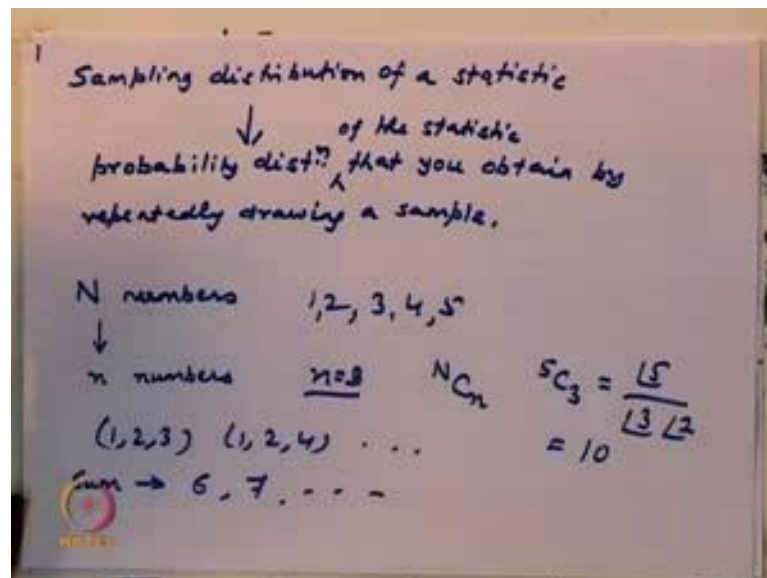


Introduction to Biostatistics
Prof. Shamik Sen
Department of Bioscience and Bioengineering
Indian Institute of Technology, Bombay

Lecture - 27
Sampling distributions and Central limit theorem Part-II

Hello and welcome to today's lecture. So, in the last class we had briefly touched upon the idea of central limit theorem right. One of the most important theories in the field of statistics, so before going to central limit theorem, I will begin by discussing about again sampling distribution right.

(Refer Slide Time: 00:36)



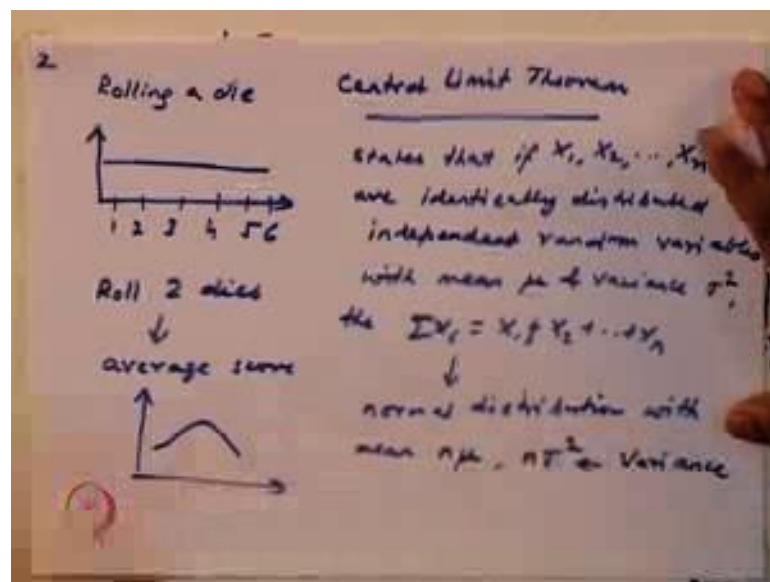
So, what is the sampling distribution and you can find the sampling distribution of a statistic. So, by this I mean you can compute the sampling distribution of a mean of a median or any other measure that you define.

So, what is the sampling distribution of the statistics? So, a sampling distribution of the statistic is the probability distribution that you obtain. So, probability distribution of the statistic that you obtain by repeatedly drawing a sample, in other words imagine you have a population of n numbers. Let us say n numbers and from this you draw a sample of n numbers. So, based on this for example, if my 5 numbers are 1 2 3 4 5 and let us say my n is equal to 3 right. So, in how many ways can I draw 3 numbers from n? So, it can be n c n, the different ways in which I can draw 3 numbers from 5 right. So, this is 5 c 3

equal to factorial 5 by factorial 3 factorial 2 is equal to 10. So, there are 10 ways of drawing a sample of 3 and one example is 1 2 3 or one 2 4. So on and so forth.

So, from these you can let us say compute the sum of these numbers in that case this will be 6 this will be 7 and so on and so forth. So, you can ask the question what is the probability distribution of the sum that you obtain. And this is the sampling distribution. So, ideally the concept of sampling distribution is very important, because you do not know the; you want to find out some information about the population based on the sample that you draw. So, central limit theorem is that theorem which allows us to correlate the sample mean or median or any other metric with the population metric.

(Refer Slide Time: 03:30)



So, what does central limit theorem state. So, I had discussed this brief example in last class right, where if you are interested in rolling a die and checking out it is value that you get right. So, the distribution would be fully flat. So, each of these values is 1 2 3 4 5 or 6 this distribution is flat, but if you roll 2 dice together and you compute the average score that you get. So, you will see begin to see that slowly your distribution will look closer and closer like a normal distribution. And with increase with the averaging process, and this is what central limit theorem says. So, what does the central limit theorem state?

So your central limit theorem states that if x_1, x_2, \dots, x_n are identically distributed, independent random variables, with mean μ and variance σ^2 . Then summation of x_i that is $x_1 + x_2 + \dots + x_n$, this is also a random variable, it will follow

a normal distribution with mean n times μ and variance n sigma square. So, this is your variance. And this is very easy to prove.

(Refer Slide Time: 06:00)

Handwritten mathematical derivation on a whiteboard:

$$X = X_1 + X_2 + \dots + X_n$$

$$E(X) = E(X_1 + X_2 + \dots + X_n)$$

$$= E(X_1) + E(X_2) + \dots + E(X_n)$$

$$= n\mu + \underbrace{\mu + \dots + \mu}_{n \text{ times}}$$

$$= n\mu$$

$$\text{Var}(X) = \sum \text{Var}(X_i)$$

$$= n\sigma^2$$

$$\downarrow$$

$$(\sigma\sqrt{n})^2$$

Annotations on the right side of the whiteboard:

- $\sum X_i \rightarrow$ normally distributed with mean $n\mu$ & variance $n\sigma^2$
- $\frac{\sum X_i - n\mu}{\sigma\sqrt{n}} \rightarrow$ standard normal variable

So, I can write x as x_1 plus x_2 plus x_n . So, since each of these variables are independent, the expectation of x is e of x_1 plus x_2 plus x_n , because each of these random variables are independent. We can break it down into e of x_1 plus e of x_2 So on and so forth, plus e of x_n . So, this is going to be n . So, μ plus μ plus μ , this is n times. So, that is why this will give you the expectation of this x is simply equal to $n\mu$. And by the same token I can write variance of x is equal to summation of variance of x_i by again because these random variables are independent and so, I will get n times sigma square. So, this variance is n times sigma square

So, what in other words which would mean that if you have, I can say that summation x_i is normally distributed with mean $n\mu$ and variance n sigma square. So, then I can say that $\frac{\sum x_i - n\mu}{\sigma\sqrt{n}}$ is nothing, but sigma root of n whole square right. So, this is approximately equal to a standard normal variable.

(Refer Slide Time: 07:57)

4

$$P(\sum X_i < x) \approx P\left[\frac{\sum X_i - n\mu}{\sigma\sqrt{n}} < z\right]$$

Insurance Company \rightarrow 25000 policy holders.
Yearly claim of a policy holder is a RV with
a mean of 320 $\&$ a s.d. of 540.
The chance that the total claim disbursed by
the company is greater than 8.3 billion dollars.

$X_i \rightarrow$ RV of claim of i^{th} individual
 $\hookrightarrow \mu = 320 \quad \sigma = 540$

CLT $\sum X_i \rightarrow$ normal distribution with
mean = $320 \times 25000 = 8$ million
 $\sigma = \sigma\sqrt{n} = 540\sqrt{25000} \approx 8.5 \times 10^4$

So, in other words in other words probability of summation x_i , let us say less equal less than x you can approximate is the probability of summation x_i minus $n\mu$ by σ root n less than x this is the approximation.

So, let us first try out few examples, how this theory can be used. Imagine you have a insurance company which has 25000 policy holders. So, the yearly claim of a policyholder is a random variable with a mean of 320 right and standard deviation of 540. So, this is the background information. So, we want to compute the chance that the total claim disbursed by the company disbursed by the company is greater than 8.3 million dollars.

So, this is what we want to compute the chance that the total claim is greater than 8.3 billion dollars. So, if you have to make use of the central limit theorem. So, what we are given is. So, each x_i is a random variable of claim of i^{th} individual. So, we want to. So, the total. So, the total claim disbursed by the company is going to be summation of x_i . So, each x_i has a mean equal to 320 and σ equal to 540 dollars. So, my summation $\sum x_i$ must follow a normal distribution. So, this is by central limit theorem. What I can say is summation $\sum x_i$ follows a normal distribution with mean equal to 320 into the total number of policyholders. So, for each x_i you have mean as 320. So, for summation $\sum x_i$ the mean has to be 320 into 25000.

So, if you plug in the numbers this will come to be 8 million and standard deviation will be equal to I can say sigma star equal to sigma into root of n. So, which is 540 into root of 25000 so, this will approximately come to be 8.5 into 10 to the power 4.

(Refer Slide Time: 12:07)

$$P[X_i > 8.3 \text{ million}]$$

$$= P\left[\frac{X_i - \mu}{\sigma\sqrt{n}} > \frac{8.3 \times 10^6 - 8 \times 10^6}{8.5 \times 10^4}\right]$$

$$= P[Z > 3.5] \approx 0.00023$$

Ex: Bridge
 ↳ the weight it can withstand without any damage is a random variable with mean of 400 tons & $\sigma = 40$ tons.
 Car wt. → RV with mean 3 & s.d. = 0.3.
 How many cars would have to be on the bridge for probability of damage = 0.1.

So, we are asked to calculate probability that summation x_i will be greater than 8.3 million. So, we can approximate this as the probability summation x_i minus $n\mu$ by $\sigma\sqrt{n}$ is greater than 8.3×10^6 for million minus. So, $n\mu$ is going to be 8×10^6 that we found and $\sigma\sqrt{n}$ is going to be 8.5×10^4 to the power 4. So, this is equal to probability of. So, this is your standard normal variable now. So, this is z greater than. So, if you plug in the values this will give you a value of 3.5 and this comes to approximately 0.00023.

So, as you can see that where there is very little chance, that the claim is going to be more than 8.5 million. So, this is how you make use of central limit theorem, to gain an idea of or to make a prediction of calculating the chance that this particular insurance company will have to disperse this once. And these are calculations which the company does at its end to figure out how much the insurance amount should charge. So, that it makes a certain amount of profit.

So, let us go to another example. So, think of a bridge. So, each bridge has a certain capacity of how much weight it can bear. So, imagine for a given bridge, so the weight it can bear. The weight it can withstand without any damage is a random variable with

mean of with mean of 400 tons and sigma equal to 40 tons. So, cars fly on the bridge on an average. So, the car waits imagine, it is a random variable with mean 3 and with standard deviation 0.3.

So, the question is how many cars would have to be on the bridge for probability of damage equal to 0.1. So, what this question tries to understand is how many cars is the you know can the bridge bear. So, there is damage of 0.1. So, how do we go about this problem? So, what you have to calculate what we have been asked to calculate is.

(Refer Slide Time: 16:20)

$x_i = \text{wt. of } i^{\text{th}} \text{ car} \rightarrow \mu=3$
 $\sigma=0.3$
 $n \text{ cars on the bridge}$
 $\hookrightarrow \text{load on the bridge} = \sum x_i$
 $\sum x_i \rightarrow \text{R.V. with mean } 3n \text{ \& s.d. } 0.3\sqrt{n}$
 Damage on the car
 $\hookrightarrow 3n \rightarrow \text{total load on the car}$
 > 400
 $y = \sum x_i - W \rightarrow \text{R.V.} \rightarrow \text{mean} \rightarrow 3n - 400$
 $\text{S.d.} \rightarrow$
 $\text{Var}(y) = \text{Var}(\sum x_i) + \text{Var}(W)$
 $= n \times 0.3^2 + 40^2 = 0.09n + 1600$

So, let us say x_i is weight of i th car. So, if there are n cars on the bridge. So, the load on the bridge is equal to. So, since x , so the load on the bridge is equal to summation of x_i right. So, what can I say about summation of x_i .

So, each of these x_i is a random variable with μ is equal to 3 and σ is equal to 0.3. So, I can say summation is $\sum x_i$ is a random variable with mean 3 times n . If you have n cars 3 is the mean for each car. So, mean will be $3n$ and standard deviation will be 0.3 into root n 0.3 into root n . So, when I want to calculate the damage on the car. So, from the damage on the car we are thinking of a situation that $3n$ which is the total load on the car. This is greater than 400. 400 is the weight which can be withstood by the bridge. So, $\sum x_i - w$, so this is the effective weight mismatch between the weight of that cars on the bridge and the weight that can be withstood is a random variable, with mean this has a mean of $3n - 400$ and a standard deviation.

So, how do we calculate the standard deviation here you have 2 random variables. So, you define a random variable. Let us say y as summation of x_i minus w mean is this and you know. So, since these 2 are independent of each other. So, your variance of y will be variance of summation x_i plus variance of w and variance of summation x_i is nothing, but n times sigma square and this is whatever has been given to you which is 40 tons. So, this is 40 standard deviations. So, 40 square, so this is n into 0.3 square plus 40 square equal to $0.09n$ plus 1600. So, variance is this. So, y which is defined as the difference between the net weight from the cars on the bridge to the weight that can be withstood is this. So, y is the random variable, which has a mean of $3n$ minus 400 and a variance of $0.09n$ plus 1600.

(Refer Slide Time: 20:01)

The image shows a whiteboard with handwritten mathematical steps:

$$P\{\sum x_i - w \geq 0\} \approx 0.1$$

$$\downarrow$$

$$P\left\{\frac{\sum x_i - (3n - 400)}{\sqrt{0.09n + 1600}} \geq \frac{0 - (3n - 400)}{\sqrt{0.09n + 1600}}\right\} = 0.1$$

$$P\left\{Z \geq -\frac{(3n - 400)}{\sqrt{0.09n + 1600}}\right\} = 0.1 = P\{Z \geq 1.28\}$$

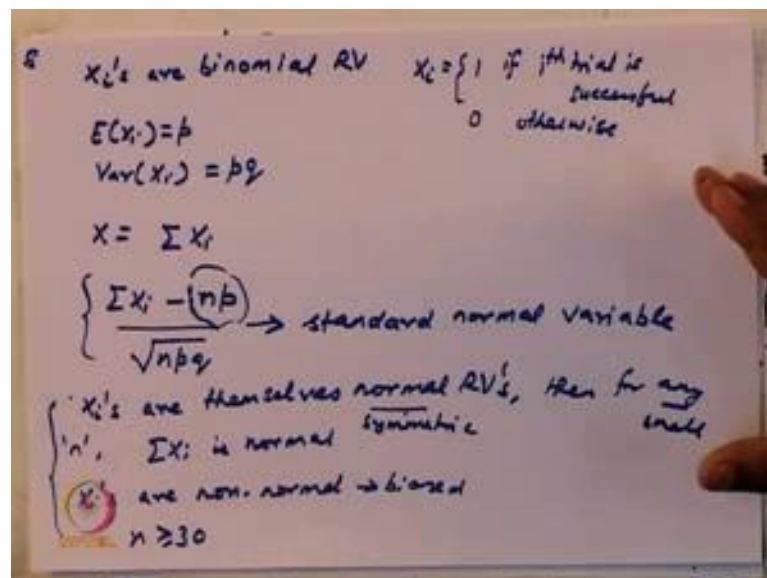
$$\frac{400 - 3n}{\sqrt{0.09n + 1600}} \leq 1.28 \rightarrow n \approx 117$$

So we have been asked to calculate probability summation x_i minus w greater equal to 0. Why? Because if summation x_i is less than w there is no damage on the bridge. Only when summation x_i begins to exceed w , then there any chance of the bridge. So, we have been asked to calculate this is roughly equal to 0.1. So, probability of summation x_i minus w , this can be written as probability of z summation x_i minus of $3n$ minus 400. So, this is the net mean minus root of $0.09n$ plus 1600. Because the variance was $0.09n$ plus 1600 is greater equal to 0, p of z greater or equal to 0. So, p of this is basically greater equal to 0 minus of $3n$ minus 400 by root of $0.09n$ plus 1600. And this we are given approximately equal to 0.1

So, probability of z greater equal to minus of 3 n minus 400 by root of 0.09 n plus 1600 is equal to 0.1. So, if you look up the table, this is roughly p of z greater equal to 1.28. So, this would mean that 400 minus 3 n by root of 0.09 n plus 1600 is less or equal to 1.28. And from here you can find out the value of n. So, this is one equation you take a square and determine the value of n. So, this n will be ordered 117. So, this allows if there are roughly 120 cars then there is a chance that 10 there is a 10 percent chance that the bridge will get damaged.

So, here we made use of central limit theorem, but we first defined our random variable y as the difference between the total weight of the cars on the bridge minus the weight that could be withstood because we are calculating the chance of damage.

(Refer Slide Time: 22:41)



So, imagine your x. So, x_i s are binomial random variables, x can be of any distribution x can follow any distribution. Let us assume that x_i s have binomial random variables. Then we know that e of x_i . So, we can call x_i one if ith trial is successful equal to 0 otherwise. So, we have previously demonstrated that e of x_i is p and variance of x_i is equal to p q. So, if we define a random variable x which is summation of x_i . So, in that case again by you invoking central limit theorem, I can say summation x_i minus. So, x is summation of x_i . So, summation x_i will have a mean of n p and a variance of n p q. So, so n p q is going to be the variance of this x and n p is going to be the mean of this x. So, this means that summation x_i minus n p by root of n p q is a standard normal variable.

So, one thing is, there is a requirement on the value of n depending on what distribution these x_i 's follow. So, generally if x_i 's, if x_i 's are themselves normal random variables then for any n summation x_i is going to be normal, but if x_i 's are non-normal, and maybe even biased, in that case you need to have n greater or equal to 30 for this distribution to be a normal distribution. So, for a binomial random variable, we have summation x_i follows summation x_i follows normal with mean of np and variance of npq . So, I will conclude my lecture today by saying that we would we can briefly remember about central limit theorem and how we can use central limit theorem to estimate certain things given the sample distribution.

So, as we saw that summation x_i will follow a normal distribution with mean of n times μ and variance of $\sigma^2 n$ which means the standard deviation will be $\sigma \sqrt{n}$. And depending on the type of distributions, so if n if you are underlying x_i 's are themselves normal random variables or even symmetric, then for any or small values of $n \geq 2$ or 3 you will have summation x_i or \bar{x} or summation x_i by n follow normal distribution, but if x_i 's are non-normal in nature, particularly they are an asymmetric then you need n greater equal to 30.

With that I conclude my lecture today. I look forward to meet you again in next class.