \int_0^1 y \frac{x+y}{x+1/2} dy \\ &= \frac{1}{x+1/2} \left[\int_0^1 yx dy + \int_0^1 y^2 dy \right] \\ &= \frac{1}{x+1/2} \left[\frac{x}{2} + \frac{1}{3} \right]. \quad \square \end{aligned}$$

For simplicity, consider $E[g(Y)|x]$, which is in general a function of x . Let us

denote it as

$$\phi(x) = E[g(Y)|x]$$

$\therefore \phi(x)$ can also be written as
 $\phi(x) = E[g(Y)|X]$

\therefore For a given $w \in \Omega$

$$\phi(X(w)) = E[g(Y)|X(w)]$$

\therefore If $X(w) = x$, then we have

$$\phi(x) = E[g(Y)|x]$$

Note that in general $\phi(x) = E[g(Y)|X]$ can be viewed as a random variable.

So our learning here is as follows:

Conditional Expectation $E[g(Y)|X]$ is a random variable.

It might be sometimes relevant to ask, "What is the expectation of $\phi(x)$ ". The answer may surprise you. We do the calculation for $g(Y) = Y$. Assume $f_X(x) > 0, \forall x$

Set

$$\begin{aligned}
 E[\phi(x)] &= E[E[Y|x]] \\
 &= E[\phi(x)] = \int_{-\infty}^{\infty} \phi(x) f_X(x) dx \\
 &= \int_{-\infty}^{\infty} E[Y|x] f_X(x) dx \\
 &= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} y f_{Y|x}(y|x) dy \right] f_X(x) dx \\
 &= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} y f_{Y|x}(y|x) dy \right] f_X(x) dx \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y f_{Y|x}(y|x) f_X(x) dy dx \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y \cdot \frac{f_{X,Y}(x,y)}{f_X(x)} \cdot f_X(x) dy dx \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y f_{X,Y}(x,y) dy dx \\
 &= E[Y].
 \end{aligned}$$

$$\therefore \boxed{E[E[Y|x]] = E[Y]}$$

(The reader can prove it for the discrete case).

We shall now define the notion of conditional variance of Y when the random variable X takes a particular value say x .

$$\text{Var}(Y|X=x) = E[Y^2|X=x] - [E[Y|X=x]]^2.$$

In way we discussed before the case of the conditional expectation, the conditional variance can also be viewed as a r.v. $\text{Var}(Y|X)$, whose value at any $w \in \Omega$, is computed as

$$\boxed{\text{Var}(Y|X)(w) = \text{Var}(Y|X=w)}$$

$X(w)$ can of course take the value of x .

We end our discussion with the following interesting result.

Theorem 8.4 Let X and Y be random variables. Then

$$\text{Var}(Y) = E[\text{Var}(Y|X)] + \text{Var}[E[Y|X]].$$

$$\begin{aligned} \text{Proof: } E[\text{Var}(Y|X)] &= E[E[Y^2|X] - (E[Y|X])^2] \\ &= E[E[Y^2|X]] - E[(E[Y|X])^2] \\ &= E[Y^2] - E[E[Y|X]^2] \end{aligned}$$

\therefore In fact $E[E[Y^2|X]] = E[Y^2]$, can be proved from the proof style of showing $E[E[Y|X]] = E[Y]$

$$\begin{aligned} \therefore E[\text{Var}(Y|X)] &= E[Y^2] - (E[Y])^2 - E[E[Y|X]] \\ &\quad - E[(E[Y|X])^2] + (E[Y])^2 \end{aligned}$$

$$\begin{aligned} \therefore E[\text{Var}(Y|X)] &= E[\text{Var}(Y)] - [E[(E[Y|X])^2] \\ &\quad - (E(E[Y|X]))^2] \end{aligned}$$

$$= \text{Var}(Y) - \text{Var}(E[Y|X])$$

$$\therefore \text{Var}(Y) = E[\text{Var}(Y|X)] + \text{Var}(E[Y|X]).$$

Lecture 9: Transformation of Random Variables & their Distributions

Suppose we have a random variable X and another random Y is defined as $Y = g(x)$. If we assume just for the sake of it that X is a continuous random variable with a p.d.f. Sometimes specially in statistics, we need to know the density of Y .

We can also have more than one such random variables, X_1, X_2, \dots, X_n and suppose we know their joint pmf/pdf. Then we might want to know the pmf/pdf of

$$Y = g(x_1, x_2, \dots, x_n)$$

In this chapter we shall show how can we compute $F_Y(y)$. There are several approaches and as we will see what kind of technique to use depends on the problem at hand.

Section 1: Distribution function Approach

One of the simplest approaches to is to compute the distribution function of Y and then differentiate it to get the p.d.f

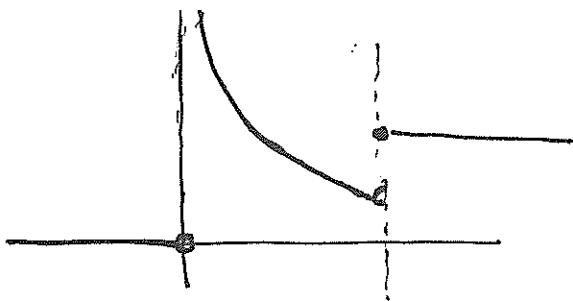
$$\text{We have } F_Y(y) = P(Y \leq y) = P(g(x_1, \dots, x_n) \leq y)$$

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

E.g. 1. Let $X \sim \text{Uniform}(0, 1)$, then if $Y = g(x) = x^2$. We shall find the distribution of Y .

$$\begin{aligned} F_Y(y) &= P(Y \leq y) = P(x^2 \leq y) = \int_{\{x: x^2 \leq y\}} f_X(x) dx \\ &= \int_0^{\sqrt{y}} dx \quad (0 < y < 1) \\ &= \sqrt{y}, \quad \text{when } 0 < y < 1. \end{aligned}$$

$$\therefore F_Y(y) = \begin{cases} 0, & y \leq 0 \\ \sqrt{y}, & 0 < y < 1 \\ 1, & y \geq 1 \end{cases}$$



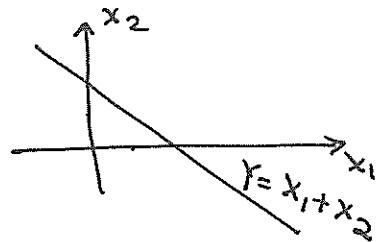
$$\therefore f_Y(y) = \begin{cases} \frac{1}{2\sqrt{y}}, & 0 < y < 1 \\ 0, & \text{otherwise} \end{cases}$$

E.g. 2: If X_1 is a continuous r.v. & X_2 is a continuous r.v. whose joint density is given as

From
Miller & Miller
Freund's
Mathematical

$$f_{X_1, X_2}(x_1, x_2) = \begin{cases} 6e^{-3x_1 - 2x_2} & ; x_1 > 0, x_2 > 0 \\ 0 & \text{elsewhere} \end{cases}$$

Find the distribution of $Y = X_1 + X_2$ and also its probability density



$$\begin{aligned} F_Y(y) &= P(Y \leq y) \\ &= P(X_1 + X_2 \leq y) \\ &= P(x_1 \leq y - x_2, x_2 \leq y) \\ &= \int_0^y \int_0^{y-x_2} 6e^{-3x_1 - 2x_2} dx_1 dx_2 \\ &= 1 + 2e^{-3y} - 3e^{-2y}. \quad (\text{when } y > 0) \end{aligned}$$

Thus $f_Y(y) = 6(e^{-2y} - e^{-3y})$, for $y > 0$

& $f_Y(y) = 0$, elsewhere.

Section 2: Transformation Approach : Single r.v.

So in this section we consider the case $Y = u(X)$, and given a probability mass/density of X , we would like to know what is the probability distribution of Y .

Consider the following simple example from Miller and Miller []

Eg 1: X be the r.v. denoting the number of heads in a toss of four coins a fair coin four times.

Find the distribution of $Y = \frac{1}{1+x}$

[E.g. from
Miller and
Miller]
Freund's Mathematical
Statistics

The distribution of X is given as

x	$f_X(x)$
0	$\frac{1}{16}$
1	$\frac{4}{16}$
2	$\frac{6}{16}$
3	$\frac{4}{16}$
4	$\frac{1}{16}$

So Y takes the values, $1, -\frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \frac{1}{5}$
corresponding to $x=0, x=1, x=2, x=3, x=4$, respectively

y	Y	$g_Y(y)$
1	$\frac{1}{16}$	$\frac{1}{16}$
$-\frac{1}{2}$	$\frac{4}{16}$	$\frac{4}{16}$
$\frac{1}{3}$	$\frac{6}{16}$	$\frac{6}{16}$
$\frac{1}{4}$	$\frac{4}{16}$	$\frac{4}{16}$
$\frac{1}{5}$	$\frac{1}{16}$	$\frac{1}{16}$

So we have $g_Y(y) = f_X(\frac{1}{y} - 1)$. But observe that

$$x = \frac{1}{y} - 1.$$

For the continuous case however we need to rely on the following result.

Theorem 9.1: Let X be a random variable, which is continuous one with p.d.f $f_X(x)$. Let $y = u(x)$ be a function of x , which is differentiable and monotone (either increasing or decreasing) over $\text{Range}(x)$ the set $\{x \in \text{Range } x : f_X(x) \neq 0\}$. Then for these values of x , the function u has an inverse w , for which $x = w(y)$. Then the probability density is given of $Y = u(X)$ is given as

$$g_Y(y) = f_X[w(y)] |w'(y)|, \text{ provided } w'(y) \neq 0$$

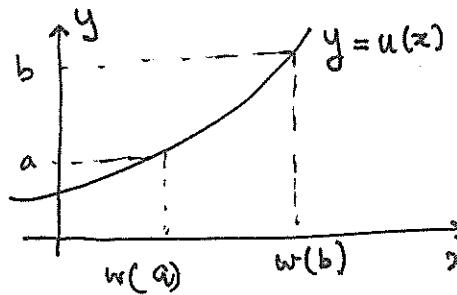
and $g_Y(y) = 0$, elsewhere. We assume u and w are differentiable

and $g_Y(y) = 0$, elsewhere. We assume u and w are differentiable and $g_Y(y) = 0$, elsewhere. We assume u and w are differentiable

and $g_Y(y) = 0$, elsewhere. We assume u and w are differentiable

and $g_Y(y) = 0$, elsewhere. We assume u and w are differentiable

and $g_Y(y) = 0$, elsewhere. We assume u and w are differentiable



$$\begin{aligned} \text{Now } P[a < Y < b] &= P[w(a) < x < w(b)] \\ &= \int_{w(a)}^{w(b)} f_X(x) dx \end{aligned}$$

Now substitute in the integral, $x = w(y)$. \therefore when $x = w(a)$, we have $y = a$, & when $x = w(b)$, we have $y = b$.

\therefore When $x = w(a)$, we have $y = a$, & when $x = w(b)$, we have $y = b$.

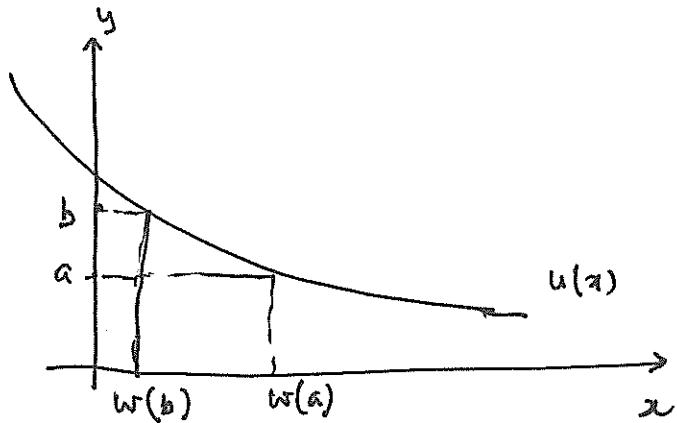
$$y = b \text{ and } dx = w'(y) dy.$$

$$\therefore P[a < x < b] = \int_a^b f_X(w(y)) w'(y) dy$$

Thus the probability density is given by the function

$$g_Y(y) = f_X(w(y)) w'(y).$$

Let us consider now u to be decreasing



In this case we have

$$\begin{aligned} P[a < Y < b] &= P[w(b) < X < w(a)] \\ &= \int_{w(b)}^{w(a)} f_X(x) dx \\ &= \int_{\infty}^a f_X(w(y)) w'(y) dy \\ &= - \int_a^b f_X(w(y)) w'(y) dy \end{aligned}$$

In this case we have

$$g_Y(y) = - f_X(w(y)) w'(y) dy$$

$$\text{In fact. } w'(y) = \frac{dx}{dy} = \frac{1}{\frac{dy}{dx}} \text{ i.e. } w'(y) = \frac{1}{u'(x)}$$

In fact when u is increasing $w'(y)$ is positive and when u is decreasing $-w'(y)$ is positive. Hence compactly we write

$$g_Y(y) = f_X[w(y)] |w'(y)|$$

$$\text{or } g_Y(y) = f_X(u^{-1}(y)) \left| \frac{dx}{dy} \right| ; \left(\text{Here } w \equiv u^{-1} \right)$$

$$g_Y(y) = f_X(u^{-1}(y)) \left| \frac{dx}{dy} \right|$$

$u'(x) \neq 0$
 $\Leftrightarrow x = u^{-1}(y)$.

E.g #2: Let $X \sim N(0, 1)$. Then find the density of

$$Z = u(X) = X^2.$$

Note that the function $Z = x^2$, is decreasing when $x < 0$ and increasing when $x > 0$. So conditions of the Theorem a.1 are not met.

The key to this is to take an additional step and have a random variable on whose range the given transformation This brings us to the case what happens if $u(x)$ behaves differently on different parts of $\{x \in \text{Ran}(x) : u(x) \neq 0\}$.

Assume that we can partition the set $\{x \in \text{Ran}(x) : u(x) \neq 0\}$ into a finite partition $\{A_n\}_{n=1}^k$, i.e.

$$\bigcup_{n=1}^k A_n = \{x \in \text{Ran}(x) : u(x) \neq 0\}$$

$$\Leftrightarrow A_i \cap A_j = \emptyset, \forall i \neq j, i, j = 1, \dots, k.$$

So each A_n , the function u is either increasing or decreasing.

Thus define

$$u_n(x) = \begin{cases} u(x), & x \in A_n \\ 0, & x \text{ is otherwise.} \end{cases}$$

$\therefore u_n$ has a unique inverse in A_n . Thus

$$y = u_n(x) \Leftrightarrow x = g_{u_n^{-1}}(y), \text{ for } x \in A_n$$

Then we can get the pdf on the individual A_n and then sum

i.e.

$$g_Y(y) = \sum_{n=1}^k f_X(g_{u_n^{-1}}(y)) \left| \frac{d}{dy} u_n^{-1}(y) \right|$$

So in our particular case we have

$$\{x \in \mathbb{R} : x^2 \neq 0\} = \{x \in \mathbb{R}^*, x < 0\} \cup \{x \in \mathbb{R} : x^* > 0\}$$

\downarrow
 \downarrow

A_1
A₂

So on $(-\infty, 0)$ we have $x_1 = -\sqrt{y}$ & on $(0, \infty)$ we have $x_2 = +\sqrt{y}$. Here $y > 0$. as $y = x^2 > 0$.

$$\begin{aligned} g_Y(y) &= f_X(-\sqrt{y}) \frac{1}{2\sqrt{y}} + f_X(\sqrt{y}) \frac{1}{2\sqrt{y}} \\ &= \frac{1}{\sqrt{2\pi}} e^{-y/2} \frac{1}{2\sqrt{y}} + \frac{1}{2\sqrt{y}} \cdot \frac{1}{\sqrt{2\pi}} e^{-y/2} \\ &= \frac{1}{2} \left(\frac{1}{\sqrt{2\pi y}} e^{-y/2} \right) + \frac{1}{2} \left(\frac{1}{\sqrt{2\pi y}} e^{-y/2} \right) \\ &\therefore g_Y(y) = \boxed{\frac{1}{\sqrt{2\pi y}} e^{-y/2}, \quad y > 0} \end{aligned}$$

Of course $g_Y(y) = 0$, if $y \leq 0$.

In fact g_Y is often said to be the pdf of a χ^2 -distri random variable with degrees of freedom 1. ↴ (more on this later).

Section 3: Distribution of the sum and difference of two r.v.

In this section we will be concerned with the distribution of the sum and differences of two random variables. These notions would play an important role in sampling theory. So we begin by presenting the analysis in forms of theorems, so that the reader appreciates the importance of these results. We present these results from "Introduction to the Theory of Statistics" by Mood, Graybill, Boyce.

Obsv:
 $Y \sim \text{Gamma}(\frac{1}{2}, \frac{1}{2})$

χ^2 is call
 χ^2 -with
one d
of free

Theorem 9.1: Let X and Y be two continuous random variables, with joint pdf $f_{X,Y}(x,y)$. Let $Z = X+Y$ & $V = X-Y$, then we have

$$\text{A) } f_Z(z) = \int_{-\infty}^{\infty} f_{X,Y}(x, z-x) dx = \int_{-\infty}^{\infty} f_{X,Y}(z-y, y) dy$$

$$\text{B) } f_V(v) = \int_{-\infty}^{\infty} f_{X,Y}(x, x-v) dx = \int_{-\infty}^{\infty} f_{X,Y}(v+y, y) dy$$

We will just prove the first part of case A). The rest can be proved in an analogous way.

Proof:

$$\begin{aligned} F_Z(z) &= P(Z \leq z) \\ &= P(X+Y \leq z) \\ &= \iint_{\substack{x+y \leq z}} f_{X,Y}(x,y) dx dy \\ &= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{z-x} f_{X,Y}(x,y) dy \right] dx \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^z \left[f_X(x, y) du \right] dx \quad [\text{By substituting } y=u-x] \\ &= \int_{-\infty}^z \left[\int_{-\infty}^{\infty} f_{X,Y}(x, u-x) dx \right] du \\ \therefore f_Z(z) &= \frac{dF_Z(z)}{dz} = \frac{d}{dz} \int_{-\infty}^z \left[\int_{-\infty}^{\infty} f_{X,Y}(x, u-x) dx \right] du \\ &= \int_{-\infty}^{\infty} f_{X,Y}(x, z-x) dx \quad [\text{Use the fundamental theorem of Calculus but for the improper integrals.}] \end{aligned}$$

Corollary 9.1: Let X and Y are now assumed to be independent, and $Z = X+Y$.

Then,

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(z-x) dx = \int_{-\infty}^{\infty} f_X(z-y) f_Y(y) dy$$

Section 4: The distribution of products and ratio

Theorem 9.2: Let X and Y are continuous random variables with joint p.d.f $f_{X,Y}(x,y)$ and let $Z = XY$ and $V = \frac{X}{Y}$.

Then;

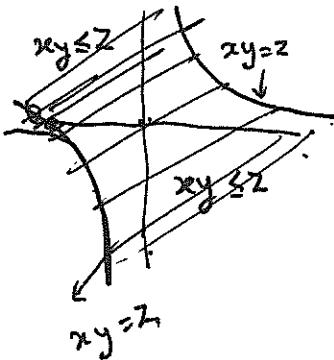
$$\begin{aligned} a) f_Z(z) &= \int_{-\infty}^{\infty} \frac{1}{|x|} f_{X,Y}\left(x, \frac{z}{x}\right) dx \\ &= \int_{-\infty}^{\infty} \frac{1}{|y|} f_{X,Y}\left(\frac{z}{y}, y\right) dy \end{aligned}$$

$$b) f_V(u) = \int_{-\infty}^{\infty} |y| f_{X,Y}(uy, y) dy$$

We shall only prove the first part of a), the rest can be similarly worked out by the reader.

$$F_Z(z) = P(Z \leq z) = \iint_{xy \leq z} f_{X,Y}(x, y) dx dy$$

$$= \int_{-\infty}^0 \left[\int_{z/x}^{\infty} f_{X,Y}(x, y) dy \right] dx + \int_0^{\infty} \left[\int_{-\infty}^{z/x} f_{X,Y}(x, y) dy \right] dx$$



Let us now substituting $u = xy$, we have

$$\begin{aligned}
 F_Z(z) &= \int_{-\infty}^0 \left[\int_z^{-\infty} f_{X,Y}(x, \frac{u}{x}) \frac{dx}{x} du \right] dx \\
 &\quad + \int_0^\infty \left[\int_{-\infty}^z f_{X,Y}(x, \frac{u}{x}) \frac{du}{x} \right] dx \\
 &= \int_{-\infty}^z \left[\int_{-\infty}^0 \frac{1}{|x|} f_{X,Y}(x, \frac{u}{x}) dx \right] du \\
 &\quad + \int_{-\infty}^z \left[\int_0^\infty \frac{1}{|x|} f_{X,Y}(x, \frac{u}{x}) dx \right] du \\
 &= \int_{-\infty}^z \left[\int_{-\infty}^{\infty} \frac{1}{|x|} f_{X,Y}(x, \frac{u}{x}) dx \right] du.
 \end{aligned}$$

Hence

$$\begin{aligned}
 f_Z(z) &= \frac{dF_Z(z)}{dz} \\
 &= \int_{-\infty}^{\infty} \frac{1}{|x|} f_{X,Y}(x, \frac{u}{x}) dx.
 \end{aligned}$$

Example: Let X and Y are independent uniform random variables in $(0,1)$. Compute the distribution of $Z = XY$.

We have $f_Z(z) = \int_{-\infty}^{\infty} \frac{1}{|x|} f_{X,Y}(x, \frac{z}{x}) dx$

$$\begin{aligned}
 &= \int_{-\infty}^{\infty} \frac{1}{|x|} f_X(x) f_Y(\frac{z}{x}) dx, \quad \text{by independence of } X \text{ and } Y. \\
 &\approx \int_0^1 \frac{1}{x} \cdot \text{Now } f_Y(\frac{z}{x}) = 1, \text{ if } 0 < \frac{z}{x} < 1 \\
 &\quad > f_X(x) = 1, \text{ if } 0 < x < 1
 \end{aligned}$$

$$\therefore f_Y\left(\frac{z}{x}\right) = I_{(0,1)}\left(\frac{z}{x}\right)$$

$$f_X(x) = I_{(0,1)}(x)$$

$$I_A(x) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{if } x \notin A \end{cases}$$

Indicator function or
characteristic function

$$\text{Now } I_{(0,1)}(x) I_{(0,1)}\left(\frac{z}{x}\right)$$

$$= I_{(0,1)}(x) I_{(z,1)}(x) \quad (\text{Try to establish this})$$

$$\therefore f_Z(z) = \int_{-\infty}^{\infty} \frac{1}{|x|} I_{(0,1)}(x) I_{(0,1)}\left(\frac{z}{x}\right) dx$$

$$= \int_{-\infty}^{\infty} \frac{1}{|x|} I_{(0,1)}(x) I_{(z,1)}(x) dx$$

$$= I_{(0,1)}(x) \int_{-\infty}^{\infty}$$

$$= I_{(0,1)}(x) \int_{-\infty}^{\infty} \frac{1}{|x|} I_{(z,1)}(x) dx$$

$$= I_{(0,1)}(x) \int_z^1 \frac{1}{x} dx$$

$$\boxed{f_Z(z) = -\log z I_{(0,1)}(z)}$$

$$\text{i.e. } \boxed{f_Z(z) = -\log z, \quad z \in (0,1)}$$

$$f_Z(z) = 0, \quad \text{otherwise.}$$

Section 5: Further Transformation Techniques

Suppose now we have the following transformation.

$$Y = u(x_1, x_2).$$

How are we going to get the pmf or pdf in this case.

Our key idea is the following. Suppose we keep x_2 fixed and vary x_1 , then y becomes a function of x_1 . Suppose under this fixed x_2 , the function $u(x_1, x_2)$ is increasing or decreasing and thus has an inverse, i.e. x_1 can be written as a function of y , i.e. $x_1 = w_1(y, x_2)$. So

using Theorem 9.1 we write,

$$g_Y(y, x_2) = f_{x_1, x_2}(x_1, x_2) \left| \frac{\partial x_1}{\partial y} \right|$$

or we may have that $u(x_1, x_2)$ satisfies the requirements of Theorem 9.1 if x_1 is kept fixed and thus we will have

$$g_Y(x_1, y) = f_{x_1, x_2}(x_1, x_2) \left| \frac{\partial x_2}{\partial y} \right|.$$

Whatever be the case, we can from $g_Y(x_1, y)$ or $g_Y(y, x_2)$ can be integrated out with respect to x_1 or x_2 , to get the marginal distribution of Y .

Let us see an example

Let x_1 and x_2 be two random variables which are independent and each follow uniform distribution in $(0, 1)$.

Let us fix set $Y = x_1 + x_2$. Compute the pdf of Y . Then

Y varies from 0 to 2 as x_1 and x_2 vary from 0 to 1.

If x_2 is kept fixed we have

$$\therefore x_1 = Y - x_2.$$

$$\begin{aligned} \therefore g_Y(y, x_2) &= I_{(0,1)}(y-x_2) I_{(0,1)}(x_2) |1| \\ &= I_{(0,1)}(y-x_2) I_{(0,1)}(x_2). \end{aligned}$$

$$\begin{aligned} g_Y(y) &= \int_{-\infty}^{\infty} g_Y(y, x_2) dx_2 \\ &= \int_{-\infty}^{\infty} I_{(0,1)}(y-x_2) I_{(0,1)}(x_2) dx_2 \end{aligned}$$

When $y - x_2 \in (0, 1)$, then $0 < y - x_2 < 1$

$\therefore x_2 < y \quad \& \quad y < x_2 + 1$. Now if we set A to
~~set $y \in (0, 1)$, then if $y > 1$, & $y < 2$, then $y - x_2 < 1$~~
 $\Rightarrow y-1 < x_2 < 1$

Now observe that

$$\begin{aligned} I_{(0,1)}(y-x_2) I_{(0,1)}(x_2) \\ = I_{(0,y)}(x_2) I_{(0,1)}(y) + I_{(y-1,1)}(x_2) I_{[1,2]}(y) \end{aligned}$$

$$\begin{aligned} \therefore g_Y(y) &= \int_{-\infty}^{\infty} \left[I_{(0,y)}(x_2) I_{(0,1)}(y) + I_{(y-1,1)}(x_2) I_{[1,2]}(y) \right] dx_2 \\ &= I_{(0,1)}(y) \int_0^y I_{(0,y)}(x_2) dx_2 + I_{[1,2]}(y) \int_{y-1}^1 I_{(y-1,1)}(x_2) dx_2 \\ &= y I_{(0,1)}(y) + I_{[1,2]}(y)(2-y) \\ \therefore g_Y(y) &= \begin{cases} y & \text{if } 0 < y < 1 \\ 2-y & \text{if } 1 \leq y < 2 \end{cases} \end{aligned}$$

& $g_Y(y) = 0$, elsewhere.

Now what happens if we have

$$Y_1 = u_1(x_1, x_2)$$

$$Y_2 = u_2(x_1, x_2)$$

Then if we know the joint distribution of x_1, x_2 can we find the joint distribution of Y_1, Y_2 . The following theorem tells us how to do it.

Theorem 9.2 :

Let x_1 and x_2 are continuous random variables, with joint pdf $f_{x_1, x_2}(x_1, x_2)$. If the transformations $y_1 = u_1(x_1, x_2)$ and $y_2 = u_2(x_1, x_2)$ are bijective mappings over the set $\{(x_1, x_2) \in \text{Range}(x_1) \times \text{Range}(x_2) : f_{x_1, x_2}(x_1, x_2) \neq 0\}$.

Then over the above set we can uniquely write

$$x_1 = w_1(y_1, y_2)$$

$$x_2 = w_2(y_1, y_2)$$

Then the joint distribution $Y_1 = u_1(x_1, x_2)$ & $Y_2 = u_2(x_1, x_2)$

is given as

$$g_{Y_1, Y_2}(y_1, y_2) = f(w_1(y_1, y_2), w_2(y_1, y_2)) |\det J|$$

where J is the Jacobian matrix of the map

$$\phi(x_1, x_2) = \begin{pmatrix} w_1(y_1, y_2) \\ w_2(y_1, y_2) \end{pmatrix}$$

$$\circ \phi(y_1, y_2) = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} w_1(y_1, y_2) \\ w_2(y_1, y_2) \end{pmatrix}$$

given as

$$J = \begin{bmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} \end{bmatrix}.$$

Jacobian matrix
is the derivative
of the map
 $\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^2$

Let us provide some interesting examples.

Eg (A): Let x_1 & x_2 are independent standard normal variables.

$$\text{Let } Y_1 = x_1 + x_2 \quad \& \quad Y_2 = \frac{x_1}{x_2}$$

Find $g_{Y_1, Y_2}(y_1, y_2)$.

$$\text{Here } y_1 = u_1(x_1, x_2) = x_1 + x_2.$$

$$\& \quad y_2 = u_2(x_1, x_2) = x_1/x_2.$$

Then,

$$x_1 = w_1(y_1, y_2) = \frac{y_1 y_2}{1 + y_2}$$

$$x_2 = w_2(y_1, y_2) = \frac{y_1}{1 + y_2}$$

$$\text{do} \quad J = \begin{bmatrix} \frac{y_2}{1+y_2}, & \frac{y_1}{(1+y_2)^2} \\ \frac{1}{1+y_2}, & -\frac{y_1}{(1+y_2)^2} \end{bmatrix}$$

$$\therefore \det J = -\frac{y_1}{(1+y_2)^2} \quad \therefore |\det J| = \frac{|y_1|}{(1+y_2)^2}$$

$$\therefore g_{Y_1, Y_2}(y_1, y_2) = \frac{1}{2\pi} \left[e^{-\frac{1}{2} \left[\frac{(y_1 y_2)^2}{(1+y_2)^2} + \frac{y_1^2}{(1+y_2)^2} \right]} \right] \frac{|y_1|}{(1+y_2)^2}$$

$$\therefore g_{Y_1, Y_2}(y_1, y_2) = \frac{1}{2\pi} \frac{|y_1|}{(1+y_2)^2} \left[e^{-\frac{1}{2} \left[\frac{(1+y_2^2) y_1^2}{[1+y_2]^2} \right]} \right].$$

Question: Can you find the marginal distribution of Y_2 .

E.g (B): Let $x_1 \sim \text{Gamma}(n_1, \lambda_1)$

$x_2 \sim \text{Gamma}(n_2, \lambda_2)$

and x_1 and x_2 are independent. Then

find the joint distribution of $y_1 = x_1 + x_2$ with

$$y_2 = \frac{x_1}{x_2} \quad \therefore$$

Find the marginal distribution of y_1 .

Our first approach would be to compute, first the joint distribution of $g_{Y_1, Y_2}(y_1, y_2)$ and then compute out the marginal;

Here

$$x_1 = w_1(y_1, y_2) = \frac{y_1 y_2}{1+y_2}$$

$$x_2 = w_2(y_1, y_2) = \frac{y_1}{1+y_2}$$

$$\therefore |\det J| = \frac{y_1}{(1+y_2)^2} \quad (\text{note } y_1 \text{ and } y_2 \text{ are positive as } x_1 \in (0, \infty) \text{ and } x_2 \in (0, \infty))$$

For $y_1 \in (0, \infty)$ and $y_2 \in (0, \infty)$ we have

$$\therefore g_{Y_1, Y_2}(y_1, y_2) = \frac{y_1}{(1+y_2)^2} \frac{1}{\Gamma(n_1)} \frac{1}{\Gamma(n_2)} \lambda^{n_1+n_2} \left(\frac{y_1 y_2}{1+y_2}\right)^{n_1-1} \left(\frac{y_1}{1+y_2}\right)^{n_2-1}$$

$$\begin{aligned} \therefore g_{Y_1, Y_2}(y_1, y_2) &= \frac{y_1}{(1+y_2)^2 \Gamma(n_1) \Gamma(n_2)} \lambda^{n_1+n_2} \left(\frac{y_1 y_2}{1+y_2}\right)^{n_1-1} \left(\frac{y_1}{1+y_2}\right)^{n_2-1} \\ &= \frac{\lambda^{n_1+n_2}}{\Gamma(n_1) \Gamma(n_2)} y_1^{n_1+n_2-1} e^{-\lambda y_1} \frac{y_2^{n_1-1}}{(1+y_2)^{n_1+n_2}} \end{aligned}$$

$$\text{Now } B(n_1, n_2) = \frac{\Gamma(n_1) \Gamma(n_2)}{\Gamma(n_1+n_2)} \Rightarrow B(n_1, n_2) \Gamma(n_1+n_2) = \Gamma(n_1) \Gamma(n_2)$$

$$\therefore \Rightarrow g_{Y_1, Y_2}(y_1, y_2) = \underbrace{\left(\frac{\lambda^{n_1+n_2}}{\Gamma(n_1+n_2)} \cdot y_1^{n_1+n_2-1} e^{-\lambda y_1} \right)}_{g_{Y_1}(y_1)} \times \underbrace{\left(\frac{1}{B(n_1, n_2)} \frac{y_2^{n_1-1}}{(1+y_2)^{n_1+n_2}} \right)}_{g_{Y_2}(y_2)}$$

(16)

Thus Y_1 and Y_2 are independent random variables, with

$$Y_1 \sim \text{Gamma}(n_1+n_2, \lambda)$$

If $n_1 = n_2 = 1$, Compute $g_{Y_2}(y_2)$.

Section 6: The mgf technique.

When we have $Y = u(x_1, x_2)$, then, sometimes it is useful to compute the mgf of Y and see if it matches with mgf of any distribution known to us. This is truly effective when x_1, x_2 are independent. In fact this works pretty well even if we have a larger number of independent variables.

For example if we have x_1, x_2, \dots, x_n as independent random variables, then, if

$$Y = x_1 + \dots + x_n$$

$$\text{Then } m_Y(t) = \prod_{i=1}^n m_{x_i}(t). \quad (\text{Prove it in the assignment})$$

Consider again for example that x_1, x_2, \dots, x_n are independent random variables with

$$x_i \sim \text{Poisson}(\lambda_i)$$
$$\therefore m_{x_i}(t) = e^{\lambda_i(e^t - 1)}$$

$$\text{Let } Y = x_1 + \dots + x_n.$$

$$\begin{aligned}\therefore m_Y(t) &= \prod_{i=1}^n m_{x_i}(t) = \prod_{i=1}^n e^{\lambda_i(e^t - 1)} \\ &= e^{(\lambda_1 + \dots + \lambda_n)(e^t - 1)}\end{aligned}$$

$$\text{Then } Y \sim \text{Poisson}(\lambda_1 + \dots + \lambda_n).$$

Now consider another example.

Let x_1, x_2 are standard normal variables.

Find the joint distribution of Y_1 and Y_2

$$\text{where } Y_1 = x_1 + x_2 \quad \& \quad Y_2 = x_2 - x_1$$

$$\begin{aligned} m_{Y_1 Y_2}(t_1, t_2) &= E[e^{t_1 Y_1 + t_2 Y_2}] \\ &= E[e^{t_1(x_1+x_2) + t_2(x_2-x_1)}] \\ &= E[e^{(t_1-t_2)x_1 + (t_1+t_2)x_2}] \\ &= E[e^{(t_1-t_2)x_1}] E[e^{(t_1+t_2)x_2}] \\ &= E[e^{(t_1-t_2)x_1}] E[e^{(t_1+t_2)x_2}] \end{aligned}$$

$(\because x_1, x_2$
are indept.)

$$\begin{aligned} m_{Y_1 Y_2}(t_1, t_2) &= m_{X_1}(t_1-t_2) m_{X_2}(t_1+t_2) \\ &= e^{\frac{1}{2}(t_1-t_2)^2} e^{\frac{1}{2}(t_1+t_2)^2} \\ &= e^{t_1^2} e^{t_2^2} \\ &= e^{\frac{2t_1^2}{2}} e^{\frac{2t_2^2}{2}} \\ &= e^{t_1^2} e^{t_2^2} \\ &= m_{Y_1}(t_1) m_{Y_2}(t_2) \quad (\text{See next page}) \end{aligned}$$

Now $m_{Y_1 Y_2}(t_1, t_2) = m_{Y_1}(t_1) m_{Y_2}(t_2)$. We shall establish this fact.
when x_1, x_2 are independent

~~Assume~~ Let $Y_1 = g(x_1, x_2) u_1(x_1, x_2)$
 $Y_2 = h(x_1, x_2) u_2(x_1, x_2)$.

$$\begin{aligned} m_{Y_1 Y_2}(t_1, t_2) &= E[e^{t_1 Y_1 + t_2 Y_2}] \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{t_1 u_1(x_1, x_2) + t_2 u_2(x_1, x_2)} f_{x_1, x_2}(x_1, x_2) dx_1 dx_2 \end{aligned}$$

(18)

$$m_{Y_1, Y_2}(t_1, t_2) = e^{2t_1^2/2} e^{2t_2^2/2} \longrightarrow (\approx)$$

Now set $t_2 = 0$

$$m_{Y_1, Y_2}(t_1, 0) = E[e^{t_1 Y_1}] = m_{Y_1}(t_1)$$

$$\text{But } m_{Y_1, Y_2}(t_1, 0) = e^{2t_1^2/2}$$

$$\therefore m_{Y_1}(t_1) = e^{2t_1^2/2}$$

$$\Rightarrow Y_1 \sim N(0, 2)$$

Similarly setting $t_1 = 0$, we have $\cancel{Y_2 \sim N(0, 2)}$
 i.e. $m_{Y_2}(t_2) = e^{2t_2^2/2}$
 $\therefore Y_2 \sim N(0, 2)$.

$$\therefore m_{Y_1, Y_2}(t_1, t_2) = m_{Y_1}(t_1) m_{Y_2}(t_2).$$

Thus Y_1 & Y_2 are both independent and identically distributed normal random variable with mean 0 and variance 2

— x —

[End of the Probability Part]

Lecture 10: Random Sampling and Sampling Distribution

Sect 1: Sampling

Statistics concerns itself with the study of several aspects of a population, in which chance plays a role. The population can be human or otherwise. In statistics we are more concerned with target population in the study rather than whole population. Suppose we are looking for the sleeping patterns or sleeping hours of all adults in Kanpur in the age group 30-45. Then that particular segment is the target population. In most cases the target population is large enough so that collecting individual data is unviable and thus we need to seek a sub-population or a subset of a population called a "sample" on which we will carry out exhaustive studies/measurements. However to prevent bias creeping into the formation of a sample it is always advisable to randomize the process of hand drawing a sample.

Let there be a target population then we can consider two approaches to draw a random sample from it. The first one is sampling with replacement, where the sampled members of the population are returned back, while in the second approach they are not, since if we sample with replacement from a population of 100 a sample of each 10, each time, then in ten such draws we have exhausted the population.

So from a theoretical framework the theory of sampling or random sampling we consider a random sampling with replacement. Thus once a value is noted it becomes eligible again to be considered.

Consider any random variable X , (say height, weight, etc) which represents a characteristic of a population. Let X , follows a distribution with pmf/pdf, given as $f_X(x)$. Thus by the population we shall now mean $f_X(x)$, the pmf/pdf.

A random sample is a collection of $n \in \mathbb{N}$, identically and independently distributed random variables each having pmf/pdf

$f_{X_i}(x) = f_X(x)$, $\forall i=1, 2, \dots, n$. The random sample is represented as

$$x_1, x_2, \dots, x_n \quad (n: \text{size of the sample})$$

Thus for each i , $X_i : \Omega \rightarrow \mathbb{R}$, a random variable. Now what is this Ω ? $\omega \in \Omega$ is a collection n -individuals drawn from the population. Thus $X_i(\omega) = x_i$, is the value of characteristic of the i -th member chosen. So ω is in "some sense", the chosen physical sample. What happens is the following. Suppose $X \sim N(\mu, \sigma^2)$ but μ and σ^2 are not known though we do know that our characteristic follows normal distribution. By the vehicle of random sampling we will estimate these parameters.

Let x_1, \dots, x_n be a sample of size n . The following are two important parameter or statistic associated with a the random sample one the sample mean and sample variance, which are themselves, random variable given as.

$$\bar{X}_n = \frac{\sum_{i=1}^n x_i}{n}$$

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Why we have divided $(n-1)$ will be clear soon. It will be dealt in more detail when we will study point estimation.

$$E(\bar{X}) = \frac{1}{n} \sum_{i=1}^n E(x_i), \quad \text{Suppose } E(x_i) = \mu$$

Since x_1, \dots, x_n are iid r.v.'s we have $E(x_i) = \mu$, $\forall i$

$$\therefore E(\bar{X}) = \mu$$

$$\text{while } \text{Var}(\bar{X}) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n x_i\right) = \frac{1}{n} \sum \text{Var}(x_i).$$

Denoting $\text{Var}(x_i) = \sigma^2$ we have

$$\text{Var}(\bar{x}) = \sigma_{\bar{x}}^2 = \frac{n}{n^2} \sigma^2 = \frac{1}{n} \sigma^2$$

$$\therefore \text{Std.dev of } \bar{x} = \sigma_{\bar{x}} = \sqrt{\text{Var}(\bar{x})} = \frac{\sigma}{\sqrt{n}}$$

$\frac{\sigma}{\sqrt{n}}$ is often referred to as the standard error associated with ~~the~~ the sample

Our aim in this chapter is to know the distribution of \bar{x} and s^2 depending on the density of the population from which the sample is drawn. The following result is important enough to be stated as a theorem.

Theorem 10.1: Let x_1, \dots, x_n be a random sample from a population with density $f_x(\cdot)$, ~~and~~ Then if $n > 1$, we have

$$E[S^2] = \sigma^2$$

where σ^2 is the population variance. Further

$$\text{Var}(S^2) = \frac{1}{n} \left(\mu_4 - \frac{n-3}{n-1} \sigma^4 \right), \text{ for } n > 1$$

where $\mu_4 = E[X^4]$, the ~~ext~~ 4th population moment.

Proof: We shall only compute the mean. We know that

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

We shall first show that

$$\sum_{i=1}^n (x_i - \mu)^2 = \sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2$$

$$\begin{aligned} \sum_{i=1}^n (x_i - \mu)^2 &= \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - \mu)^2 \\ &= \sum_{i=1}^n (x_i - \bar{x})^2 + 2(\bar{x} - \mu) \sum_{i=1}^n (x_i - \bar{x}) + \sum_{i=1}^n (\bar{x} - \mu)^2 \\ &= \sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2. \quad (\because \sum_{i=1}^n (x_i - \bar{x}) \\ &\quad = \sum_{i=1}^n x_i - n\bar{x} = n\bar{x} - n\bar{x} = 0) \end{aligned}$$

$$\begin{aligned}
 E[S^2] &= E\left[\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2\right] \\
 &= \frac{1}{n-1} E\left[\sum_{i=1}^n (x_i - \mu)^2 - n(\bar{x} - \mu)^2\right] \\
 &= \frac{1}{n-1} \sum_{i=1}^n E[(x_i - \mu)^2] - \frac{n}{n-1} E[(\bar{x} - \mu)^2] \\
 &= \frac{n\sigma^2}{n-1} - \frac{n}{n-1} \text{Var}(\bar{x}) \\
 &= \frac{n\sigma^2}{n-1} - \frac{n}{n-1} \frac{\sigma^2}{n} \\
 &= \frac{n\sigma^2}{n-1} - \frac{\sigma^2}{n-1} \\
 &= \frac{(n-1)\sigma^2}{(n-1)} = \sigma^2 \\
 \therefore \boxed{E[S^2] = \sigma^2}
 \end{aligned}$$

So the sample mean and sample variance acts as a kind of estimator of the population mean and variance. However to compute $E(\bar{x})$ or $E(S^2)$, we need to know the distribution of \bar{x} and S^2 when the population distribution. The probability distributions is known. The probability distribution of \bar{x} and S^2 is known as Sampling distributions. This is what we discuss next.

See 2: Sampling Distributions

To begin with we will consider a population which obeys the normal distribution. Then we will seek in such a case the distribution of \bar{x} and S^2 .

Theorem 10.2: Let x_1, \dots, x_n be a random sample of n iid random variables, with mean \bar{x} . Let each $x_i \sim N(\mu, \sigma^2)$, $i=1, 2, \dots, n$. Then $\bar{x} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$

Proof: We will use the mgf technique

$$\begin{aligned}
 m_{\bar{x}}(t) &= E[e^{t\bar{x}}] \\
 &= E\left[e^{t \frac{\sum x_i}{n}}\right] \\
 &= E\left[e^{t \frac{\sum t x_i}{n}}\right] \\
 &= E\left[\prod_{i=1}^n e^{\frac{tx_i}{n}}\right] \\
 &= \prod_{i=1}^n E\left[e^{\frac{tx_i}{n}}\right], \text{ by independence} \\
 &= \prod_{i=1}^n m_{x_i}\left(\frac{t}{n}\right) \\
 &= \prod_{i=1}^n \left(\frac{\mu t}{n} + \frac{1}{2} \frac{\sigma^2 t^2}{n^2}\right) \\
 &= \prod_{i=1}^n e^{\left(\mu t + \frac{1}{2} \frac{\sigma^2 t^2}{n}\right)} \\
 &= e^{\left(\mu t + \frac{1}{2} \frac{\sigma^2 t^2}{n}\right)}
 \end{aligned}$$

$$\therefore \bar{x} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \quad \square$$

Now we shall focus on how to find the distribution of

$$S^2 = \frac{\sum_{i=1}^{n-1} (x_i - \bar{x})^2}{n-1}$$

when we draw a random sample of size n from a normal population.

A very special type of probability distribution that will play a crucial role in computing the probability distribution of S^2 is the so called chi-square distribution or χ^2 -distribution.

A random variable or rather a continuous random variable X is said to be a χ^2 random variable (chi-square r.v.) if X has the density

$$f_X(x) = \frac{1}{\Gamma(k/2)} \left(\frac{1}{2}\right)^{k/2} x^{k/2-1} e^{-\frac{1}{2}x}, \quad x \in (0, \infty)$$

$$f_X(x) = 0, \quad \text{otherwise.}$$

This a Gamma distribution with $r = k/2$ & $\lambda = 1/2x$. k is called the degrees of freedom of the χ^2 -random variable X .

So using our information about Gamma distribution we can write. If $X \sim \chi^2(k)$, (ie χ^2 with k -degrees of freedom)

$$E(X) = \frac{r}{\lambda} = \frac{k/2}{1/2} = k$$

$$\text{Var}(X) = \frac{r}{\lambda^2} = \frac{\left(\frac{k}{2}\right)}{\left(\frac{1}{2}\right)^2} = 2k$$

$$m_X(t) = \left(\frac{\lambda}{\lambda-t}\right)^r, \quad t < \lambda$$

$$= \left(\frac{\frac{1}{2}}{\frac{1}{2}-t}\right)^{k/2}; \quad t < \frac{1}{2}, \text{ for}$$

$$= \left[\frac{1}{1-2t}\right]^{k/2}; \quad t < \frac{1}{2}$$

We now state the following important theorem

Theorem 10.3 : Let $X_i, i=1, 2, \dots, k$ are continuous random variables and $X_i \sim N(\mu_i, \sigma_i^2)$ and are independent. Then

$$U = \sum_{i=1}^{2k} \left(\frac{X_i - \mu_i}{\sigma_i} \right)^2$$

follows χ^2 distribution with k -degrees of freedom.

Proof: We will again use the mgf technique which validates its powers. Let us write

$$Z_i = \left(\frac{X_i - \mu_i}{\sigma_i} \right)$$

then $Z_i \sim N(0, 1)$, i.e. Z_1, \dots, Z_k are iid r.v.s

$$\therefore U = \sum_{i=1}^k Z_i^2.$$

Hence

$$\begin{aligned} m_U(t) &= E[e^{tU}] \\ &= E[e^{t \sum_{i=1}^k Z_i^2}] \\ &= E[e^{\prod_{i=1}^k t Z_i^2}] \\ &= \prod_{i=1}^k E[e^{t Z_i^2}] \quad (\text{by independence}) \end{aligned}$$

$$\begin{aligned} \text{Now } E[e^{t Z_i^2}] &= \int_{-\infty}^{\infty} e^{t z_i^2} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} z_i^2} dz \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(1-2t)z^2} dz \end{aligned}$$

$$\text{Set } (\sqrt{1-2t})z = v \Rightarrow z^2 = \frac{v^2}{1-2t}, \quad \text{for } t < \frac{1}{2}$$

$$\therefore dv = \sqrt{1-2t} dz \Rightarrow dz = \frac{dv}{\sqrt{1-2t}}$$

$$\begin{aligned} \therefore E[e^{tZ_i^2}] &= \frac{1}{\sqrt{1-2t}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dx \\ &= \frac{1}{\sqrt{1-2t}} \\ \therefore m_V(t) = \prod_{i=1}^k E[e^{tZ_i^2}] &= \left(\frac{1}{1-2t}\right)^{k/2}, \quad t < \frac{1}{2}. \end{aligned}$$

Thus V follows χ^2 with k -degrees of freedom.

Thus if X_1, \dots, X_n is a random sample of size n , where for each i , $X_i \sim N(\mu, \sigma^2)$, then

$$U = \sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma} \right)^2 \rightarrow A$$

follows a χ^2 distribution with n -degrees of freedom.

Question !! In the expression of U as given in A . If either μ or σ^2 or both are unknown, will you call U a ~~statistic~~ statistic?

We shall now state the following Theorem without proof

Theorem 10.4 Let X_1, \dots, X_n , be a random sample drawn from a standard normal distribution.

- i) Then $\bar{X} \sim N(0, \frac{1}{n})$
- ii) \bar{X} and $\sum_{i=1}^n (X_i - \bar{X})^2$ are independent
- iii) $\sum_{i=1}^n (X_i - \bar{X})^2$ has a χ^2 -distribution with $(n-1)$ degrees of freedom.

Proof: Part i) is already proved in Theorem 10.2. We shall just prove part iii) here as needed for our purpose

To prove iii) observe that

$$\sum_{i=1}^n x_i^2 = \sum_{i=1}^n (x_i - \bar{x})^2 + n\bar{x}^2 \quad (\text{The reader } \exists \text{ should prove it})$$

$$\therefore m \sum_{i=1}^n x_i^2 (t) = m \sum_{i=1}^n (x_i - \bar{x})^2 (t) + m_n \bar{x}^2 (t) \quad (\text{Again the reader should prove it})$$

$$\therefore m \sum (x_i - \bar{x})^2 (t) = \frac{m \sum x_i^2 (t)}{m_n \bar{x}^2 (t)}$$

Now as $x_i \sim N(0, 1)$, by using Theorem 10.3 we

know that $\sum x_i^2 \sim \chi^2$ (or) with n degrees of freedom

$$\therefore m \sum x_i^2 (t) = \left[\frac{1}{(1-2t)} \right]^{n/2}; \quad t < \frac{1}{2}$$

Since $\bar{x} \sim N(0, \frac{1}{n})$, thus $Y = n\bar{x}^2$ satisfies follows a χ^2 distribution with degrees of freedom 1 as $\Gamma(\frac{1}{2}) = \sqrt{\pi}$
 (The student / reader should be able to prove this)

$$\therefore m_n \bar{x}^2 (t) = \left[\frac{1}{(1-2t)} \right]^{\frac{1}{2}}; \quad t < \frac{1}{2}$$

$$\therefore m \sum (x_i - \bar{x})^2 (t) = \left(\frac{1}{1-2t} \right)^{\frac{n-1}{2}}; \quad t < \frac{1}{2}$$

$\therefore \sum (x_i - \bar{x})^2 \sim \chi^2$ (or) distribution with $(n-1)$ degrees of freedom. \square

So how do we write Theorem 10.4, if we write have x_1, \dots, x_n a random sample from a normal distribution with mean μ and σ^2 ?

Note that we just set $Z_i = \frac{x_i - \mu}{\sigma}$

$$\therefore \bar{Z} = \frac{(\bar{x} - \mu)}{\sigma}$$

and hence $\frac{\bar{x} - \mu}{\sigma} \sim N(0, \frac{1}{n})$, i.e what we

get from ~~(10.4)~~ Theorem i) of Theorem 10.4

$$\text{Now } \sum_{i=1}^n (Z_i - \bar{Z})^2 = \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma} \right)^2$$

$\therefore \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma} \right)^2 \sim \chi^2$ with $(n-1)$ degrees of freedom from iii)' in Theorem 10.4

$$\text{Now } S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\text{Hence } \frac{(n-1)S^2}{\sigma^2} = \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma} \right)^2$$

Thus $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2$ with $(n-1)$ degrees of freedom.

So how do we find the p.d.f of S^2

$$\text{Set } U = \frac{(n-1)S^2}{\sigma^2}$$

We have seen that $U \sim \chi^2$ with $(n-1)$ degrees of freedom. Set $S^2 = Y$.

$$\therefore U = \frac{(n-1)Y}{\sigma^2} \text{ a } Y = \frac{\sigma^2 U}{(n-1)}$$

From standard transformation technique, (since the transformation is linear), we have for $n > 1$. & $y > 0$ ($\because U > 0$)

$$\begin{aligned}
 f_{S^2}(y) &= f_U\left(\frac{(n-1)y}{\sigma^2}\right) \left| \frac{dy}{du} \right| \\
 &= \frac{1}{\Gamma\left(\frac{n-1}{2}\right)} \left(\frac{1}{2}\right)^{\frac{n-1}{2}} \left[\frac{(n-1)y}{\sigma^2}\right]^{\frac{n-1}{2}-1} e^{-\frac{1}{2}\frac{(n-1)y}{\sigma^2}} \frac{(n-1)}{\sigma^2} \\
 &= \frac{1}{\Gamma\left(\frac{n-1}{2}\right)} \left(\frac{1}{2}\right)^{\frac{n-1}{2}} \frac{(n-1)^{\frac{n-3}{2}}}{\sigma^{(n-3)}} \cdot y^{\frac{n-3}{2}} e^{\frac{(n-1)y}{2\sigma^2}} \cdot \frac{(n-1)}{\sigma^2} \\
 &= \frac{1}{\Gamma\left(\frac{n-1}{2}\right)} \left[\frac{(n-1)}{(2\sigma^2)} \right]^{\frac{n-1}{2}} y^{\frac{n-3}{2}} e^{\frac{(n-1)y}{2\sigma^2}}
 \end{aligned}$$

$$\therefore \boxed{f_{S^2}(y) = \frac{1}{\Gamma\left(\frac{n-1}{2}\right)} \left(\frac{n-1}{2\sigma^2}\right)^{\frac{n-1}{2}} y^{\frac{n-3}{2}} e^{\frac{(n-1)y}{2\sigma^2}}} \quad \text{for } y > 0, n > 1$$

Of course $f_{S^2}(y) = 0$, for $y \leq 0$.

Our next object of focus is the F-distribution named after the celebrated statistician Ronald Fisher. This distribution plays quite an important role in statistics. Let us see how an F-distribution pdf is constructed. Let $U \sim \chi^2$ with m -degrees of freedom and $V \sim \chi^2$ with n degrees of freedom. Let U and V be independent, then their joint density is given as

$$f_{U,V}(u,v) = \frac{1}{\Gamma\left(\frac{m}{2}\right)\Gamma\left(\frac{n}{2}\right)} u^{\frac{m-2}{2}} v^{\frac{n-2}{2}} e^{-\frac{1}{2}(u+v)},$$

$$f_{U,V}(u,v) = \frac{1}{\Gamma\left(\frac{m}{2}\right)\Gamma\left(\frac{n}{2}\right)} 2^{\frac{(m+n)}{2}} u^{\frac{m-2}{2}} v^{\frac{n-2}{2}} e^{-\frac{1}{2}(u+v)}, \quad 0 < u < \infty, 0 < v < \infty$$

where $0 \leq u \leq \infty, 0 \leq v \leq \infty$.

Let us consider transform

$$X = \frac{U}{\frac{m}{n}}$$

We shall call such a random variable X a F-variable.

Consider the transformation

$$X = \frac{U}{\frac{m}{n}} \quad \& \quad Y = V$$

We seek to compute $f_{X,Y}(x,y)$, since we know

$f_{U,V}(u,v)$. Now

$$\text{Jacobian} \quad U = \frac{m}{n} \times Y$$

$$V = Y$$

$$\therefore \text{Jacobian matrix } J = \begin{bmatrix} \frac{\partial U}{\partial x} & \frac{\partial U}{\partial y} \\ \frac{\partial V}{\partial x} & \frac{\partial V}{\partial y} \end{bmatrix} = \begin{bmatrix} \frac{m}{n}y & \frac{m}{n}x \\ 0 & 1 \end{bmatrix}$$

$$\therefore \det J = \frac{m}{n}y \Rightarrow |\det J| = \frac{m}{n}y. \quad (\text{as } m > 0, n > 0, y > 0)$$

$$\therefore f_{X,Y}(x,y) = \frac{m}{n}y \cdot \frac{1}{\Gamma\left(\frac{m}{2}\right)\Gamma\left(\frac{n}{2}\right)2^{\frac{m+n}{2}}} \left(\frac{m}{n}xy\right)^{\frac{m-2}{2}} y^{\frac{(n-2)}{2}} e^{-[\frac{m}{n}xy+y]/2}$$

where, $0 < x < \infty, 0 < y < \infty$

$$\begin{aligned} \therefore f_X(x) &= \int_0^\infty f_{X,Y}(x,y) dy \\ &= \frac{1}{\Gamma\left(\frac{m}{2}\right)\Gamma\left(\frac{n}{2}\right)2^{\frac{m+n}{2}}} \left(\frac{m}{n}\right)^{\frac{m}{2}} x^{\frac{(m-2)}{2}} \int_0^\infty y^{\frac{m+n-2}{2}} e^{-\frac{1}{2}\left[\frac{m}{n}x+1\right]y} dy \\ &\quad \text{Set } \frac{1}{2}\left[\frac{m}{n}x+1\right]y = t \Rightarrow \frac{dt}{\frac{1}{2}\left[\frac{m}{n}x+1\right]} = dy \end{aligned}$$

$$f_X(x) = \frac{\Gamma\left(\frac{m+n}{2}\right)}{\Gamma\left(\frac{m}{2}\right)\Gamma\left(\frac{n}{2}\right)} \left(\frac{m}{n}\right)^{m/2} \frac{x^{\frac{(m-2)}{2}}}{\left[\frac{m}{n}x + 1\right]^{\frac{(m+n)}{2}}} \quad \text{--- } \cancel{\text{Step 2}}$$

Note that

$$\int_0^\infty y^{\frac{m+n-2}{2}} e^{-t} dt = \Gamma\left(\frac{m+n}{2}\right)$$

↓
(after substitution)

$$\therefore f_X(x) = \frac{\Gamma\left(\frac{m+n}{2}\right)}{\Gamma\left(\frac{m}{2}\right)\Gamma\left(\frac{n}{2}\right)} \left(\frac{m}{n}\right)^{m/2} \frac{x^{(m-2)/2}}{\left[\frac{m}{n}x + 1\right]^{\frac{(m+n)}{2}}} \quad x \in (0, \infty)$$

↓
F-distribution (We say
 $X \sim F$ with (m, n) degrees of freedom)

F-variable $X = \frac{U}{m} / \frac{V}{n}$ is often called

variance ratio.

But how is this discussion relevant to sampling distributions. This is how it is.

Let X_1, \dots, X_{m+1} be a sample from a normal distribution with mean μ_1 & variance σ^2 & Y_1, \dots, Y_{n+1} is a sample of size $(n+1)$ drawn from a normal distribution with

mean μ_2 and variance σ^2 , then we know that

$$\frac{1}{\sigma^2} \sum_{i=1}^{m+1} (X_i - \bar{X})^2 \sim \chi^2 \text{ with } m \text{ degrees of freedom}$$

and

$$\frac{1}{\sigma^2} \sum_{i=1}^{n+1} (Y_i - \bar{Y})^2 \sim \chi^2 \text{ with } n \text{ degrees of freedom}$$

~~The statistic~~

Consider the statistic $T = \frac{\frac{1}{\sigma^2} \sum_{i=1}^{m+1} (x_i - \bar{x})^2}{\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \bar{x})^2}$

$$\therefore T = \frac{\frac{1}{m} \sum_{i=1}^{m+1} (x_i - \bar{x})^2}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

Then $T \sim F$ with $m > n$ degrees of freedom.

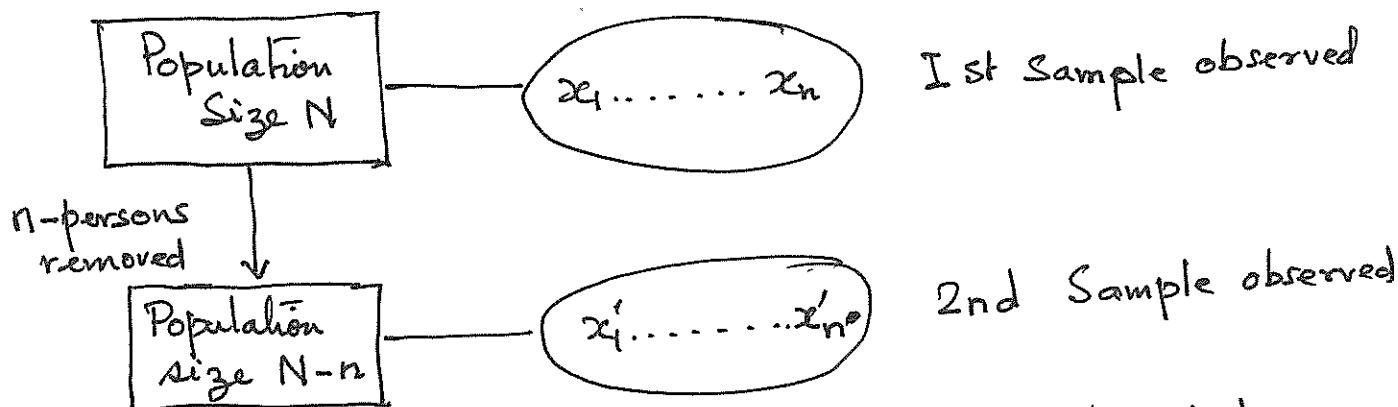
— x —

Note: Though the material here is standard, a large part of the discussion here is based on Chapter 6 of Mood, Graybill and Boes, "Introduction to the Theory of Statistics" Statistics 1974, Mc. Graw Hill.

Section 3: Sampling Without Replacement

Till now we have studied sampling with replacement, i.e. a value which appeared in the first sample can also appear in the second, third or any other sample. However when we ~~do~~ ^{don't do}, so at every step we reduce the population size and thus bring in dependency among sample observation.

Diagrammatically



Think of the urn-model, where you draw a ball but do not return to the urn

Given a population of size N , a sample where each member can appear only once ~~in a set~~ and if we consider $n < N$ members in a sample, then there are $\binom{N}{n}$ possible samples that can be drawn from such a population. What would be the standard error or $\sqrt{\text{Var}(\bar{X})}$ in this case. So we have to compute this quantity. Note that in such a scenario we can always ask the questions, ~~what~~ that what is the probability that the k -th member of the population is the j -th member of the sample given that ~~the~~ k -th member of the population will be i -th member of the sample. Note that

$$\text{Var}(\bar{X}) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \text{cov}(X_i, X_j)$$

So if we consider sampling without replacement, then

$$\text{Cov}(x_i, x_j) = -\frac{\sigma^2}{N-1}, \text{ if } i \neq j$$

Now observe that

$$\text{Cov}(x_i, x_j) = E(x_i x_j) - E(x_i)E(x_j).$$

~~Now considering that we draw the sample from a discrete distribution~~

$$\begin{aligned} E(x_i, x_j) &= \sum_{k=1}^m \sum_{l=1}^m \bar{x}_k \bar{x}_l P(X_i = \bar{x}_k \text{ and } X_j = \bar{x}_l) \\ &= \sum_{k=1}^m \sum_{l=1}^m \bar{x}_k \bar{x}_l P(X_i = \bar{x}_k) P(X_j = \bar{x}_l | X_i = \bar{x}_k) \\ &= \sum_{k=1}^m \bar{x}_k P(X_i = \bar{x}_k) \sum_{l=1}^n \bar{x}_l P(X_j = \bar{x}_l | X_i = \bar{x}_k) \end{aligned}$$

Now we have used the summation formula because we have a sample population of size N and hence, the sample observations are values of a random variable from a discrete distribution. Note this idea of sampling without replacement makes sense only in that case.

$$P(X_j = \bar{x}_l | X_i = \bar{x}_k) = ??$$

How to compute this. We have to now make clear certain issues, before we proceed. Let us assume that values of the characteristic of population members are m -values given as $\bar{x}_1, \dots, \bar{x}_m$. We will assume that n_j of the N population members take the value \bar{x}_j (Think of how you construct the histogram)

$$\therefore \cancel{P(x_i)} \quad P(X_i = \bar{x}_j) = \frac{n_j}{N}$$

Now in our case thus $P(X_i = \xi_k) = \frac{n_k}{N}$

Further,

$$P(X_j = \xi_l | X_i = \xi_k) = \frac{n_l}{N-1} \quad \text{if } k \neq l$$

$$P(X_j = \xi_l | X_i = \xi_k) = \frac{n_{k-1}}{N-1} \quad \text{if } k = l.$$

$$\begin{aligned} \therefore \sum_{i=1}^m \xi_i P(X_j = \xi_l | X_i = \xi_k) \\ = \sum_{l \neq k} \xi_l \frac{n_l}{N-1} + \xi_k \frac{n_{k-1}}{N-1} \end{aligned}$$

$$\begin{aligned} \therefore E(X_i, X_j) &= \sum_{k=1}^m \xi_k \cancel{\left(\frac{n_k}{N} \right)} \left(\sum_{l \neq k} \xi_l \frac{n_l}{N-1} + \xi_k \frac{n_{k-1}}{N-1} \right) \\ &= \sum_{k=1}^m \xi_k \frac{n_k}{N} \left(\sum_{l=1}^m \xi_l \frac{n_l}{N-1} - \frac{\xi_k}{N-1} \right) \end{aligned}$$

To proceed further we need more information. But for now we have

$$E(X_i, X_j) = \frac{1}{N(N-1)} \left(\sum_{k=1}^m \sum_{j=1}^l \xi_k \xi_l n_k n_l - \sum_{k=1}^m \underline{\xi_k n_k} \right)$$

Observe that if we look at the sum of all n -observed values of the sample and call it \bar{x} , we have

$$\bar{x} = x_1 + x_2 + \dots + x_n$$

So the observed value

$$\hat{\bar{x}} = x_1 + x_2 + \dots + x_n$$

$$\hat{\bar{x}} = \sum_{i=1}^m n_i \xi_i \quad \textcircled{17}$$

Hence

$$\hat{\sigma}^2 = \left(\sum_{j=1}^m n_j \xi_j \right)^2 \\ = \sum_{k=1}^m \sum_{l=1}^m \xi_k \xi_l n_k n_l$$

Also observe that the ^{population} sample variance is calculated as

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

\therefore The observed value of σ^2 which we denote as $\hat{\sigma}^2$ is given as

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N x_i^2 - \hat{\mu}^2$$

where $\hat{\mu} = \frac{\sum x_i}{N}$, is the observed population mean

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{k=1}^m (\xi_k^2 n_k)$$

$$\hat{\sigma}^2 = \frac{1}{N} \left(\sum_{k=1}^m n_k \xi_k^2 \right) - \hat{\mu}^2$$

$$\Rightarrow \hat{\sigma}^2 + \hat{\mu}^2 = \frac{1}{N} \left(\sum_{k=1}^m n_k \xi_k^2 \right)$$

Thus

$$\Rightarrow E(x_i x_j) = \frac{\hat{\sigma}^2}{N(N-1)} - \frac{1}{N(N-1)} \left(\sum_{k=1}^m n_k \xi_k^2 \right)$$

$$= \frac{\hat{\sigma}^2}{N(N-1)} - \frac{(\hat{\sigma}^2 + \hat{\mu}^2)}{N-1}$$

$$= \frac{N^2 \hat{\mu}^2}{N(N-1)} - \frac{(\hat{\sigma}^2 + \hat{\mu}^2)}{N-1}$$

$$\begin{aligned}
 E(x_i x_j) &= \frac{N \hat{\mu}^2}{N-1} - \frac{(\hat{\sigma}^2 + \hat{\mu}^2)}{N-1} \\
 &= \frac{-\hat{\sigma}^2 + (N-1)\hat{\mu}^2}{(N-1)} \\
 &= -\frac{\hat{\sigma}^2}{N-1} + \hat{\mu}^2
 \end{aligned}$$

$$\begin{aligned}
 \therefore \text{Cov}(x_i, x_j) &= E(x_i x_j) - E(x_i)E(x_j) \\
 &= -\frac{\hat{\sigma}^2}{N-1} + \hat{\mu}^2 - \hat{\mu}^2 \\
 &= -\frac{\hat{\sigma}^2}{N-1}
 \end{aligned}$$

$$\text{Cov}(x_i, x_j) = -\frac{\hat{\sigma}^2}{N-1}$$

$\hat{\sigma}^2 = \text{observed population mean.}$

Now we have to compute $\text{Var}(\bar{x})$, and thus

$$\begin{aligned}
 \text{Var}(\bar{x}) &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \text{Cov}(x_i, x_j) \\
 &\Rightarrow \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \\
 &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(x_i) + \frac{1}{n^2} \sum_{i=1}^n \sum_{i \neq j} \text{Cov}(x_i, x_j) \\
 &= \frac{1}{n^2} \cdot n \hat{\sigma}^2 + \frac{1}{n^2} n(n-1) \left[-\frac{\hat{\sigma}^2}{N-1} \right] \\
 &= \frac{\hat{\sigma}^2}{n} \left(1 - \frac{n-1}{N-1} \right)
 \end{aligned}$$

$$\text{Var}(\bar{x}) = \frac{\hat{\sigma}^2}{n} \left(1 - \frac{n-1}{N-1} \right)$$

The term

$$\left[1 - \frac{n-1}{N-1} \right]$$

is called the finite population correction.

The above discussion is largely based on Chapter 7 of the Mathematical Statistics and Data Analysis, by John. A. Rice (Cengage Learning: 2007).

Lecture 11: Parametric Estimation

In this lecture we are going to devise methods to find estimators of population parameters satisfying certain criterion. A population parameter is estimated using what is called a sample statistic.

A sample statistic is a function of a random sample. Let us consider a population $f_x(\cdot; \theta)$, where θ denotes the population parameter.

If θ is known then, the distribution of X is completely known. The idea of estimation comes in when θ is not known. In fact θ can be a vector or a scalar. If $X \sim \text{Poisson}(\lambda)$, then $\theta = \lambda$, a real scalar, whereas if $X \sim N(\mu, \sigma^2)$, then $\theta = (\theta_1, \theta_2) \in \mathbb{R}^2$ with $\theta_1 = \mu$ & $\theta_2 = \sigma^2$.

For the moment assume $\theta \in \mathbb{R}$. Then having a random sample from x_1, \dots, x_n of size n from $f_x(\cdot, \theta)$, we construct a statistic

Symbolically given as

$$T = l(x_1, x_2, \dots, x_n)$$

We say that the random variable T is an unbiased estimator of θ if

$$E(T) = \theta.$$

Now the expectation of T is calculated based on the sampling distribution of T .

The expression $E(T) - \theta$ is called the bias in the estimation if $E(T) - \theta \neq 0$.

If $\theta \in \mathbb{R}^k$ say i.e. $\theta = (\theta_1, \dots, \theta_k)$, then we can have

k different ^{statistics} estimators, $T_i = l_i(x_1, \dots, x_n)$, $i=1, \dots, k$

which can acts as estimator of θ .

Note that \bar{x} , the sample mean & s^2 the sample variance have been proved in the last chapter as unbiased estimators of population ^{mean} variance and population variance respectively.

Such an approach to get an approximation for θ is often called point estimation. The statistic T when viewed as a random variable is called an estimator and the specific value of T is called an estimate. Sometimes $\hat{\theta}$ is denoted as an estimate of θ . There are several approaches to find an estimator. We would also sometimes refer $\hat{\theta}$ as an estimator of θ corresponding to the estimate $\hat{\theta}$.

Section 1: Method of Moments

Let x_1, \dots, x_n be a random sample drawn from a population of size n
 $f_x(x; \theta_1, \dots, \theta_k)$, where $\theta \in \mathbb{R}^k$, with $\theta = (\theta_1, \dots, \theta_k)$. Let us recall the r -th moment of X , i.e

$$\mu'_r = E[X^r]$$

$$\therefore \mu'_r = \int x^r f_x(x, \theta_1, \dots, \theta_k) dx$$

We can write

$$\mu'_r = \mu'_r(\theta_1, \dots, \theta_k)$$

The r -th sample moment

$$M'_r = \frac{x_1^r + x_2^r + \dots + x_n^r}{n} = \frac{1}{n} \sum_{i=1}^n x_i^r$$

To find $\theta = (\theta_1, \dots, \theta_k)$, set up the k -equations, in $\theta_1, \dots, \theta_k$

$$\mu'_j(\theta_1, \dots, \theta_k) = M'_j, \quad j=1, \dots, k$$

Suppose $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k$ are the solutions to the equations. If say $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k$ is unique, then we can consider $\hat{\theta}_1, \dots, \hat{\theta}_k$ to be the estimates of $\theta_1, \dots, \theta_k$.

E.g 1: x_1, \dots, x_n be a random sample from a exponential population

$$f_x(x, \theta) = \theta e^{-\theta x}; \quad \theta > 0, \quad x \in (0, \infty).$$

Thus we will solve ~~one~~ one equation

$$\therefore M'_1 = \mu'_1 = \mu'_1(\theta) = \frac{1}{\theta}$$

$\therefore \hat{\theta} = \frac{1}{M'_1}$ So if we take a particular sample
 ~~M'_1~~ $x_1 = x_1, \dots, x_n = x_n$, then $\hat{\theta} = \frac{n}{x_1 + \dots + x_n}$

So $\hat{\theta}$ is estimated by $\frac{1}{\bar{x}} = \frac{n}{x_1 + \dots + x_n}$

$$\therefore \hat{\theta} = l(\hat{\theta})(x_1, \dots, x_n) = \frac{n}{x_1 + \dots + x_n}$$

$\hat{\theta}$ is the estimator.

E.g. 2: x_1, \dots, x_n be a random sample of size n , drawn from a normal population with mean μ and σ^2 variance σ^2 .
 $\therefore \theta = (\mu, \sigma^2)$.

$$M'_1 = \mu'_1 = \sigma' \mu$$

$$M'_2 = \mu'_2 = \sigma^2 + \mu^2$$

μ is estimated by $M'_1 = \bar{x}$, i.e. \bar{x} is the unbiased estimator of μ .

σ^2 is estimated by the statistic

$$\hat{\theta}_2 = M'_2 - \frac{\mu^2}{n}$$

$$\therefore \hat{\theta}_2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \frac{\bar{x}^2}{n}$$

$$\text{Now } \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \left[\sum_{i=1}^n x_i^2 - 2\bar{x} \sum x_i + n\bar{x}^2 \right]$$

$$= \frac{1}{n} \sum x_i^2 - 2\bar{x} \underbrace{\frac{1}{n} \sum x_i}_{\bar{x}} + \frac{n\bar{x}^2}{n}$$

$$= \frac{1}{n} \sum_{i=1}^n x_i^2 - 2\bar{x}^2 + \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$$

(3)

$$\hat{\Theta} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$\hat{\Theta}$ is the estimator of σ^2 and note that through the method of moments we do not have S^2 as the estimator of σ^2 . Thus $\hat{\Theta}$ is not an un-biased estimator of σ^2 . Thus method of moments need not give us unbiased estimators always.

Section 2: Maximum Likelihood Estimator

The technique of maximum likelihood function is based on using the fact that x_1, \dots, x_n are iid random variables.

Let us first look at the simple case where $\theta \in \mathbb{R}$. Let us draw a sample of size n , from a population $f_X(x; \theta)$.

The likelihood function in this case is defined

as

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

$$L(x, \theta) = L(\theta) = f_{X_1, \dots, X_n}(x_1, \dots, x_n, \theta) \rightarrow (\text{By independence})$$

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

$$\therefore L(\theta) = \prod_{i=1}^n f_{X_i}(x_i, \theta)$$

We compute $\max_{\theta \in A} L(\theta)$, where A is the set in which θ is restricted to belong. We find that value of θ , which maximizes

$L(\theta) = L(\theta; x_1, \dots, x_n)$, for a chosen sample. Then if $\hat{\theta}$ is the maximizer for a chosen sample i.e. $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$, then $\hat{\theta}(x_1, \dots, x_n)$ is called an maximum likelihood estimator estimate of θ for the chosen sample values, $X_1 = x_1, \dots, X_n = x_n$. So the estimator of θ is $\hat{\Theta} = \hat{\theta}(x_1, \dots, x_n)$.

thus $\hat{\Theta}$ is called the maximum likelihood estimator of θ

Noting that $\ln(\text{base loge})$ is an increasing function. We know that if $L(\theta) > 0, \forall \theta \in A$, then if $\hat{\theta}$ is the minimizer, we have

$$L(\theta) \leq L(\hat{\theta}), \quad \forall \theta \in A$$

$$\Rightarrow \ln L(\theta) \leq \ln L(\hat{\theta}) \quad \forall \theta \in A$$

Thus $\hat{\theta}$ is also the maximizer of $\ln L(\theta)$, over A . In most cases the set A is an open interval. Thus to compute $\hat{\theta}$ we can use the equation

$$\frac{\partial}{\partial \theta} \ln L(\theta) = 0$$

$$\frac{\partial}{\partial \theta} \ln L(\theta) = 0$$

$$\boxed{\frac{\partial}{\partial \theta} \ln L(\theta) = 0} \rightarrow (\ast)$$

Then also check the second order condition $\left. \frac{\partial^2}{\partial \theta^2} \ln L(\theta) \right|_{\theta=\hat{\theta}} < 0$,

where $\left. \frac{\partial}{\partial \theta} \ln L(\theta) \right|_{\theta=\hat{\theta}} = 0$, $\hat{\theta}$ solves (\ast) . Such a $\hat{\theta}$ will provide an estimator or maximum likelihood estimator using the expression

$$\hat{\theta} = \hat{\theta}(x_1, \dots, x_n).$$

Now suppose $\theta \in \mathbb{R}^k$, and we draw a random sample of size n x_1, \dots, x_n , from a population described by $f_x(x; \theta) = f_{x_1}(x_1, \theta_1, \dots, \theta_k) \cdots f_{x_n}(x_n, \theta_1, \dots, \theta_k)$.

In this case we have the likelihood function as

$$\begin{aligned} L(x, \theta_1, \dots, \theta_k) &= f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta_1, \dots, \theta_k) \\ &= f_{x_1}(x_1, \theta_1, \dots, \theta_k) \cdot f_2(x_2, \theta_1, \dots, \theta_k) \cdots f_n(x_n, \theta_1, \dots, \theta_k) \\ &= \prod_{i=1}^n f_{x_i}(x_i, \theta_1, \dots, \theta_k). \end{aligned}$$

as before
 \therefore Similarly, the maximum likelihood function estimation problem is given as.

$$\max_{\theta \in K} L(x, \theta_1, \dots, \theta_k), \quad \text{where } K \subset \mathbb{R}^k$$

$$\max_{\theta \in K} \ln L(x, \theta_1, \dots, \theta_k), \quad \text{where } K \subset \mathbb{R}^k$$

(5)

Eg Now if $\hat{\theta}$ is the maximizer, then as before we have $\theta = (\theta_1, \dots, \theta_n)$

$$\hat{\theta}_i = \hat{\theta}_i(x_1, \dots, x_n), \quad i=1, \dots, n$$

So there are n , maximum likelihood estimators. We shall now provide two examples. In the first case where $\theta \in \mathbb{R}$, and in the second case $\theta \in \mathbb{R}^2$.

E.g. 3: Let x_1, \dots, x_n be a random sample drawn from an exponential population $f_X(\cdot, \theta)$, and is given as

$$f_X(x, \theta) = \theta e^{-\theta x}, \quad x > 0, \quad \theta > 0.$$

$$\therefore L(\theta, x) = L(\theta) = \prod_{i=1}^n f_{X_i}(x_i, \theta)$$

$$= \prod_{i=1}^n \theta e^{-\theta x_i}$$

$$= \prod_{i=1}^n (\theta)^n e^{-\theta \sum_{i=1}^n x_i}$$

$$\therefore \ln L(\theta) = n \ln \theta + (-\theta \sum_{i=1}^n x_i)$$

The maximum likelihood estimation problem (MLE)
 $\max_{\theta > 0} L(x) \ln L(\theta)$ (P1)

Thus $\frac{\partial}{\partial \theta} \ln L(\theta) = n \cdot \frac{1}{\theta} - \sum_{i=1}^n x_i = 0$

$$\therefore \frac{n}{\theta} = \sum_{i=1}^n x_i$$

$$\therefore \frac{1}{\theta} = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}$$

$$\text{or } \theta = \frac{1}{\bar{x}}$$

Now $\frac{\partial^2}{\partial \theta^2} \ln L(\theta) = -\frac{n}{\theta^2} < 0$ } Thus $\theta = \frac{1}{\bar{x}}$ maximizes (strictly) problem P1).

$$= -n \bar{x}^2 < 0 \quad \therefore \hat{\theta} = \frac{1}{\bar{x}} \text{ is the mle estimator of } \theta$$

(6) mle: Short form for "maximum likelihood estimator"

Now for the case when there are more than one parameter, i.e. $\theta \in \mathbb{R}^k$, we first have to find $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$ satisfying

$$\frac{\partial \ln L(\theta_1, \dots, \theta_k)}{\partial \theta_j} = 0; \text{ for } j=1, \dots, k.$$

Once we find a $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k)$, we compute the Hessian matrix of $\nabla^2 \ln L(\theta_1, \dots, \theta_k)$, and see if this negative definite at $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$. In the above expression we took the liberty of writing $L(x, \theta_1, \dots, \theta_k) = L(\theta_1, \dots, \theta_k)$.

When $\theta \in \mathbb{R}^2$, i.e. $\theta = (\theta_1, \theta_2)$, then we have to check the following. First find $(\hat{\theta}_1, \hat{\theta}_2)$ satisfying

$$\left. \begin{aligned} \frac{\partial}{\partial \theta_1} \ln L(\theta_1, \theta_2) &= 0 \\ \frac{\partial}{\partial \theta_2} \ln L(\theta_1, \theta_2) &= 0 \end{aligned} \right\}$$

Then the Hessian matrix of $\ln L(\theta_1, \theta_2)$ is given as

$$\nabla^2 \ln L(\theta_1, \theta_2) = \begin{bmatrix} \frac{\partial^2}{\partial \theta_1^2} \ln L & \frac{\partial^2}{\partial \theta_1 \partial \theta_2} \ln L \\ \frac{\partial^2}{\partial \theta_2 \partial \theta_1} \ln L & \frac{\partial^2}{\partial \theta_2^2} \ln L \end{bmatrix}$$

To show that $\nabla^2 \ln L(\theta_1, \theta_2)$ is positive definite we have
 $\frac{\partial^2}{\partial \theta_1^2} \ln L < 0$ and $\det[\nabla^2 \ln L(\theta_1, \theta_2)] > 0$ evaluated at $(\hat{\theta}_1, \hat{\theta}_2)$. Or find the eigen-values which both has to be negative.

Note: If we have a diagonal matrix, which is 2×2 , then to check whether it is negative definite, just see if the diagonal elements are negative as they are the eigenvalues. (7)

E.g.: Let us consider a random sample from a normal population, with mean μ and variance σ^2 :

$$\begin{aligned}\therefore L(x, \mu, \sigma^2) &= \prod_{i=1}^n f_{X_i}(x_i, \mu, \sigma^2) \\ &= \prod_{i=1}^n \frac{1}{\sigma(\sqrt{2\pi})} e^{-\frac{1}{2\sigma^2}(x_i-\mu)^2} \\ &= \left(\frac{1}{2\pi\sigma^2}\right)^{n/2} e^{-\frac{1}{2\sigma^2}\sum_{i=1}^n(x_i-\mu)^2}\end{aligned}$$

where $\sigma > 0$, $-\infty < \mu < +\infty$

$$\therefore \ln(L(x, \mu, \sigma^2)) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

\therefore The maximum likelihood estimation problem; i.e MLE problem is

$$\max_{\begin{array}{l} \sigma > 0 \\ -\infty < \mu < +\infty \end{array}} \ln(L(x, \mu, \sigma^2)) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\begin{aligned}\therefore \frac{\partial}{\partial \mu} \ln L(x, \mu, \sigma^2) &= 0 \quad \text{and} \quad \frac{\partial}{\partial \sigma^2} \ln L(x, \mu, \sigma^2) = 0 \\ \therefore \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) &= 0 \quad \text{and} \quad -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 = 0 \\ -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 &= 0\end{aligned}$$

$$\therefore \begin{cases} \hat{\mu} = \frac{1}{n} \sum x_i = \bar{x} \\ \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \end{cases} \rightarrow \text{The MLE estimators?}$$

To confirm that these are MLE estimators we have to establish the fact that these are maximizers. The simplest way to realize this is that $\ln L(x, \mu, \sigma^2)$ is a concave function of (μ, σ^2) and hence the critical point is a global maximizer. But this is beyond the scope of this course. Hence we check the usual way stated in the previous page.

The Hessian matrix $\nabla^2 \ln L(x, \mu, \sigma^2)$ evaluated at $(\hat{\mu}, \hat{\sigma}^2)$ is given as

$$\nabla^2 \ln L(x, \hat{\mu}, \hat{\sigma}^2) = \begin{bmatrix} -\frac{n}{\sigma^2} & 0 \\ 0 & \frac{1}{\sigma^4} \left[\frac{n}{2} - n^2 \right] \end{bmatrix}$$

Thus $-\frac{n}{\sigma^2} < 0$ & $\frac{1}{\sigma^4} \left[\frac{n}{2} - n^2 \right] < 0$, hence

$$\nabla^2 \ln L(x, \hat{\mu}, \hat{\sigma}^2)$$

to be negative definite. Thus $\hat{\mu}$ & $\hat{\sigma}^2$ are the maximizers of the mle problem.

Detail computation of $\nabla^2 \ln L(x, \hat{\mu}, \hat{\sigma}^2)$:

$$\frac{\partial}{\partial \mu} \ln L(x, \mu, \sigma^2) = -\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu)$$

$$\therefore \frac{\partial^2}{\partial \mu^2} \ln L(x, \hat{\mu}, \hat{\sigma}^2) = \cancel{\sum_{i=1}^n} -\frac{n}{\hat{\sigma}^2}$$

$$\begin{aligned} \frac{\partial}{\partial \sigma^2} \frac{\partial}{\partial \mu} \ln L(x, \hat{\mu}, \hat{\sigma}^2) &= \frac{\partial}{\partial \sigma^2} \ln L(x, \hat{\mu}, \hat{\sigma}^2) \\ &= -\hat{\sigma}^{-4} \sum_{i=1}^n (x_i - \hat{\mu}) \\ &= -\hat{\sigma}^{-4} \left[\sum_{i=1}^n (x_i - \bar{x}) \right] \\ &= \hat{\sigma}^{-4} \left(\sum_{i=1}^n x_i - n\bar{x} \right) = 0 \end{aligned}$$

$$\frac{\partial}{\partial \sigma^2} \frac{\partial \ln L(x, \mu, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2$$

$$\therefore \frac{\partial}{\partial \mu} \cdot \frac{\partial}{\partial \sigma^2} \ln L(x, \hat{\mu}, \hat{\sigma}^2) = \frac{1}{2\hat{\sigma}^2} 2 - \left(\sum_{i=1}^n (x_i - \hat{\mu}) \right) = 0$$

$$\frac{\partial}{\partial \sigma^2} \left(\frac{\partial}{\partial \sigma^2} \ln L(x, \hat{\mu}, \hat{\sigma}^2) \right) = \frac{n}{2} \hat{\sigma}^{-4} - n^2 \hat{\sigma}^{-4} = \frac{1}{\hat{\sigma}^4} \left[\frac{n}{2} - n^2 \right] < 0$$

(We use $\sum_{i=1}^n (x_i - \bar{x})^2 = n\hat{\sigma}^2$) ⑨

We would like to note that second order condition only tells us that the parameter estimators that we have got are really only local minimizer albeit strict one. To establish that they are global we need to use ^{Concavity}.

The second order check also establishes the concavity of functions here.

A crash course on Concavity

Let I be an interval in \mathbb{R} and $f: I \rightarrow \mathbb{R}$.

The function $f: I \rightarrow \mathbb{R}$ is called concave if for any $x, y \in I$ & $\lambda \in [0, 1]$

$$f(\lambda y + (1-\lambda)x) \geq \lambda f(y) + (1-\lambda)f(x). \quad \rightarrow (A)$$

Note: If $x \geq y$ are in I ; then for any $\lambda \in [0, 1]$, $\lambda y + (1-\lambda)x \in I$

• Thus if f is differentiable, then from (A), by Taylor's Theorem.

$$f(x) + \lambda f'(x)(y-x) + o(\lambda) \geq \lambda f(y) + (1-\lambda)f(x), \quad \forall \lambda \in (0, 1)$$

where $\frac{o(\lambda)}{\lambda} \rightarrow 0$ as $\lambda \rightarrow 0$.

$$\therefore \lambda f'(x)(y-x) + o(\lambda) \geq \lambda (f(y) - f(x))$$

$$\therefore \text{As } \lambda \rightarrow 0, \quad f'(x)(y-x) \geq f(y) - f(x)$$

If $f'(x) = 0 \Rightarrow f(x) \geq f(y) \Rightarrow x$ is a global maximum maximizer

Let $f: I \times I \rightarrow \mathbb{R}$, then f is concave if for any $(x_1, x_2) \in I \times I$ and $(y_1, y_2) \in I_1 \times I_2$ & $\lambda \in (0, 1)$

$$f(\lambda y_1 + (1-\lambda)x_1, \lambda y_2 + (1-\lambda)x_2) \geq \lambda f(y_1, y_2) + (1-\lambda)f(x_1, x_2)$$

Again using Taylor's Theorem in 2-dimensions we have $\forall (x_1, x_2) \in I \times I$

$$\frac{\partial f}{\partial x_1}(y_1 - x_1) + \frac{\partial f}{\partial x_2}(y_2 - x_2) \geq f(y_1, y_2) - f(x_1, x_2).$$

$$\therefore \frac{\partial f}{\partial x_1} = 0 \text{ & } \frac{\partial f}{\partial x_2} = 0, \quad f(x_1, x_2) \geq f(y_1, y_2), \quad \forall (y_1, y_2) \in I \times I$$

$\Rightarrow (x_1, x_2)$ is the global maximizer.

E.g.: $-\alpha(x-\mu)^2 = f(x),$
 $\alpha > 0, \mu \in \mathbb{R}$
is concave in x

$f(x) = -\ln x$
 $x > 0, i.e. x \in \mathbb{R}^+$
is concave in x .
(1m)

If we fix x , then
 f is concave on I if
 $f''(x) \leq 0, \forall x \in I$
 f is concave on $I \times I$ if
 $\nabla^2 f$ is negative definite on $I \times I$

In our study in the example the log likelihood functions can be shown to be concave.

Section 4 Fisher Information & Cramer-Rao Inequality

The maximum likelihood estimation technique, is very popular, and the reason for this is follows. Given a sample observation,

$x_1 = x_1, \dots, x_n = x_n$ the likelihood function $L(x, \theta)$

provides us the "likelihood" or the frequency of the occurrence of the observation x_1, x_2, \dots, x_n . Since the parameter θ is unknown the maximum Suppose we ask the question: What is the distribution for which the observation observations x_1, \dots, x_n occur the maximum number of time? This can be done by finding the θ , which maximizes the joint pdf, since it is the parameter θ , which specifies the distribution. We shall now see how the likelihood function can be used to build some important measures related to estimators.

We shall begin with the notion of Fisher Information, but let us first write down the two basic assumptions we need. We will ~~first~~ study only the single parameter case.

Assumptions

a) The pdf $f_X(x, \theta)$ is twice continuously differentiable as a function of the parameter θ .

3) We can differentiate ~~as~~ with respect to θ under the integral sign for the integral $\int f_X(x, \theta) dx$. In fact we will assume that we can differentiate twice under the integral sign with respect to θ .

$$\begin{aligned} & \frac{\partial}{\partial \theta} \int_A f_X(x, \theta) dx \\ &= \int_A \frac{\partial}{\partial \theta} f_X(x, \theta) dx \end{aligned}$$

$A \subseteq \mathbb{R}$. c)

The parameter θ belongs to an open interval.

The Fisher information function is built on the idea of a Fisher score function.

Let X_1, \dots, X_n be a random sample from a population $f_X(x, \theta)$.
 The Fisher score function for the i th observation is given as the r.v. $F(X_i, \theta)$
 whose value for the observation $x=x_i$ is given as $\bar{F}(x_i, \theta) = \frac{\partial}{\partial \theta} \ln f_X(x_i, \theta)$

$$\therefore F(x_i, \theta) = \frac{1}{f_X(x_i, \theta)} \cdot \frac{\partial}{\partial \theta} f_X(x_i, \theta)$$

$$\begin{aligned}\therefore E[F(x_i, \theta)] &= \int_{-\infty}^{\infty} \frac{1}{f_X(x_i, \theta)} \frac{\partial}{\partial \theta} f_X(x_i, \theta) f_X(x_i, \theta) dx \\ &= \int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} f_X(x_i, \theta) dx \\ &= \frac{\partial}{\partial \theta} \int_{-\infty}^{\infty} f_X(x_i, \theta) dx = \frac{\partial}{\partial \theta} (1) = 0.\end{aligned}$$

$$\therefore E[\bar{F}(x_i, \theta)] = 0$$

The total Fisher score function for a sample of size n is given as the r.v.

$$\bar{F}_n(\theta) = \sum_{i=1}^n \bar{F}(x_i, \theta)$$

$$\boxed{\bar{F}_n(\theta) = \sum_{i=1}^n \bar{F}(x_i, \theta)}$$

$$\left[\text{Note } E(\bar{F}_n(\theta)) = 0 \right]$$

The first Fisher information about the i -th observation is denoted as $I(\theta)$ and is given as

$$\boxed{I_i(\theta) = \text{Var}(F(x_i, \theta))}$$

$$\begin{aligned}
 I_i(\theta) &= \text{Var}(F(x_i, \theta)) \\
 &= \sqrt{E[(F(x_i, \theta))^2]} \quad (\because E(F(x_i, \theta)) = 0) \\
 &= \sqrt{E\left[\left(\frac{\partial}{\partial \theta} \ln f_x(x_i, \theta)\right)^2\right]} \\
 &= \int_{-\infty}^{\infty} \left[\frac{\partial}{\partial \theta} \ln f_x(x_i, \theta) \right]^2 f_x(x, \theta) dx \\
 &= \int_{-\infty}^{\infty} \frac{1}{f_x(x, \theta)} \cdot \frac{1}{f_x^2(x_i, \theta)} \left[\frac{\partial}{\partial \theta} f_x(x_i, \theta) \right]^2 f_x(x, \theta) dx \\
 &= \int_{-\infty}^{\infty} \frac{1}{f_x(x_i, \theta)} \left[\frac{\partial}{\partial \theta} f_x(x_i, \theta) \right]^2 dx
 \end{aligned}$$

Usually in texts one writes

$$I_i(\theta) = E\left[\left(\frac{\partial}{\partial \theta} \ln f_x(x_i, \theta)\right)^2\right]$$

For a random sample of size n , the Fisher information is given as.

E.g 1: Consider a random sample drawn from of size n drawn from an exponential population, $f_x(x, \theta)$, given as

$$\therefore f_x(x, \theta) = \theta e^{-\theta x}; \quad x > 0, \theta > 0$$

Thus

$$\begin{aligned}
 F(x_i, \theta) &= \frac{\partial}{\partial \theta} \ln(\theta e^{-\theta x_i}) \\
 &= \frac{1}{\theta e^{-\theta x_i}} \left[\frac{\partial}{\partial \theta} (\theta e^{-\theta x_i}) \right] \\
 &= \frac{1}{\theta e^{-\theta x_i}} \left[e^{-\theta x_i} + \theta e^{-\theta x_i} (-x_i) \right]
 \end{aligned}$$

(13)

$$\therefore \bar{F}(x_i, \theta) = \frac{1}{\theta} \left[1 - e^{-\theta x_i} \right]$$

$$= \frac{(1 - \theta x_i)}{\theta}$$

$$\therefore \bar{F}(x_i, \theta) = \left[\frac{1 - \theta x_i}{\theta} \right]$$

$$\therefore I(\theta) = E \left[(\bar{F}(x_i, \theta))^2 \right]$$

$$= \int_0^\infty \frac{(1 - \theta x_i)^2}{\theta^2} \cdot \theta e^{-\theta x_i} dx$$

$$= \int_0^\infty \frac{(1 - \theta x_i)^2}{\theta^2} \cdot \theta e^{-\theta x_i} dx$$

$$= \frac{1}{\theta} \int_0^\infty (1 - \theta x_i)^2 e^{-\theta x_i} dx$$

$$= \frac{1}{\theta} \int_0^\infty (1 - 2\theta x_i + \theta^2 x_i^2) \cdot e^{-\theta x_i} dx$$

$$= \frac{1}{\theta} \int_0^\infty e^{-\theta x_i} dx - \frac{2}{\theta} \int_0^\infty x_i \cdot \theta e^{-\theta x_i} dx + \cancel{\frac{1}{\theta} \int_0^\infty x_i^2 \cdot \theta e^{-\theta x_i} dx}$$

$$= \frac{1}{\theta^2} \int_0^\infty \theta e^{-\theta x_i} dx - \frac{2}{\theta} E(x_i) + \cancel{\frac{1}{\theta} E(x_i^2)}$$

$$= \frac{1}{\theta^2} - \frac{2}{\theta} \cdot \frac{1}{\theta} + \cancel{\frac{1}{\theta}} \left[\text{Var}(x_i) + (E(x_i))^2 \right]$$

$$= \frac{1}{\theta^2} - \frac{2}{\theta^2} + \cancel{\frac{1}{\theta}} \left[\frac{1}{\theta^2} + \frac{1}{\theta^2} \right]$$

$$= \frac{1}{\theta^2} - \frac{2}{\theta^2} + \cancel{\frac{1}{\theta^2}} \left[\frac{2}{\theta^2} \right] = -\frac{1}{\theta^2} + \frac{2}{\theta^4}$$

$$= \frac{2}{\theta^4} - \frac{1}{\theta^2} = \frac{1}{\theta^2} \left[\frac{2}{\theta^2} - 1 \right]$$

So the Fisher information for the whole sample is given as

$$I_n(\theta) = \text{Var}(F_n(\theta)) \\ = \text{Var}(\theta F(x_1, \theta) + \dots + F(x_n, \theta))$$

Since x_1, \dots, x_n are independent, $F(x_1, \theta), \dots, F(x_n, \theta)$ are independent. Thus

$$I_n(\theta) = \text{Var}(F_{\theta}(x_1, \theta)) + \dots + \text{Var}(F(x_n, \theta)) \\ = \sum I_i(\theta)$$

$$\therefore I_n(\theta) = \sum_{i=1}^n I_i(\theta)$$

But each x_i 's are having the same distribution.

$$\boxed{I_n(\theta) = n I_i(\theta)}$$

The Cramer-Rao Inequality

A random sampling is called regular if its Fisher information is continuous, strictly positive and bounded for all θ in the given range.

For example if we consider a random sample from x_1, \dots, x_n from an exponential distribution, then

$$I_n(\theta) = n I_i(\theta) \\ = \frac{n}{\theta^2} > 0$$

$I_n(\theta)$ is of course continuous on $\theta > 0$. However, $I_n(\theta)$ is not bounded above, as a function of θ . But as $I_n(\theta)$ is continuous and positive we can consider the sampling from an exponential distribution is ^{not} regular.

Theorem II.1: Cramer-Rao Bound (Actually a lower-bound)

Let $T = \hat{\theta}(x_1, \dots, x_n)$ be a statistic which is an estimator of θ of the population $f(x, \theta)$. Let the bias $b_n(\theta) = E[\hat{\theta}] - \theta$ be continuously differentiable. Then

$$\text{Var}(\hat{\theta}) \geq \frac{(1 + b'_n(\theta))^2}{I_n(\theta)}$$

Proof:

$$\theta + b_n(\theta) = E[\hat{\theta}] = \int_{\mathbb{R}^n} \hat{\theta}_n(x_1, \dots, x_n) f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta) dx_1 \dots dx_n$$

Now differentiation under the integral sign we have.

$$\begin{aligned} 1 + b'_n(\theta) &= \int_{\mathbb{R}^n} \hat{\theta}_n(x_1, \dots, x_n) \frac{\partial f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta)}{\partial \theta} dx_1 \dots dx_n \\ &= \int_{\mathbb{R}^n} \hat{\theta}(x_1, \dots, x_n) \frac{\frac{\partial f_x(x_1, \dots, x_n, \theta)}{\partial x_i}}{f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta)} \frac{\partial}{\partial \theta} f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta) dx_1 \dots dx_n \\ &= \int_{\mathbb{R}^n} \hat{\theta}(x_1, \dots, x_n) F_n(\theta) f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta) dx_1 \dots dx_n \end{aligned}$$

$$\begin{aligned} F_n(\theta) &= \sum_{i=1}^n F(x_i, \theta) = \sum \frac{1}{f_x(x_i, \theta)} \frac{\partial}{\partial \theta} f_x(x_i, \theta) \\ &= \sum_{i=1}^n \frac{1}{f_{x_i}(x_i, \theta)} \frac{\partial}{\partial \theta} f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta) \\ &= \sum_{i=1}^n \frac{\partial}{\partial \theta} \ln f_x(x_i, \theta) \\ &= \frac{\partial}{\partial \theta} \left[\sum_{i=1}^n \ln f_{x_i}(x_i, \theta) \right] \\ &= \frac{\partial}{\partial \theta} \left[\ln \frac{f_x(x_1, \theta)}{\prod_{i=1}^n f_{x_i}(x_i, \theta)} \right] = \frac{\partial}{\partial \theta} \ln f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta) \end{aligned}$$

(By independence). ↓

(16)

$$\therefore F_n(\theta) = \frac{1}{\int_{x_1 \dots x_n} f_{x_1 \dots x_n}(x_1 \dots x_n, \theta)} \frac{\partial}{\partial \theta} f_{x_1 \dots x_n}(x_1 \dots x_n, \theta)$$

$$\therefore 1 + b'_n(\theta) = E[\hat{\theta} F_n(\theta)] = \text{cov}(\hat{\theta}, F_n(\theta))$$

$$\text{Note that } \text{cov}(\hat{\theta}, F_n(\theta)) = E[\hat{\theta} F_n(\theta)] = -E[\hat{\theta}] E[F_n(\theta)]$$

$$= E[\hat{\theta} F_n(\theta)] \quad (\because E[F_n(\theta)] = 0 \text{ as desired earlier})$$

$\therefore \text{let } P_n^2 = \frac{\text{square of correlation coefficient of } \hat{\theta}, F_n(\theta) \text{ and } F_n(\theta)}{1}$

$$\therefore P_n^2 = \frac{[\text{cov}(\hat{\theta}, F_n(\theta))]^2}{\text{Var}(\hat{\theta}) \text{Var}(F_n(\theta))}$$

But $P_n^2 \leq 1$ as $P_n \in [-1, +1]$

$$\Rightarrow \frac{[\text{cov}(\hat{\theta}, F_n(\theta))]^2}{\text{Var}(\hat{\theta}) \text{Var}(F_n(\theta))} \leq 1.$$

Hence $\text{Var}(\hat{\theta}) \geq \frac{[\text{cov}(\hat{\theta}, F_n(\theta))]^2}{\text{Var}(F_n(\theta))}$

$$\Rightarrow \boxed{\text{Var}(\hat{\theta}) \geq \frac{(1 + b'_n(\theta))^2}{I_n(\theta)}} \quad \square$$

This is called the Cramer-Rao Lower bound. Now if $\hat{\theta}$ is an unbiased estimator, then $b'_n(\theta) = 0 \Rightarrow$ for any θ in the given range. Thus we have

$$\boxed{\text{Var}(\hat{\theta}) \geq \frac{1}{I_n(\theta)}}$$

This is called the Cramer-Rao Inequality. An unbiased estimator $\hat{\theta}$ is called efficient if

$$\boxed{\text{Var}(\hat{\theta}) = \frac{1}{I_n(\theta)}}$$

Let us finish our discussion with an example by drawing a sample of size n from a normal population $N(\mu, \sigma^2)$, where μ is unknown but σ^2 is ~~not~~ known.

We shall show that $\bar{X}_n = \bar{X} = \frac{x_1 + \dots + x_n}{n}$ is indeed an efficient

estimator of μ . \bar{X}_n is known to be unbiased as we have

already shown. Now $\text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$. We have

$$f_x(x, \mu) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}},$$

$$\therefore \ln f_{x_i}(x_i, \mu) = -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{(x_i - \mu)^2}{2\sigma^2}$$

In the random-variable form we have

$$\ln f_{x_i}(x_i, \mu) = -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{(x_i - \mu)^2}{2\sigma^2}$$

$$F(x_i, \theta) = \frac{x_i - \theta\mu}{\sigma^2}$$

$$I_n(\theta) = n I_i(\theta) = n \text{Var}(F(x_i, \mu)) = n E[(F(x_i, \mu))^2]$$

$$= n E[(F(x_i, \mu))^2]$$

$$= n E\left[\frac{(x_i - \mu)^2}{\sigma^4}\right]$$

$$= \frac{n}{\sigma^4} E[(x_i - \mu)^2]$$

$$= \frac{n}{\sigma^4} \cdot \text{Var}(x_i)$$

$$= \frac{n}{\sigma^4} \cdot \sigma^2 = \frac{n}{\sigma^2}$$

$$\text{Now } \text{Var}(\bar{X}_n) = \frac{\sigma^2}{n} = \frac{1}{I_n(\theta)}, \text{ proving that}$$

\bar{X}_n is indeed an efficient estimator.

The part on Fisher Information is partially based on Chapter 1, of the book titled: Mathematical Statistics: Asymptotic minimax theory. by A. Korostelev and O. Korosteleva, AMS-2011.

If we now consider $\hat{\sigma}^2 = \hat{\theta}(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ as the maximum likelihood estimator, or $b_n(\sigma^2)$ in the normal distribution, then observe that

$$b_n(\hat{\sigma}^2) = -\frac{\hat{\sigma}^2}{n}$$

Thus $\frac{d b_n(\sigma^2)}{d\sigma^2} = -\frac{1}{n} = b'_n(\sigma^2)$

$$\therefore (1 + b'_n(\sigma^2)) = -\frac{1}{n}. \text{ Now we shall compute}$$

$$F(x_i, \sigma^2) = \frac{\partial}{\partial \sigma^2} \ln f_X(x_i, \sigma^2)$$

$$= \frac{\partial}{\partial \sigma^2} \left[-\frac{1}{2} \ln(2\pi\sigma^2) - \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$= \frac{\partial}{\partial \sigma^2} \left[-\frac{1}{2} \ln 2\pi - \frac{1}{2} \ln \sigma^2 - \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$= -\frac{1}{2} \cdot \frac{\partial}{\partial \sigma^2} \ln \sigma^2 - \frac{\partial}{\partial \sigma^2} \frac{(x_i - \mu)^2}{2\sigma^2}$$

$$= -\frac{1}{2} \frac{\partial}{\partial \sigma^2} \ln \sigma^2 - \frac{1}{2} \frac{\partial}{\partial \sigma^2} \cdot \frac{(x_i - \mu)^2}{2\sigma^2}$$

$$= -\frac{1}{2} \left[\frac{\partial}{\partial \sigma^2} \ln \sigma^2 + \frac{(x_i - \mu)^2}{2\sigma^2} \frac{\partial}{\partial \sigma^2} \left(\frac{1}{\sigma^2} \right) \right]$$

$$= -\frac{1}{2} \left[\frac{1}{\sigma^2} - \frac{(x_i - \mu)^2}{\sigma^4} \right]$$

$$= \frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2} + \frac{1}{2\sigma^4}$$

Now Fisher information for the i -th sample

$$\begin{aligned} I_i(\theta) &= \text{Var}(F_i(x, \theta)) \\ &= \text{Var}\left(\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2} + \frac{1}{2\sigma^4}\right) \\ &= \text{Var}\left(\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}\right) \\ &= \frac{1}{4} \text{Var}\left[\frac{(x_i - \mu)^2}{\sigma^2}\right] \\ &= \frac{1}{4} \text{Var}\left[\left(\frac{x_i - \mu}{\sigma}\right)^2\right] \\ &= \frac{1}{4} \text{Var}(z_i^2) \end{aligned}$$

Now $z_i^2 \sim \chi^2$ with one degree of freedom.

$$\text{Now } \text{Var}(z_i^2) = 2$$

$$\therefore I_i(\theta) = \frac{1}{4} \cdot 2 = \frac{1}{2}.$$

$$\text{Thus } I_n(\theta) = \frac{n}{2}$$

Hence Cramer-Rao lower-bound tells us that,

$$\boxed{\text{Var}(\hat{\sigma}^2) \geq \frac{\left(1 - \frac{1}{n}\right)^2}{\frac{n}{2}}}$$

Lecture -12 : Interval Estimation

Section 1: Confidence Intervals

What we have learned till now is point estimation. In order to estimate the value of a population parameter, we just draw a random sample of a reasonable size and construct a statistic which can best estimate that population parameter, in some sense, either as unbiased or as the best maximum likelihood estimator. But if we are drawing a sample from a probability density function the probability that a value of the statistic will coincide with population parameter is zero. This motivates us to take the view that it might be better to be able to construct an interval around the statistic so that the population parameter value falls in that interval with certain probability.

Let us draw a sample from the population with pdf $f_X(x, \theta)$, and let x_1, \dots, x_n be a random sample of size n drawn from it, and let $\hat{\Theta}$ be a statistic, $\hat{\Theta} = \hat{\Theta}(x_1, \dots, x_n)$ constructed to estimate, θ . We may for example want to know, if

$$\hat{\Theta} - c < \theta < \hat{\Theta} + c$$

with probability say γ ; i.e.

$$P(\hat{\Theta} - c < \theta < \hat{\Theta} + c) = \gamma$$

Then $(\hat{\Theta} - c, \hat{\Theta} + c)$ forms a confidence interval associated with θ . This is more formally defined as follows.

Definition 1: If $\hat{\Theta}_1$ and $\hat{\Theta}_2$ are values of the random variables $\hat{\Theta}_1, \hat{\Theta}_2$, such that which are themselves also statistics, such that,

$$P(\hat{\Theta}_1 < \theta < \hat{\Theta}_2) P(\hat{\Theta}_1 < \theta < \hat{\Theta}_2) = 1 - \alpha$$

for some specified probability $1 - \alpha$. We shall refer to the interval $(\hat{\Theta}_1, \hat{\Theta}_2)$ as a $(1 - \alpha) 100\%$ confidence interval for the population parameter θ . The value $1 - \alpha$ is called degree of confidence & $\hat{\Theta}_1$ & $\hat{\Theta}_2$ the upper and lower confidence limits.

limits. So if $\alpha = 0.05$, we say the degree of confidence is 0.95 and we get a 95% confidence interval. It is also important to note at the very beginning that the confidence intervals for θ need not be unique as we will just see.

Section 2: Interval Estimation of means

In this section we will study first the interval estimation of the mean μ of a normal population whose variance σ^2 is known. Next we will see how to deal the case when σ^2 is unknown. It will be important to note how $\hat{\theta}_1$ and $\hat{\theta}_2$ will be very different in these two scenarios.

Case 1: σ^2 is known.

Let x_1, \dots, x_n be a random sample of size n drawn from a population $N(\mu, \sigma^2)$, where σ^2 is known but μ is not. In order to estimate the $(1-\alpha)100\%$ confidence interval of μ we proceed as follows.

$$\text{Note that } E(\bar{x}) = \mu \quad \text{and} \quad \text{Var}(\bar{x}) = \frac{\sigma^2}{n}$$

$$\text{where } \bar{x} = \frac{x_1 + \dots + x_n}{n}.$$

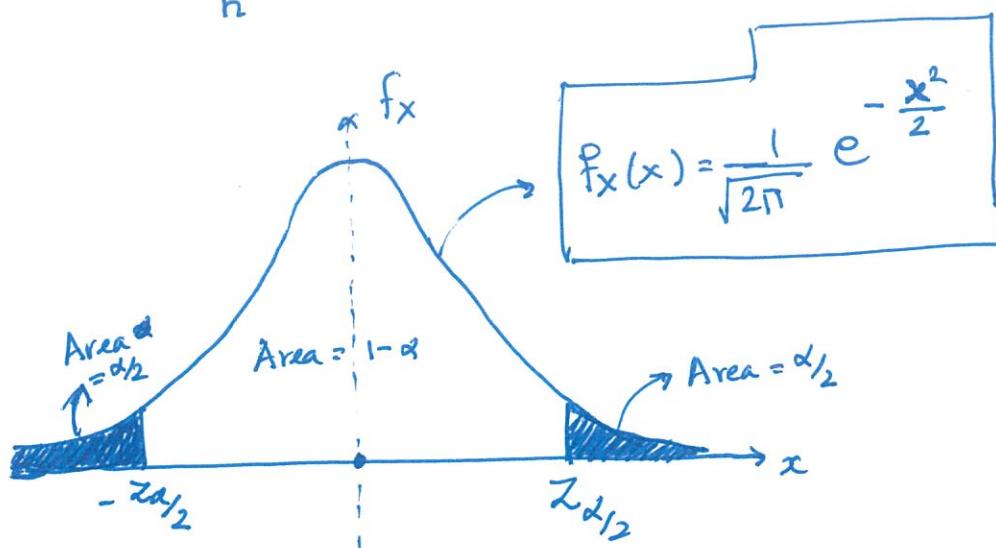


Fig 1: The point $z_{\alpha/2}$

Let us look at diagram above and it is seen how the point $Z_{\alpha/2}$ is defined; i.e.

$$f_{Z_{\alpha/2}} \int_{Z_{\alpha/2}}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = \alpha/2$$

\therefore From diagram if $Z \sim N(0,1)$ then

$$P(-Z_{\alpha/2} < Z < Z_{\alpha/2}) = 1 - \alpha \rightarrow \textcircled{4}$$

\therefore Since we have drawn a normal population

$$\bar{X} \sim N(\mu, \frac{\sigma^2}{n}) \quad [\text{Check Chapter Lecture 10}]$$

$$\therefore \text{Set } \Sigma = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0,1)$$

Hence from $\textcircled{4}$ we have

$$P\left(-Z_{\alpha/2} < \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} < Z_{\alpha/2}\right) = 1 - \alpha$$

$$P\left(\bar{X} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

$$\text{Here } \hat{\theta}_1 = \bar{X} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \quad \& \quad \hat{\theta}_2 = \bar{X} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}.$$

Now we will again consider studying a normal population but this time we shall consider that the population variance σ^2 is unknown. In this situation we have to use sample variance s^2 in order to try to develop a statistic using it. Let us set

$$T = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}}$$

Let us assume that the sample is small, i.e. $n < 30$, then it better to proceed by noting that $T \sim t$ distribution with $n-1$ degrees of freedom (see Lecture 10). To find out the $(1-\alpha)100\%$ confidence interval we must have two points $q_{1-\alpha/2}$, such that

$$P(q_{1-\alpha/2} < T < q_{\alpha/2}) = 1 - \alpha$$

$$\Rightarrow P\left(q_{1-\alpha/2} < \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} < q_{\alpha/2}\right) = 1 - \alpha$$

$$\Rightarrow P\left(\frac{s}{\sqrt{n}} q_{1-\alpha/2} < \bar{X} - \mu < \frac{s}{\sqrt{n}} q_{\alpha/2}\right) = 1 - \alpha$$

$$\Rightarrow P\left(\bar{X} - q_{\alpha/2} \frac{s}{\sqrt{n}} < \mu < \bar{X} - q_{1-\alpha/2} \frac{s}{\sqrt{n}}\right) = 1 - \alpha$$

Thus in this case the confidence interval for the mean is given by

$$\left(\bar{X} - q_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{X} - q_{1-\alpha/2} \frac{s}{\sqrt{n}}\right)$$

Hence the length of this interval = $(q_{\alpha/2} - q_{1-\alpha/2}) \frac{s}{\sqrt{n}}$

$$\text{or } L = (q_{\alpha/2} - q_{1-\alpha/2}) \frac{s}{\sqrt{n}}.$$

From a more practical point of view we want $(q_{l_2} - q_{l_1})$ to be as small as possible. However note that L is in general a random variable. However for a given sample size of size say n ; we need to actually solve the optimization problem

$$\min L = (q_{l_2} - q_{l_1}) \frac{S}{\sqrt{n}} \longrightarrow (E)$$

Subject to

$$\int_{q_{l_1}}^{q_{l_2}} f_T(q_l) dq_l = 1 - \alpha \longrightarrow (F).$$

where f_T is the pdf of the t-distribution with $(n-1)$ degrees of freedom. Let us rewrite (F) using the dummy variable t , i.e.

$$\int_{q_{l_1}}^{q_{l_2}} f_T(t) dt = 1 - \alpha. \longrightarrow (F')$$

Since α is fixed from a theoretical point of view one can always express q_{l_2} as a function of q_{l_1} . Let us first differentiate

(F') with respect to q_{l_1} . In order to do so we need to use the Leibnitz formula for differentiating under the integral sign. We shall write it in the following box.

Leibnitz rule : For differentiating under the integral sign.

Let $a(x)$ and $b(x)$ are differentiable function. Then consider the integral

$$\begin{aligned} & \int_{a(x)}^{b(x)} \varphi(x, t) dt \\ \therefore \frac{d}{dx} \int_{a(x)}^{b(x)} \varphi(x, t) dt &= \int_{a(x)}^{b(x)} \frac{\partial}{\partial x} \varphi(x, t) dt + \cancel{\int_{a(x)}^{b(x)} \varphi(x, b(x)) \frac{db}{dx}} - \varphi(x, a(x)) \frac{da}{dx} \end{aligned}$$

Note that if we go back to (F') , f_T depends on t and not on q_1 . Thus we will have

$$\frac{d}{dq_1} \int_{q_1}^{q_2} f_T(t) dt = 0$$

By Leibniz rule we

$$\int_{q_1}^{q_2} \left(\frac{\partial}{\partial q_1} f_T(t) dt + f_T(q_2) \frac{dq_2}{dq_1} - f_T(q_1) \frac{dq_1}{dq_1} \right) = 0$$

\Rightarrow As $\frac{\partial}{\partial q_1} f_T(t) = 0$ we have from the above expression

$$f_T(q_2) \frac{dq_2}{dq_1} - f_T(q_1) = 0 \longrightarrow (G)$$

Now when we replace q_2 as a function of q_1 in L then the problem become unconstrained and to minimize L we first differentiate with respect to q_1

$$\therefore \left(\frac{dq_2}{dq_1} - 1 \right) \frac{s}{\sqrt{n}} = 0$$

$$\Rightarrow \left(\frac{f_T(q_1)}{f_T(q_2)} - 1 \right) \frac{s}{\sqrt{n}} = 0$$

$$\text{or } f_T(q_1) = f_T(q_2)$$

Since s_0 , q_1 is the minimizer then $f_T(q_1) = f_T(q_2)$.

Of course one solution is $q_1 = q_2$, but that will give $1-\alpha = 0$ from (F') , which is a contradiction.

But as the graph of t distribution is symmetric we have $q_1 = -q_2$ as the solution. So from the statistical tables we have to find q_2 such that

$$P(T > q_2) = \frac{\alpha}{2}, \text{ where } T \sim t_{(n-1 \text{ degrees of freedom})}$$

In general one writes $q_2 = t_{\alpha/2, n-1}$.

What happens if we have a large sample. Let us draw a sample from a normal distribution. Then using what is known as "central limit theorem", we can deduce that if

$$Z_n = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}}$$

where n is the sample size, then as $n \rightarrow \infty$, the distribution function of Z_n converges to that of the standard normal distribution $N(0,1)$. Thus for n very large $Z_n \sim N(0,1)$ "approximately". So I leave it to the reader to deduce that the confidence $(1-\alpha)100\%$ confidence interval of μ is given as

$$\left(\bar{x} - z_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{s}{\sqrt{n}} \right)$$

Of course the next question one may ask is what about the variance. If we are studying a normal population with mean μ and variance σ^2 , what is the confidence interval for σ^2 , assuming of course σ^2 is not known. This is what we will study in the next section.

3. Interval Estimation for the Variance

As before let us draw a sample x_1, \dots, x_n from a normal distribution with mean μ and variance σ^2 . Here the variance is unknown. In order to proceed we will introduce what is known as a pivotal quantity.

Defn 3.1 (Pivotal quantity)

Let x_1, \dots, x_n be a random sample of size n drawn from a density $f(\cdot, \theta)$, where θ as before is the parameter of the distribution. Let $Q = q_r(x_1, \dots, x_n, \theta)$, i.e. Q is a function of the sample observations and θ . If Q is a probability density that does not depend on θ , then Q is called a pivotal quantity.

E.g.: In the examples of estimating mean $Z_n = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$ is a pivotal quantity.

So the key idea of estimation is to find $q_{l_1} < q_{l_2}$ such that

$$P[q_{l_1} < Q < q_{l_2}] = 1 - \alpha.$$

In the particular case of estimating variance, the pivotal quantity is

$$Q = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sigma^2} = \frac{(n-1)S^2}{\sigma^2}$$

Now $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2$ with $(n-1)$ degrees of freedom.

$$\therefore P[q_{l_1} < Q < q_{l_2}] = 1 - \alpha$$

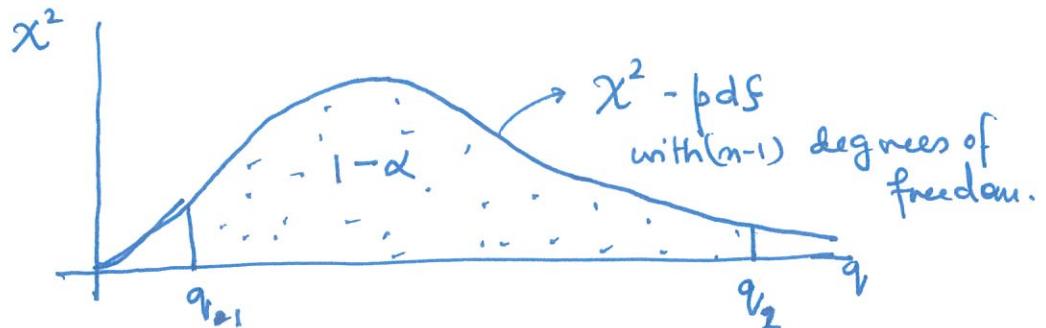
$$\Rightarrow P\left[q_{l_1} < \frac{(n-1)S^2}{\sigma^2} < q_{l_2}\right] = 1 - \alpha \quad (\text{since } S^2 \text{ is the observed value of the r.v. } S^2)$$

$$P \left[\frac{(n-1)s^2}{q_{\alpha/2}} < \sigma^2 < \frac{(n-1)s^2}{q_{1-\alpha/2}} \right] = 1 - \alpha$$

\therefore The $(1-\alpha)100\%$ confidence intervals are then given as

$$\left(\frac{(n-1)s^2}{q_{\alpha/2}}, \frac{(n-1)s^2}{q_{1-\alpha/2}} \right)$$

Remember that χ^2 with $(n-1)$ degrees of freedom is not a symmetric distribution



The key idea is to find $q_{\alpha/2}$ and $q_{1-\alpha/2}$ such that

$$P[Q < q_{\alpha/2}] = \alpha/2 \quad \& \quad P[Q > q_{1-\alpha/2}] = \alpha/2.$$

This is equal called equal tails estimation. In practice one should figure this out from the statistical tables related to χ^2 -distribution but as before one may also obtain this by numerically solving the optimization problem.

$$\min_{\sigma^2} \min L = \frac{(n-1)s^2}{\sigma^2} \left[\frac{1}{q_{\alpha/2}} - \frac{1}{q_{1-\alpha/2}} \right]$$

Subject to $\int_{q_{\alpha/2}}^{q_{1-\alpha/2}} f_Q(q) dq = 1 - \alpha.$

⑨

4. Estimation of the Confidence Interval for the difference of the means of two population

In this case we will be drawing samples from two normal populations.

Sample 1: will be drawn from $N(\mu_1, \sigma_1^2)$

Sample 2: will be drawn from $N(\mu_2, \sigma_2^2)$

We will assume that σ_1^2 and σ_2^2 are known. We want to estimate the difference of the means μ_1 and μ_2 , i.e. $\mu_1 - \mu_2$ to know how different these two normal populations are. Size of sample 1 is n_1 .

The pivotal quantity in this case is

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

where \bar{x}_1 is the mean of the first sample and \bar{x}_2 is the mean of the second sample. One can show that

$$Z \sim N(0, 1)$$

Observe that

$$E(\bar{x}_1 - \bar{x}_2) = E(\bar{x}_1) - E(\bar{x}_2) = \mu_1 - \mu_2$$

$$\text{Var}(\bar{x}_1 - \bar{x}_2) = \text{Var}(\bar{x}_1) + \text{Var}(\bar{x}_2) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$$

As, \bar{x}_1 & \bar{x}_2 are independent

It can be proved using any approach of finding the sampling distribution of $\bar{x}_1 - \bar{x}_2$, that

$$\bar{x}_1 - \bar{x}_2 \sim N\left(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right).$$

Thus $Z \sim N(0, 1)$.

Hence if we are looking for a $(1-\alpha)100\%$ interval confidence interval then

$$(\bar{x}_1 - \bar{x}_2) - Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} < \mu_1 - \mu_2 < (\bar{x}_1 - \bar{x}_2) + Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

If the samples are large and $\sigma_1^2 + \sigma_2^2$ are not known then replace them by s_1^2 and s_2^2 in the expression of Z and proceed using the central limit theorem.

In that case we have

$$(\bar{x}_1 - \bar{x}_2) - Z_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} < \mu_1 - \mu_2 < (\bar{x}_1 - \bar{x}_2) + Z_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

However if the sample size is small i.e. ($n_1 < 30$ or $n_2 < 30$ or both), then the confidence interval can be used if both population is assumed to have the same variance $\sigma^2 = \sigma_1^2 = \sigma_2^2$ even though its exact value is not known. Then the key is to consider the pooled sample variance

$$s_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}$$

This is nothing but an weighted average of s_1^2 & s_2^2 ; i.e

$$s_p^2 = \frac{\sum_{i=1}^{n_1} (x_i^1 - \bar{x}_1)^2 + \sum_{i=1}^{n_2} (x_i^2 - \bar{x}_2)^2}{n_1 + n_2 - 2}$$

You will show in the assignments that $E(s_p^2) = \sigma^2$

Construct a new random variable

$$Y = \frac{(n_1-1)s_1^2}{\sigma^2} + \frac{(n_2-1)s_2^2}{\sigma^2} = \frac{(n_1+n_2-2)s_p^2}{\sigma^2}$$

Now $\frac{(n_1-1)}{\sigma^2} s_1^2 \sim \chi^2$ with n_1-1 degrees of freedom

and $\frac{(n_2-1)}{\sigma^2} s_2^2 \sim \chi^2$ with n_2-1 degrees of freedom.

Thus $Y \sim \chi^2$ with $n_1 + n_2 - 1$ degrees of freedom. Now construct the pivotal variable

$$T = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Of course it can be shown that $\bar{x}_1 - \bar{x}_2 \sim t$ follows the t distribution with $n_1 + n_2 - 2$ degrees of freedom. Once this is known we can go ahead and write down the $(1-\alpha)100\%$ confidence intervals. But how did we get T .

Note that let us set

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

of course $Z \sim N(0, 1)$. It can also be shown that Z and Y are independent. (Note that we have seen one can show that \bar{x} and S^2 are independent).

$$\therefore T = \frac{Z}{\sqrt{\frac{Y}{n_1 + n_2 - 2}}}$$

Hence $T \sim t$ with $(n_1 + n_2 - 2)$ degrees of freedom. This is known from our study of Sampling distributions. The fact that $T \sim t$ with $(n_1 + n_2 - 2)$ degrees of freedom precisely needs the fact that Z and Y are independent

$$\therefore T = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Hence that $(1-\alpha)100\%$ the confidence interval is given as

$$(\bar{x}_1 - \bar{x}_2) - t_{\alpha/2, n-1} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} < \mu_1 - \mu_2 < (\bar{x}_1 - \bar{x}_2) + t_{\alpha/2, n-1} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

Lecture 13 : Hypothesis Testing: (Classical Theory)

Sect: The Basics

In many situations one is called upon to take a decision about a statistical statement, either quantitative or qualitative, in nature about parameters of some distribution representing the populations at hand. The decision will be to either accept or reject that statement; or hypothesis. Such a statement is usually called the null-hypothesis denoted as H_0 . Associated with a null hypothesis is an alternative hypothesis denoted as H_1 , which we accept if we reject H_0 .

A simple way to think as follows. Let a population is described by a pdf $f_x(\cdot, \theta)$, and we are to test the hypothesis (null-hypo.)

$$H_0: \theta = \theta_0$$

Against the alternative (say).

$$H_1: \theta > \theta_0 \text{ (or say } \theta \neq \theta_0\text{)}$$

The idea is to generate a random sample from the population $f_x(\cdot, \theta)$ and based on some statistic which an estimator of θ or some related function of the observation to make an informed decision about whether to accept or reject H_0 . Such a situation appears pretty often during clinical trials of drugs. Let a company X makes a new drug β for a disease for which a drug α exists and for which the average recovery rate of patients be μ_α . Let However the company X, claims that the new drug β , β is better for which it has forecasted an average recovery rate of μ_β . If the drug β be really good better than α , then

we must have to demonstrate that $\mu_p > \mu_a$, far enough to sample (c.e. patient under trial) observations. However in medicine trials this des decision problem as the following hypothesis testing problem:

$$H_0: \mu_a = \mu_p \quad (\text{Null hypothesis})$$

is tested against alternative

$$H_1: \mu_p > \mu_a.$$

Usually there are several ways to represent a hypothesis testing problem.

One is called the simple representation, like

$$H_0: \theta = \theta_0, \quad H_1: \theta = \theta_1.$$

While one can talk of a compound H_0 representation.

There are two types of error associated with it

Type-I error : Rejecting H_0 when it is true

Type-II error : Accepting H_0 when it is false,
and H_1 is true.

Note: When you reject H_0 you accept H_1 .

The basic structure of the testing procedure is as follows. Let the population under test is described by the pdf $f_X(\cdot, \theta)$, where $\theta \in \Omega$. Consider the hypothesis testing problem

$$H_0: \theta_0 \in A \quad \text{against} \quad H_1: \theta_1 \in A^c.$$

Note that $A \cup A^c = \Omega$.

Let \mathcal{D} be the space of all sample observations. The test of H_0 against H_1 is based on subset C of \mathcal{D} called the critical region of the test. It simply means one rejects H_0 if the sample observation falls in C .

The probability of a type-I error is given as follows

The rejection of H_0 , given that the observed values (x_1, x_2, \dots, x_n) lies in the critical region is called a non-randomized test.

Given a test procedure T , we shall denote henceforth the critical region. To begin with let us define the power function of a test.

Definition 13.1 Let T be a test of the null hypothesis

The power function $\Pi_T(\theta)$ of T is defined to be the probability of rejecting H_0 , when the distribution from which the sample is obtained is parametrized θ .

Thus

$$\begin{aligned}\Pi_T(\theta) &= P_{\theta} [\text{Reject } H_0] \\ &= P_{\theta} [(x_1, \dots, x_n) \in C_T]\end{aligned}$$

Example 1. Let us draw a sample from a normal population, with mean $\mu = \theta$, (unknown) and $\sigma^2 = 25$.

Consider the following hypothesis problem

$$H_0: \theta \leq 17, \text{ against } H_1: \theta > 17.$$

The test T is as follows: Reject if and only if $\bar{x} > 17 + \frac{5}{\sqrt{n}}$

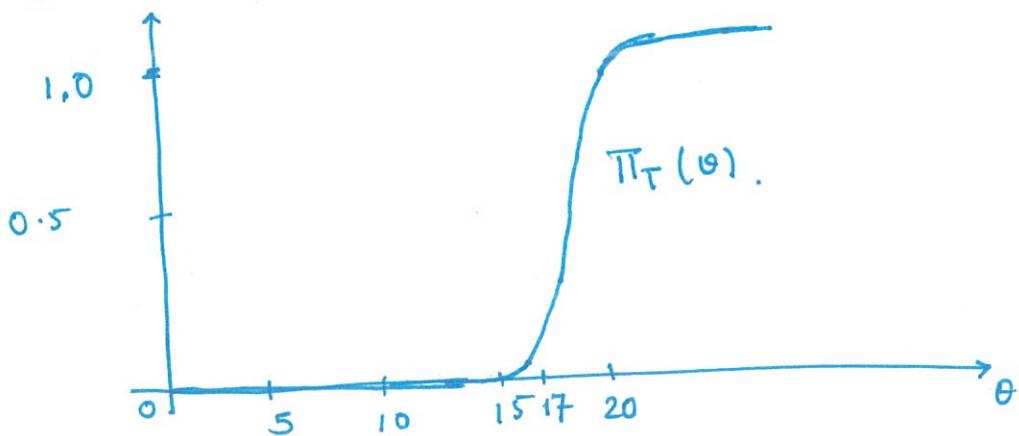
where n is the sample size. Here

$$C_T = \left\{ (x_1, \dots, x_n) : \bar{x} > 17 + \frac{5}{\sqrt{n}} \right\}$$

$$\begin{aligned}\Pi_{H_0 T}(\theta) &= P_{\theta} \left[\bar{x} > 17 + \frac{5}{\sqrt{n}} \right] \\ &= P_{\theta} \left[\frac{\bar{x} - \theta}{5/\sqrt{n}} > \frac{17 + \frac{5}{\sqrt{n}} - \theta}{5/\sqrt{n}} \right] \\ &= 1 - \Phi \left(\frac{17 + \frac{5}{\sqrt{n}} - \theta}{\frac{5}{\sqrt{n}}} \right)\end{aligned}$$

In this scenario where the hypothesis is composite, it is better to plot the graph of $\pi_T(\theta)$, to understand how effective is this particular test. We present is the graph from the book "Introduction to the Theory of Statistics" by Mood, Graybill and Boes, from where the above example was taken.

They considered $n = 25$ and we have



So let us see what do we observe. When $\theta \leq 16$, it is clear we accept H_0 . Of course if $\theta \geq 19$ say we clearly reject H_0 as $\pi_T(\theta) > 0.5$. But if $17 < \theta < 18$, we may also accept H_0 or reject H_0 as $\pi_T(\theta)$ is near 0.5 in some cases but may accept for $\theta > 17$ but very near it. Following the approach of Mood, Graybill and Boes let us define the notion of the size of a test.

The Size of a Test: Let the null hypothesis be given as

$H_0 : \theta \in \Omega_0$, where $\Omega_0 \subseteq \Omega$, where Ω is called the parameter space. The size of a test T associated with H_0 is defined as follows

$$\begin{aligned} \text{Size of Test} &= \sup_{\theta \in \Omega_0} \pi_T(\theta) \\ &= \sup_{\theta \in \Omega_0} P_\theta [\text{Reject } H_0] \end{aligned}$$

- For a non-randomized test the size of test is also referred as size of the critical region.

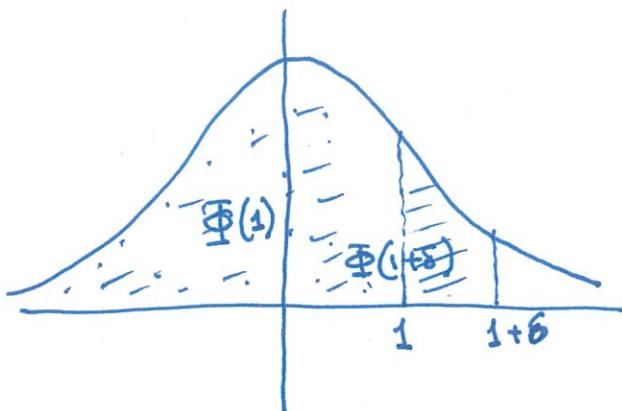
Let us try to find out the size of the test for the previous example. Just to recall, let us state that the random sample x_1, \dots, x_n be a random sample of size n drawn from a normal population from an μ with an unknown mean μ and variance $\sigma^2 = 25$. Our test T ; if we recall was as follows.

Test T : Reject H_0 if $\bar{x} > 17 + \frac{5}{\sqrt{n}}$

Here $\Omega_0 = \{\theta : \mu = \theta, \theta \leq 17\}$, (Remember we are talking about the mean μ)

$$\pi_T(\theta) = 1 - \Phi\left(\frac{17 + \frac{5}{\sqrt{n}} - \theta}{\frac{5}{\sqrt{n}}}\right)$$

$$\begin{aligned} \text{Size of the test} &:= \sup_{\theta \leq 17} \left[1 - \Phi\left(\frac{17 + \frac{5}{\sqrt{n}} - \theta}{\frac{5}{\sqrt{n}}}\right) \right] \\ &= 1 + \sup_{\theta \leq 17} \left[-\Phi\left(\frac{17 + \frac{5}{\sqrt{n}} - \theta}{\frac{5}{\sqrt{n}}}\right) \right] \\ &= 1 - \inf_{\theta \leq 17} \left[\Phi\left(\frac{17 + \frac{5}{\sqrt{n}} - \theta}{\frac{5}{\sqrt{n}}}\right) \right] \end{aligned} \quad \left. \right\} \text{**}$$



$$\Phi(1) < \Phi(1+\delta)$$

$$\therefore \int_{-\infty}^1 \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = \int_{-\infty}^{1+\delta} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx.$$

Observe that

$$\text{for } \theta = 17, \frac{17 + \frac{5}{\sqrt{n}} - \theta}{\frac{5}{\sqrt{n}}} = 1$$

For any other $\theta < 17$

$$\frac{17 - \theta + \frac{5}{\sqrt{n}}}{\frac{5}{\sqrt{n}}} = 1 + \left(\frac{17 - \theta}{\frac{5}{\sqrt{n}}}\right)$$

which is greater than 1.

** we have used the fact

$$-\sup(-\Phi(x)) = \inf(\Phi(x))$$

Thus

$$\begin{aligned}\text{Size of the test} &:= 1 - \inf_{\theta \leq 17} \left[\Phi \left(\frac{17 + 5/\sqrt{n} - \theta}{5/\sqrt{n}} \right) \right] \\ &= 1 - \inf \Phi \left(\frac{17 + 5/\sqrt{n} - 17}{5/\sqrt{n}} \right) \\ &= 1 - \Phi(1).\end{aligned}$$

$$\therefore \boxed{\text{Size of test} = 1 - \Phi(1) \approx 1.5} \quad (\text{Approx}).$$

Let us take a diversion and think a bit about, what is a non randomized test.

For a non-randomized test, we observe a sample and see if it is in the critical region, then reject it and accept it if it is the acceptance region. In fact the problem arises when the observed sample is in the boundary of these two regions. For example, consider the example we discussed above. When $\theta \in (17, 18)$, we were not completely sure what decision to take. For values of $\theta > 17$ but very near 17 we are tempted to accept H_0 , while for values of θ near 18, we are more willing to reject it.

~~This~~ These are precisely in such a kind of scenario we take

the help of non-randomized test to come to a decision.

A randomized test T depends on the following

Critical function

$$\Psi_T(x_1, \dots, x_n) = P[H_0 \text{ is rejected} \mid (x_1, \dots, x_n) \text{ is observed.}]$$

Once we have observed (x_1, \dots, x_n) , we first compute

$\Psi_T(x_1, \dots, x_n)$ and then carry out an auxiliary Bernoulli trial (a single trial with an outcome, which is either a success or a failure), which does not depend on (x_1, \dots, x_n)

and for which the probability of success is $\Psi_T(x_1, \dots, x_n)$.
 If we observe a success in the Bernoulli trial, then we reject H_0 .

Let us see the following example from Mood, Graybill and Boes
 - "Introduction to the Theory of Statistics".

Example 2: Let us consider a sample of size 10, i.e. x_1, x_2, \dots, x_{10}
 from a Bernoulli population described as

$$f_x(x, \theta) = \theta^x (1-\theta)^{1-x}, \quad x = 0 \text{ or } 1$$

where $x = 0$ is failure $x = 1$ is success, and θ represents the probability of success and let of course $\theta > 0$. Do not think that this Bernoulli population is just coin tossing. It could be an experiment of an archer hitting the bull's eye, or checking whether a light bulb is functioning or not.

Let $H_0 : \theta < \frac{1}{2}$, versus $H_1 : \theta \geq \frac{1}{2}$

Let us consider the following test, T :

Reject H_0 ; if $\sum_{i=1}^{10} x_i > 5$

Accept H_0 ; if $\sum_{i=1}^{10} x_i \leq 5$.

Now what happens if $\sum_{i=1}^{10} x_i = 5$. This is where randomization comes in. The most simple way is to decide through a coin toss. If the toss shows head (H), then reject, H_0 and if it shows tail, (T), then accept H_0 . In this case how do we find the critical function. Let us partition the region of possible sample observations into three zones

$$A = \left\{ (x_1, \dots, x_{10}) : \sum_{i=1}^{10} x_i < 5 \right\}$$

$$B = \left\{ (x_1, \dots, x_{10}) : \sum_{i=1}^{10} x_i = 5 \right\}$$

and

$$C = \left\{ (x_1, \dots, x_{10}) : \sum_{i=1}^{10} x_i > 5 \right\}$$

$$\Psi_T(x_1, \dots, x_{10}) = \begin{cases} 1 & : \text{if } (x_1, \dots, x_{10}) \in C \\ \frac{1}{2} & : \text{if } (x_1, \dots, x_{10}) \in B \\ 0 & : \text{if } (x_1, \dots, x_{10}) \in A. \end{cases}$$

This means if $(x_1, \dots, x_{10}) \in C$, H_0 will be rejected with probability 1. If $(x_1, \dots, x_{10}) \in B$ it will be rejected with probability $\frac{1}{2}$. If $(x_1, \dots, x_{10}) \in A$, then H_0 will be accepted.

Section 2: Testing Simple Hypothesis against Simple Alternative.

In this section we will study the design for testing a simple hypothesis against a simple alternative. More precisely a very general way of writing, a simple hypothesis versus a simple random alternative. Given a random sample x_1, \dots, x_n we need to decide in which of the two populations, $f_x^0(\cdot)$ or $f_x^1(\cdot)$ ^{this} random sample comes from.

Thus we have

$$H_0: x_i \sim f_x^0(\cdot) \quad \text{against} \quad H_1: x_i \sim f_x^1(\cdot)$$

Suppose we just have one observation $x_1 = x_1$, then if $f_x^0(x_1) > f_x^1(x_1)$ we accept H_0 and reject it if $f_x^1(x_1) > f_x^0(x_1)$. This simple idea leads us to what is known as the likelihood ratio test.

- Simple Likelihood Ratio Test: Given the above H_0 versus H_1 , the simple likelihood test T is given as.

Reject H_0 : If $\lambda < k \} \quad (k > 0)$

Accept H_0 : If $\lambda > k \}$

If $\lambda = k$, then either we accept H_0 or reject H_0 , may be using a randomized test.

We have

$$\lambda = \lambda(x_1 \dots x_n) = \frac{\prod_{i=1}^n f_X^0(x_i)}{\prod_{i=1}^n f_X^1(x_i)} = \frac{L_0(x_1 \dots x_n)}{L_1(x_1 \dots x_n)}$$

and $k \geq 0$. Thus λ is the ratio of the Likelihood functions. This leads us to the following section.

Section 3: Most powerful test and Neyman-Pearson Lemma

Consider the following testing test of hypothesis problem:

$$H_0: \theta = \theta_0 \quad \text{versus} \quad H_1: \theta = \theta_1$$

↓
Null hypothesis

↓
Alternative hypothesis

Corresponding to any test T , of H_0 against H_1 , the power function is denoted by $\pi_T(\theta) = P_\theta(\text{Reject } H_0)$.

$$\text{Now } \pi_T(\theta_0) = P_{\theta_0}(\text{Reject } H_0)$$

$$\text{i.e. } \pi_T(\theta_0) = P(\text{Reject } H_0 \mid H_0 \text{ is true})$$

This is the probability of type-I error and this should be reasonably small.

While, $\pi_T(\theta_1) = P_{\theta_1}(\text{Reject } H_0 \mid H_0 \text{ is false})$

Of course we want $\pi_T(\theta_1)$ to be large. Further we can write

$$\pi_T(\theta_1) = 1 - P(\text{Accept } H_0 \mid H_0 \text{ is false})$$

$$\Rightarrow P(\text{Accept } H_0 \mid H_0 \text{ is false}) = 1 - \pi_T(\theta_1)$$

↓
Prob of Type-II error.

We denote it by

$$\beta_T(\theta_1) = 1 - \pi_T(\theta_1)$$

Ideally both $\pi_T(\theta_0)$ & $\beta_T(\theta_1)$ should be small. In practice one often fixes the value of $\pi_T(\theta_0)$ to a number say $0 < \alpha < 1$, i.e. $\pi_T(\theta_0) = \alpha$, often called the size of the test and thus try to find a test which minimizes $\beta_T(\theta_1)$. Such a test is often called the most powerful test, which is given we define below.

A test T^* for testing $H_0 : \theta = \theta_0$, versus $H_1 : \theta = \theta_1$, is called the most powerful test if $\pi_{T^*}(\theta_0) = \alpha$ and for any test T with $\pi_T(\theta_0) \leq \alpha$ we have

$$\beta_{T^*}(\theta_1) \leq \beta_T(\theta_1)$$

The following result of Neyman and Pearson shows us how to devise the most powerful test.

Theorem 1: Neyman-Pearson Lemma [Mood's, Graybill and Boes. Introduction to the theory of Stat.]

Let x_1, \dots, x_n be a random sample of size n from a population $f_X(x, \theta)$, where θ is one of the two values θ_0 or θ_1 .

Let us test the following hypotheses

$$H_0: \theta = \theta_1 \text{ versus } H_1: \theta = \theta_2$$

Further let us set $0 < \alpha < 1$; be fixed

Let $k^* > 0$ and $C^* \subset \Omega_n$ where Ω_n is the set of all possible sample values from $f_X(x, \theta)$. Assume that the following relationship holds. facts hold.

$$i) P_{\theta_0}((x_1, \dots, x_n) \in C^*) = \alpha$$

$$ii) \text{ Let } \lambda = \frac{L_0(\theta_0)}{L_1(\theta_1)}$$

$$\text{Let } \lambda = \frac{L_0(x_1, x_2, \dots, x_n, \theta_0)}{L_1(x_1, x_2, \dots, x_n, \theta_1)} = \frac{L_0}{L_1}$$

where $L_0 = L(x_1, \dots, x_n, \theta_0)$ \Rightarrow Likelihood function
 $L_1 = L(x_1, \dots, x_n, \theta_1)$

$$\begin{aligned} \lambda &\leq k^* & \text{if } (x_1, \dots, x_n) \in C^* \\ \lambda &\geq k^* & \text{if } (x_1, \dots, x_n) \in (C^*)^c; \quad [(C^*)^c = \Omega_n \setminus C^*] \end{aligned}$$

Then the test T^* corresponding to the critical region $C_{T^*} = C^*$
 is a most powerful test of size α for testing $H_0: \theta = \theta_0$
 against $H_1: \theta = \theta_1$.

Proof: (This can be skipped in first reading).

We shall for simplicity consider any other test T , such that
 $\pi_T(\theta_0) = \alpha$. Let the critical region for T be given by D .

$$\therefore \int \int \dots \int_{n\text{-fold integral over } C^*} L_0 d\sigma = \int \int \dots \int_D L_0 d\sigma = \alpha$$

Since $\pi_{T^*}(\theta_0) = \int \int \dots \int_{C^*} L_0 d\sigma = \alpha$, where T^* is the test corresponding to C^*

$$\text{Here } d\sigma = dx_1 \dots dx_n.$$

$$C^* = (C \cap D) \cup (C \cap D^c) \quad \text{and as } (C \cap D) \cap (C \cap D^c) = \emptyset$$

we have

$$\iiint_{C^*} L_0 dx = \iiint_{C^* \cap D} L_0 dx + \iiint_{C^* \cap D^c} L_0 dx = \alpha$$

Similarly

$$\iiint_{\text{D}} L_0 dx = \iiint_{C^* \cap D} L_0 dx + \iiint_{(C^*)^c \cap D} L_0 dx = \alpha$$

From the above two equations we have

$$\iiint_{C^* \cap D^c} L_0 dx = \iiint_{(C^*)^c \cap D} L_0 dx$$

Now inside C^* we have $\lambda = \frac{L_0}{L} \leq k^* \text{ or } L \geq \frac{L_0}{k^*}$

(Think why k^* has to be positive!)

$$\therefore \iiint_{C^* \cap D^c} L_1 dx \geq \iiint_{C^* \cap D^c} \frac{L_0}{k^*} dx = \iiint_{(C^*)^c \cap D} \frac{L_0}{k^*} dx \geq \iiint_{(C^*)^c \cap D} L_1 dx \quad (\because \frac{L_0}{k^*} \geq L \text{ on } (C^*)^c)$$

$$\therefore \iiint_{C^* \cap D^c} L_1 dx \geq \iiint_{(C^*)^c \cap D} L_1 dx$$

Now

$$\begin{aligned} \iiint_{C^*} L_1 dx &= \iiint_{C^* \cap D} L_1 dx + \iiint_{C^* \cap D^c} L_1 dx \\ &\geq \iiint_{(C^*)^c \cap D} L_1 dx + \iiint_{C^* \cap D^c} L_1 dx = \iiint_D L_1 dx. \end{aligned}$$

This shows that

$$\begin{aligned} \text{N} & \pi_{T^*}(\theta_1) \geq \pi_T(\theta_1), \\ \Rightarrow & \beta_{T^*}(\theta_1) \leq \beta_T(\theta_1) \end{aligned}$$

Showing that T^* is the most powerful test. This proves the result.

Example: Let x_1, \dots, x_n be a random sample of size n from an exponential distribution given as

$$f_X(x, \theta) = \theta e^{-\theta x}, \quad \theta > 0, \quad x > 0$$

$$\therefore L_0 = L(x_1, \dots, x_n | \theta_0) = (\theta_0)^n e^{-(\theta_0 \sum x_i)}, \quad \text{while}$$

$$L_1 = L(x_1, \dots, x_n | \theta_1) = (\theta_1)^n e^{-(\theta_1 \sum x_i)}.$$

Our aim is to develop the most powerful test of size α for the hypothesis testing problem

$$H_0: \theta = \theta_0 \text{ against } H_1: \theta = \theta_1, \quad (\theta_1 > \theta_0)$$

To find the most powerful test we have to effect find the k^* . Let us see how we do it in this case.

The idea is reject H_0 if $\lambda \leq k^*$ (from the Neyman-Pearson lemma)

$$\text{Thus } \frac{L_0}{L_1} \leq k^* \Rightarrow \left(\frac{\theta_0}{\theta_1} \right)^n e^{[-(\theta_0 - \theta_1) \sum x_i]} \stackrel{(\#)}{\leq} k^*$$

$$\Rightarrow \sum x_i \leq \frac{1}{\theta_1 - \theta_0} \log_e \left[\left(\frac{\theta_1}{\theta_0} \right)^n k^* \right]$$

This is obtained by taking logarithm on both sides of $(\#)$.

$$\text{Set } \frac{1}{\theta_1 - \theta_0} \log_e \left[\left(\frac{\theta_1}{\theta_0} \right)^n k^* \right] = k \text{ (say).}$$

Thus

$$\frac{L_0}{L} \leq k^* \Rightarrow \sum_{i=1}^n x_i \leq k.$$

$$\begin{aligned}\therefore \alpha &= P_{\theta_0} \left[(x_1, \dots, x_n) \in C^* \right] \\ &= P_{\theta_0} \left[\sum_{i=1}^n x_i \leq k^* \right]\end{aligned}$$

Now since $x_i \sim \exp(\theta)$ we know from our study of sampling distribution we know that

$$\sum_{i=1}^n x_i \sim \text{Gamma}(n, \theta)$$

\therefore If we set $Y = \sum_{i=1}^n x_i$, then

$$P_{\theta_0} \left[\sum_{i=1}^n x_i \leq k \right] = P_{\theta_0} [Y \leq k] = \int_0^k \frac{1}{\Gamma(n)} \theta_0^n y^{n-1} e^{-y \theta_0} dy = \alpha$$

Once sample values are observed, then

$$\int_0^k \frac{1}{\Gamma(n)} \theta_0^n y^{n-1} e^{-y \theta_0} dy = \alpha$$

is an equation in k from which k can be computed in fact by using the tables. Once the k is known, - our test T^* :

Reject H_0 : if $\sum x_i \leq k$.

By Neyman-Pearson's lemma, we see that T^* is the most powerful test.

—x—

The course JTSO201A formally ends here