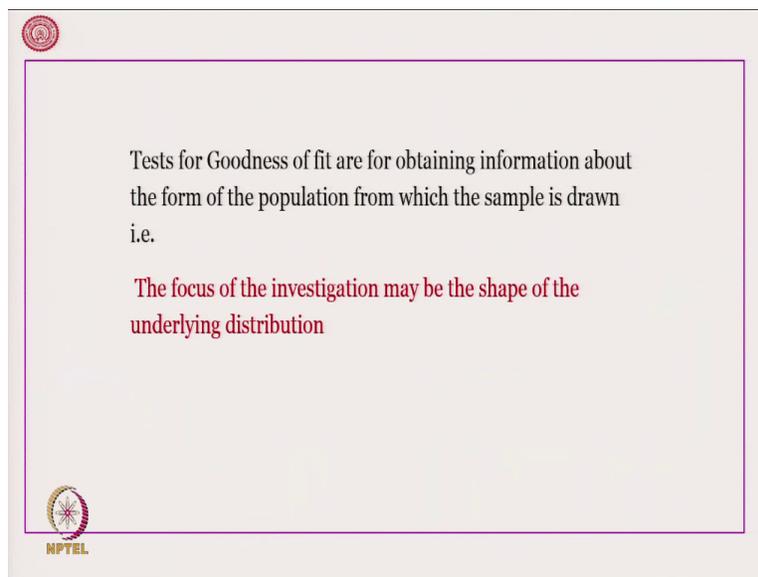


**Nonparametric Statistical Inference**  
**Professor. Niladri Chatterjee**  
**Department of Mathematics**  
**Indian Institute of Technology, Delhi**  
**Lecture No. 06**

Welcome students to MOOC's series of lectures on Nonparametric Statistical Inference, this is lecture number 6. As I said at the end of last class that in this lecture, we shall study some tests for goodness of fit.

(Refer Slide Time: 0:41)



Now, the question is what are test for goodness of fit? Basically this is about obtaining information about the form of the population from which the sample is drawn. And the focus of the investigation may be the shape of the underlying distribution.

(Refer Slide Time: 1:02)

**Example 1**

Suppose one wishes to check if a coin is unbiased from the outcomes of  $n$  tosses of the coin.

We know that if the number of tosses is  $n$  then the shape of the distribution will depend upon the value of  $p$  which is the probability of obtaining say, a head in a single toss of the coin.

If the coin is unbiased, i.e.  $p = 0.5$  we know the distribution is symmetric about the Mean.

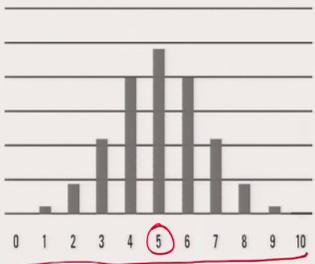


Let me now illustrate with respect to one example. Suppose one wishes to check if a coin is unbiased from the outcomes of  $n$  tosses of the coin. We know that if the number of tosses is  $n$ , then the shape of the distribution will depend upon the value  $p$ ; which is the probability of obtaining a head in a single toss of the coin. Or, in other words we call it the probability of success. If the coin is unbiased that is  $p$  is equal to 0.5; we know that the distribution is symmetric about the mean.

(Refer Slide Time: 1:50)

**Example 1**

When  $p = 0.5$  get a bar graph of probabilities as follows:



Number of Heads (k)	Probability
0	0.0010
1	0.0107
2	0.0547
3	0.2209
4	0.3770
5	0.4768
6	0.3770
7	0.2209
8	0.0547
9	0.0107
10	0.0010

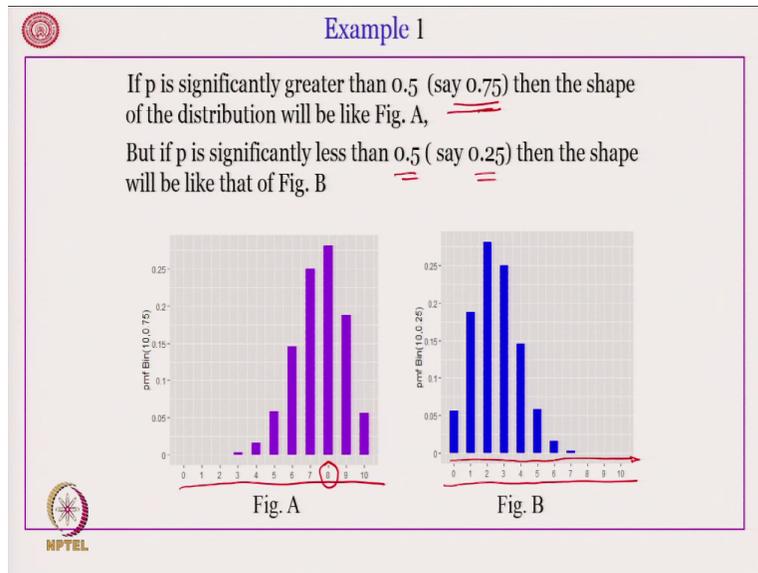
What happens if the coin is not unbiased?



For illustration let us look at the graph that we have plotted; here  $n$  is equal to 10,  $p$  is equal to 0.5. And therefore mean

is 5 and we can see that the distribution is symmetrical in the range 0 to 10, for the binomial 10 comma 0.5 variable. The question is what happens if the coin is not unbiased?

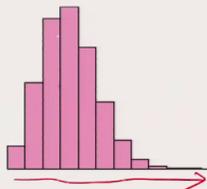
(Refer Slide Time: 2:25)



So, if the coin is not unbiased,  $P$  is different from 0.5, it maybe higher, it maybe lower than that. So, let us consider a  $p$  which is say 0.75. Then for the same  $n$  that is 10, we get a distribution which will look like this; which is given in the figure A. And as you can understand that since the value of  $p$  is equal to 0.75, the probability of getting more heads is higher than getting a lesser number of heads. And consequently this is the type of distribution that one can obtain, with the mode value at 8. On the other hand, if the value of  $p$  is significantly lower than 0.5; so for this illustration we have considered it to be 0.25.

And we can see that the bar chart describing the probabilities look like this, which is not symmetric around the mean at all. And we can see that it has a skewed distribution.

(Refer Slide Time: 3:41)



Therefore if we get a discrete distribution like this:

We may like to test if the data is from Binomial Distribution with  $p = 0.25$  with appropriate value of  $n$

NPTEL

Therefore, suppose we get a discrete distribution from the observed sample something similar to this; which we can see you skewed on this side. Therefore, it is quite natural that we like to test whether it is coming from a binomial distribution with  $p$  is equal to 0.25, for an appropriate value of  $n$ . And that is what testing of hypothesis is all about.

(Refer Slide Time: 4:15)

### Example 2

Suppose we want to test whether a dice is unbiased i.e. each of the six numbers is equally likely i.e. probability of each of the six numbers is  $1/6$ .

That is to test the null hypothesis

$$H_0: p_1 = p_2 = p_3 = p_4 = p_5 = p_6 = 1/6$$

where,  $p_j = P(X=j)$  is the random variable indicating the outcome of a single roll of the dice

A major test in this regard is Chi-Square test, which has several subtle variations

NPTEL

Let me give you another example. Suppose we want to test whether a dice is unbiased, that is each of the 6 faces of the die will have equal probability. Therefore, the hypothesis that we are testing is  $p_1$  is equal to  $p_2$  is equal to say up to  $p_6$ ; and each one of them is  $1/6$ . Where,  $P_j$  is

probability that  $X$  is taking the value  $j$ ; where  $X$  is the random variable indicating the outcome of a single roll of the dice. Question is how to test that? Let us understand that if the coin is unbiased; it does not mean that in 100 tosses we shall get exactly 50 heads and 50 tails. That may not happen, but there will be deviations from this obtained value from the expected value  $np$ .

However, if the deviation is too much, then we would say that it is not acceptable; and or in other words we are going to reject the null hypothesis. So, corresponding test we need to do, the major test in this regard is Chi-Square test, which has several subtle variations; but, we shall study the basic Chi-Square test today. So, the idea of frequency Chi-Square is like this.

(Refer Slide Time: 5:59)

**Frequency Chi-Square  $\chi^2$**

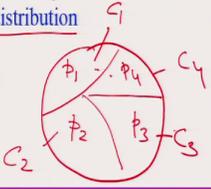
Suppose a population consists of mutually exclusive classes with the proportion of members belonging to the  $j^{\text{th}}$  class being  $p_j$ ,  $j = 1, 2, \dots, k$ .

Note that  $\sum_{j=1}^k p_j = 1$

Suppose a random sample of size  $n$  is taken from the population.

The probability that there are  $n_j$  many objects from the class  $j$ ,  $j = 1, 2, \dots, k$ , can be obtained using the Multinomial distribution

If  $k = 2$ , it is a Binomial distribution



NPTEL

Suppose a population consists of mutually exclusive classes with the proportion of members belonging to the  $j^{\text{th}}$  class is  $p_j$ . What does it mean? It means that suppose this is the entire population; it is divided into several classes. Suppose in this case I have taken four classes; so  $p_1$  is the proportion of elements belonging to class 1.  $p_2$  is the proportion of elements belonging to class 2, like that for class 3 and class 4. And we have to understand that they are mutually exclusive; that means that no element can belong to two of these classes at the same time. Now,

$$\sum_{j=1}^k p_j = 1$$

where  $k$  is the number of classes.

Now, suppose we have taken  $n$  samples from the population; the probability that there are  $n_j$  many objects from the class  $j$ ,  $j$  is equal to 1 to  $k$ , can be obtained using the multinomial distribution. I hope you are familiar with multinomial, but I will explain it now. If  $k$  is equal to 2, then we get a binomial distribution.

(Refer Slide Time: 7:33)

The Density Function for Multinomial Distribution with parameters  $(n, p_1, p_2, \dots, p_k)$  at  $(n_1, n_2, \dots, n_k)$  is:

$$\frac{n!}{n_1! n_2! \dots n_k!} p_1^{n_1} p_2^{n_2} \dots p_k^{n_k} \quad \text{when } n_1 + n_2 + \dots + n_k = n$$

Handwritten notes in red on the slide:

- $n_1$  in # of elements in the sample
- when  $n_1 + n_2 + \dots + n_k = n$
- $n_k$  from class

The derivation shown is:

$$\binom{n}{n_1} p_1^{n_1} \cdot \binom{n-n_1}{n_2} p_2^{n_2} \dots \binom{n-n_1-n_2-\dots-n_k}{n_k} p_k^{n_k}$$

$$\frac{n!}{n_1! (n-n_1)!} \cdot \frac{(n-n_1)!}{n_2! (n-n_1-n_2)!} \dots \frac{1!}{n_k!}$$

$$\frac{n!}{n_1! n_2! \dots n_k!} p_1^{n_1} p_2^{n_2} \dots p_k^{n_k}$$

So, what is multinomial distribution with parameters  $p_1, p_2, p_k$ , and  $n$  which is the number of samples. The probability distribution will look like this,  $n$  factorial upon  $n_1$  factorial,  $n_2$  factorial into  $n_k$  factorial, into  $p_1$  to the power  $n_1$ ,  $p_2$  to the power  $n_2$ ,  $p_k$  to the power  $n_k$ ; where  $n_1, n_2, n_k$  is that  $n_1$  is number of elements in the sample from class  $c_1$ . Similarly,  $n_2$  from  $c_2$  and  $n_k$  from  $c_k$ ;  $n_1$  plus  $n_2$  plus  $n_k$  is equal to  $n$ , the total number of samples.

Now, how do we get that, it is very simple; so out of  $n$  we first select  $n_1$  of them that can be done in  ${}^n C_{n_1}$ . And let us put them in class number one that is  $c_1$ , so that probability is  $p_1$  to the power  $n_1$ . Now, how many are remaining?  $n$  minus  $n_1$  are there; so out of that  $n$  minus  $n_1$ , we choose  $n_2$  of them and put them in the class  $c_2$ . That probability is  $p_2$  to the power  $n_2$ ; like that finally we will have there only  $n_k$  many elements left.

So,  $n_k$   $c_k$  and that will give to the class  $K$ ; that probability is  $p_k$  to the power  $n_k$ . So, this is the basic derivation of the multinomial distribution; if we simplify, it will see that  $n$  factorial upon  $n_1$  factorial into  $n$  minus  $n_1$  factorial, multiplied by let us simplify this. So, that will give you  $n$  minus  $n_1$  factorial, upon  $n_2$  factorial into  $n$  minus  $n_1$  minus  $n_2$  factorial. Like that we will go and

here we will get 1 multiplied by, as you can understand  $p_1$  to the power  $n_1$ ,  $p_2$  to the power  $n_2$ ,  $p_k$  to the power  $n_k$ .

And therefore if we look at these cancels with this in a similar way corresponding to class 3; this is going to cancel with the numerator. So, finally we will end up with  $n$  factorial upon  $n_1$  factorial,  $n_2$  factorial up to  $n_k$  factorial, into  $p_1$  to the power  $n_1$ ,  $p_2$  to the power  $n_2$ ,  $p_k$  to the power  $n_k$ . So, that is how the multinomial density function is obtained. I hope most of you are familiar with this, but still I wanted to clarify the density function.

(Refer Slide Time: 11:23)

Observe that

$$\frac{n!}{n_1! n_2! \dots n_k!} p_1^{n_1} p_2^{n_2} \dots p_k^{n_k}$$

$$= \frac{e^{-(np_1 + np_2 + \dots + np_k)}}{e^{-(np_1 + np_2 + \dots + np_k)}} * \frac{n^n}{n^n} * \frac{n!}{n_1! n_2! \dots n_k!} p_1^{n_1} p_2^{n_2} \dots p_k^{n_k}$$

$$= \frac{e^{-np_1} e^{-np_2} \dots e^{-np_k}}{e^{-(np_1 + np_2 + \dots + np_k)}} * \frac{n^{n_1 + n_2 + \dots + n_k}}{(np_1 + np_2 + \dots + np_k)^n}$$

$$* \frac{n!}{n_1! n_2! \dots n_k!} p_1^{n_1} p_2^{n_2} \dots p_k^{n_k}$$

Now, we shall work with this formula, we write it slightly complicated way. So, we write it as  $e$  to power minus  $np_1$  plus  $np_2$  plus  $np_k$ , and divide by the same quantity. Similarly, again we multiply it by  $n$  to the power  $n$  and divide by the same quantity; and then this is the formula that we already had. Now, let us simplify this term, so we can write this as  $e$  to the power minus  $np_1$ ,  $e$  to the power minus  $np_2$ , up to  $e$  to the power minus  $np_k$ ; denominator we keep unchanged.

And numerator again this one we write,  $n$  to the power  $n_1$  plus  $n_2$  plus  $n_k$ , divided by this  $n$  to the power  $n$ . We write it as  $np_1$  plus  $np_2$  plus  $np_k$  whole to the power  $n$ , multiplied by the standard density thing;  $n$  factorial upon  $n_1$  factorial et-cetera  $p_1$  to the power  $n_1$ ,  $p_2$  to the power  $n_2$  into  $p_k$  to the power  $n_k$ . So, this is little bit of I gave jugglery that we have done, but for some purpose as it will be clear in the next slide.

(Refer Slide Time: 12:50)

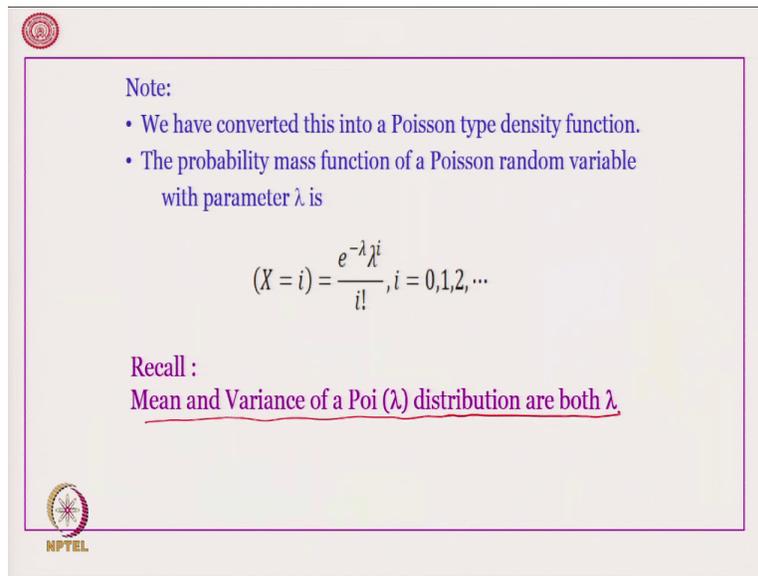
Thus we can rewrite the expression as

$$\frac{\frac{\exp(-np_1)(np_1)^{n_1}}{n_1!} * \frac{\exp(-np_2)(np_2)^{n_2}}{n_2!} * \dots * \frac{\exp(-np_k)(np_k)^{n_k}}{n_k!}}{\frac{\exp(-(np_1+np_2+\dots+np_k))(np_1+np_2+\dots+np_k)^n}{n!}}$$


So, we write the entire expression, now as e to the power minus np1 into np1 to the power n1, upon n1 factorial. So, let us go back to the slide to the power minus np1, from here then n to the power n1 multiplied by p1 to the power n1. This we write it as np1 to the power n1, then there is a n1 factorial which comes here; and therefore entire expression as, e to the power minus , into np1 to the power n1, upon n1 factorial, into like that up to e to the power minus npk, npk to the power nk upon nk factorial.

This whole thing divided by e to the power minus np1 plus np2 plus npk, multiplied by np1 plus np2 plus npk whole to the power n. This is in the denominator divided by n factorial, this comes from here.

(Refer Slide Time: 14:10)



Note:

- We have converted this into a Poisson type density function.
- The probability mass function of a Poisson random variable with parameter  $\lambda$  is

$$(X = i) = \frac{e^{-\lambda} \lambda^i}{i!}, i = 0, 1, 2, \dots$$

Recall :  
Mean and Variance of a Poi ( $\lambda$ ) distribution are both  $\lambda$



Therefore, what we have done? So, we convert the whole thing into a Poisson type density function, and I am sure all of you are familiar with Poisson. A Poisson random variable  $X$  takes the value  $i$  with the probability,

$$: \frac{e^{-\lambda} \lambda^i}{i!}, i = 0, 1, 2, \dots$$

where  $\lambda$  is the parameter of the distribution. Now, recall that the mean and variance of a Poisson  $\lambda$  distribution are both  $\lambda$ .

(Refer Slide Time: 14:59)



The limiting distribution of a Poisson random variable  $X$  with Parameter  $\theta$  as  $\theta \rightarrow \infty$  is Normal random variable with both The parameters equal to  $\theta$ ,  $N(\theta, \theta)$

Therefore,  $\lim_{\theta \rightarrow \infty} \frac{X - \theta}{\sqrt{\theta}} = N(0,1)$

Now with respect to the above problem note the following:  
 As,  $n \rightarrow \infty, np_j \rightarrow \infty \quad \forall j = 1, 2, 3, \dots k$

If  $p_j$  is the probability of an observation falling into class  $j$ , then the random variable  $X_j$  denoting the number of observations in class  $j$  will have expected value  $np_j$ ,  $j = 1, 2, \dots k$



Now, the limiting distribution of a Poisson distribution or a Poisson random variable  $X$  with parameter  $\theta$ . As  $\theta$  goes to infinity is a normal random variable with both the parameters equal to  $\theta$ , that is normal  $\theta, \theta$ . Therefore, limit  $\theta$  going to infinity,  $X$  minus  $\theta$  upon root  $\theta$ ; this converges to normal  $0, 1$ ; because, we have standardized it by subtracting the mean and dividing it by the standard deviation; which is root over  $\theta$ .

Now, with respect to the above problem that is the multinomial distribution note that, if  $n$  goes to  $\infty$ ; then  $np_j$  goes to  $\infty$  for all  $j$ ,  $j$  is equal to  $1$  to  $k$ . Hence, the  $p_j$  is the probability of an observation falling into the class  $j$ , then the random variable  $X_j$  denoting the number of observations of class  $j$ ; will have expected value  $np_j$ . This is very clear because the proportion coming from class  $j$  is  $p_j$ , for  $j$  is equal to  $1$  to  $k$ .

(Refer Slide Time: 16:22)

With this now let us examine the joint distribution:

$$\frac{\frac{\exp(-np_1)(np_1)^{n_1}}{n_1!} * \frac{\exp(-np_2)(np_2)^{n_2}}{n_2!} * \dots * \frac{\exp(-np_k)(np_k)^{n_k}}{n_k!}}{\frac{\exp(-(np_1+np_2+\dots+np_k))(np_1+np_2+\dots+np_k)^n}{n!}}$$

Handwritten annotations in red:

- For the first term:  $Poi(np_1)$  taking the value  $n_1$
- For the second term:  $Poi(np_2)$  taking the value  $n_2$
- For the denominator:  $Poi(n)$  taking the value  $n$

NPTEL logo is visible in the bottom left corner of the slide.

With this now let us examine the expression just we have just obtained some time back. So, it is  $e^{-np_1} (np_1)^{n_1} / n_1!$  to the power minus lambda lambda power  $n_1$ , upon factorial  $n_1$ ; so, this looks like a Poisson random variable with parameter  $np_1$  taking the value  $n_1$ . Similarly, this is a Poisson random variable with  $np_2$  as the parameter or we can write it as  $Poi$  with  $np_2$ , taking the value  $n_2$ . Similarly for the  $k$ th one and in the denominator we have it is a Poisson random variable with parameter  $np_1 + np_2 + \dots + np_k$ . That means, it is a Poisson random variable with parameter  $n$  taking the value  $n$ .

(Refer Slide Time: 17:38)



Thus, the numerator is the product of  $k$  Poisson random Variables  $X_1, X_2, \dots, X_k$  taking values  $n_1, n_2, \dots, n_k$ , respectively.

This can be treated as the product of  $k$  independent Poisson random variables with respective parameter

$$np_j, j = 1, 2, \dots, k$$


Therefore, the whole numerator is the product of  $k$  Poisson random variables  $X_1, X_2, X_k$  taking values  $n_1, n_2, n_k$ . And this can be treated as a product of  $k$  independent Poisson random variables with respect parameters  $np_1, np_2$  up to  $np_k$ .

(Refer Slide Time: 18:02)



As  $np_j \rightarrow \infty \forall j$  each  $X_j$  can be approximated as  $N(np_j, np_j)$ ,  $j = 1, 2, \dots, K$

Since sum of independent Poisson distributions  $Poi(\lambda_j)$  is also Poisson with parameter  $= \sum_i \lambda_i$  (i.e. sum of the individual parameters) the denominator is actually the distribution of

$$\sum_{j=1}^k X_j \Rightarrow Poi\left(\sum_{j=1}^k np_j\right) = Poi(n)$$


Now, as  $np_j$  goes to infinity for all  $j$ ; because if we take large samples each  $X_j$  can be approximated as normal with  $np_j$  comma  $np_j$ . This we have just shown that as  $\theta$  goes to infinity, Poisson converges to normal with  $\theta$  comma  $\theta$ ; here  $\theta$  is  $np_j$  for the  $j$ th class.

Now, since sum of independent Poisson distributions, what the  $i$ th one is Poisson with lambda  $i$ , is also a Poisson distribution with parameters

$$\sum_i \lambda_i$$

That is the sum of individual parameters, the denominator of the earlier distribution is equal to  $\sum X_j$  that is Poisson with parameter  $n$ , as we have discussed just now.

(Refer Slide Time: 18:56)

Therefore, the whole expression

$$\frac{\frac{\exp(-np_1)(np_1)^{n_1}}{n_1!} * \frac{\exp(-np_2)(np_2)^{n_2}}{n_2!} * \dots * \frac{\exp(-np_k)(np_k)^{n_k}}{n_k!}}{\frac{\exp(-(np_1 + np_2 + \dots + np_k))(np_1 + np_2 + \dots + np_k)^n}{n!}}$$

can be interpreted as:



$P(X_j = n_j)$  for  $j = 1, 2, \dots, k$ , subject to the condition that

$$\sum_{j=1}^k X_j = n$$

where  $X_j \sim \text{Poi}(np_j)$   
 i.e.

$$P(X_1 = n_1, X_2 = n_2, \dots, X_k = n_k \mid \sum_{j=1}^k X_j = n)$$

The above reduces the degrees of freedom from  $k$  to  $k-1$



Therefore, the entire expression can be interpreted as probability  $X_j$  is equal to  $n_j$ , for  $j$  is equal to 1 to  $k$ , subject to the condition that  $\sum X_j$  is equal to  $n$ ; because the sum of total observations has to be  $n$ ; where  $X_j$  is equal to Poisson of with  $np_j$ . That is we are looking at the probability of the event  $X_1$  is equal to  $n_1$ ,  $X_2$  is equal to  $n_2$ ,  $X_k$  is equal to  $n_k$ ; subject to the condition that  $\sum X_j$  is equal to  $n$ . Since all the  $k$  values are not independent, there is a constraint there; therefore the degrees of freedom is reduced from  $k$  to  $k - 1$ .

(Refer Slide Time: 19:52)

Since as  $n \rightarrow \infty$  each  $X_j$  may be approximated as  $N(np_j, np_j)$ .  
 Once the observed value of  $X_j$  is  $n_j$ , we can write  $n_j$  is distributed as  $N(np_j, np_j)$  subject to the condition that

Or

$$\sum_{j=1}^k n_j = n$$

$\frac{n_j - np_j}{\sqrt{np_j}} \sim N(0,1)$  such that  $\sum_{j=1}^k \left( \frac{n_j - np_j}{\sqrt{np_j}} \right) \sqrt{np_j} = 0$

This is because

a)  $\sum_{j=1}^k n_j = n$     b)  $\sum_{j=1}^k p_j = 1$     c)  $\sum_{j=1}^k np_j = n = \sum_{j=1}^k n_j$

Now, since  $n$  goes to  $\infty$ , each  $X_j$  may be approximated as normal  $np_j$  comma  $np_j$ . And therefore once the observed value of  $X_j$  is  $n_j$ , we can write  $n_j$  is distributed as

$$N(n_j, np_j)$$

subject to the condition that

$$\sum_{j=1}^k n_j = n$$

Therefore,  $n_j$  is the observed value,  $np_j$  is the expected value divided by square root of  $np_j$ . That is going to give a normal 0, 1 distribution, subject to the condition that

$$\sum_{j=1}^k \left( \frac{n_j - np_j}{\sqrt{np_j}} \right) \sqrt{np_j} = 0$$

This is obvious because if we cancel it, then this numerator will give me 0. Why it is happening? Because  $\sum n_j$  is equal to  $n$ ,  $\sum p_j$  is equal to 1, and  $\sum np_j$  is equal to  $n$  is equal to  $n_j$ ; where  $n_j$  is the observed value of the  $X_j$ .

(Refer Slide Time: 21:00)

Now consider the Statistic

$$\sum_{j=1}^k \frac{(n_j - np_j)^2}{np_j} \leftarrow \text{Sum of square of } k \text{ std Normal variables}$$

This may be approximated as the sum of square of  $k$  standard normal distribution subject to the linear constraint that with one constraint

$$\sum_{j=1}^k \left( \frac{n_j - np_j}{\sqrt{np_j}} \right) \sqrt{np_j} = 0$$

Hence the distribution of the above statistic is  $\chi^2$  with  $k - 1$  degrees of freedom.

Now, consider the statistic  $n_j$  minus  $np_j$  whole square divided by  $np_j$ . This maybe approximated as the sum of square of  $k$  standard normal distribution subject to the linear constraint that,  $n_j$  minus  $np_j$  upon root over  $np_j$ , multiplied by root over  $np_j$ , and summation over  $K$  is equal to 0. Hence, what is this? It is the sum of square of  $k$  standard normal variables with one constraint. Therefore, this is a Chi-Square distribution with  $k$  minus 1 degrees of freedom; because we know that if there are  $m$  standard normal distribution  $x_1, x_2, x_n$ . Then  $x_1$  square plus  $x_2$  square plus  $x_n$  square is chi square with  $m$  degrees of freedom.

(Refer Slide Time: 22:20)

We can write the above as:

$$\sum_{j=1}^k \frac{(O_j - E_j)^2}{E_j} \chi^2_{(k-1)}$$

Where

- $O_j$  is the observed frequency of class  $j$  i.e.  $n_j$
- $E_j$  is the expected frequency of class  $j$  i.e.  $np_j$

So, this same thing typically we write it as

$$\sum_{j=1}^k \frac{(O_j - E_j)^2}{E_j} \sim \chi_{(k-1)}^2$$

where  $O_j$  is the observed frequency of the class  $j$ . That is  $n_j$ , the number of observation coming from the class  $j$ , and  $E_j$  is the expected frequency of the class  $j$ ; which we have just seen to be  $np_j$ . Therefore, observed minus expected whole square divided by expected, this sum for  $j$  is equal to 1 to  $k$ ; that is the number of classes is chi square with  $k$  minus one degrees of freedom.

(Refer Slide Time: 23:10)

Example

Suppose a dice is thrown 600 times and the observed frequencies are:

Value ( $j$ )	1	2	3	4	5	6
Frequency ( $O_j$ )	95	102	93	107	115	88

We want to test whether the dice is unbiased by testing the null Hypothesis

$H_0: p_1 = p_2 = p_3 = p_4 = p_5 = p_6 = 1/6$



So, let us go back to an example, suppose a dice is thrown 600 times and the observed frequencies are like this: 1 appeared 95 times, 2 appeared 102 times, 3 appears 93 times, 4 appears 107 times, 5 115 times, 6 is 88 times. So, what we can see that there is some variation in the number of observations corresponding to the 6 different faces of the die. But, does it mean that the die is not unbiased, we are not sure; we need to test that. So, what we are testing is

$$H_0: p_1 = p_2 = p_3 = p_4 = p_5 = p_6 = 1/6$$

So, this is the null hypothesis that we are testing against the dice not unbiased.

(Refer Slide Time: 24:08)

**Solution** Here  $n = 600$  ✓  
Under  $H_0$ , for each  $j = 1, 2, \dots, 6$  ✓  
 $np_j = 600 \times \frac{1}{6} = 100$  ✓  
i.e. the expected frequency  $E_j = 100 \forall j = 1, 2, \dots, 6$  ✓  
The number of classes is  $k = 6$  ✓  
Thus the statistic  $Q = \sum_{j=1}^k \frac{(O_j - E_j)^2}{E_j} \sim \chi^2_{(5)}$  ✓  
The observed value of the statistic  $Q$  is:  
$$= \frac{(95-100)^2}{100} + \frac{(102-100)^2}{100} + \frac{(93-100)^2}{100} + \frac{(107-100)^2}{100}$$
$$+ \frac{(115-100)^2}{100} + \frac{(88-100)^2}{100}$$
$$Q = \frac{5^2 + 2^2 + 7^2 + 7^2 + 15^2 + 12^2}{100} = \frac{496}{100} = 4.96$$

*This is always an one-sided test.*

NPTEL

So, what we do? We do the following  $n$  is 600, there are 6 classes, under the null hypothesis  $np_j$  is equal to 100; that is expected number of observations coming from each class is 100 that is my  $E_j$ . Therefore, the statistic is going to be

$$\sum_{j=1}^k \frac{(O_j - E_j)^2}{E_j}$$

and we know that this is going to be chi square with 5 degrees of freedom.

Now, let us compute the value of the statistic, therefore it is

$$\frac{(95-100)^2}{100} + \frac{(102-100)^2}{100} + \frac{(93-100)^2}{100} + \frac{(107-100)^2}{100}$$
$$+ \frac{(115-100)^2}{100} + \frac{(88-100)^2}{100}$$

now, you understand how to compute the statistic. So, this value is coming out to be 5 square plus 2 square plus 7 square plus 7 square plus 15 square plus 12 square; which is coming out to be 496, divided by 100, that is 4.96.

Now, we shall reject the null hypothesis, if this value is very large; therefore this is always going to be one sided test. So, how it will work?

(Refer Slide Time: 25:53)

From the chi square table shown below ( taken from the link <https://link.springer.com/content/pdf/bbm%3A978-3-319-05555-8%2F1.pdf>)

Table 3 Values of  $\chi^2_{\alpha,df}$  in a chi-square distribution with  $df$  degrees of freedom (shaded area  $P(\chi^2 > \chi^2_{\alpha,df}) = \alpha$ )



df	$\alpha = .995$	$\alpha = .990$	$\alpha = .975$	$\alpha = .950$	$\alpha = .050$	$\alpha = .025$	$\alpha = .010$	$\alpha = .005$	df
1	0.0000993	0.000157	0.000982	0.00393	2.341	5.024	6.635	7.879	1
2	0.0100	0.0201	0.0506	0.103	5.991	7.378	9.210	10.597	2
3	0.0717	0.115	0.216	0.352	7.815	9.348	11.345	12.838	3
4	0.207	0.297	0.484	0.711	9.488	11.143	13.277	14.860	4
5	0.412	0.554	0.831	1.145	11.070	12.833	15.086	16.750	5
6	0.676	0.872	1.237	1.635	12.592	14.449	16.812	18.548	6

NPTEL

So, here is a chi square distribution, there are chi square table available; I have given you one link for different degrees of freedom. Here I have just copied for 1, 2, 3, 4, 5, 6 and the cutoffs are given, which are the critical values for different alpha. We are looking at 5 percent of level of significance that means we are looking at an  $\alpha$ , such that the probability the chi square distribution is taking a value greater than this  $C\alpha$  is  $\alpha$ . So, these  $C\alpha$  values are tabulated here, we have chi square with 5 degrees of freedom; we are looking at  $\alpha$  is equal to 0.05, and the  $C\alpha$  value is given out to be 11.07.

(Refer Slide Time: 26:54)



From the chi square table shown below ( taken from the link  
<https://link.springer.com/content/pdf/bbm%3A978-3-319-05555-8%2F1.pdf>

Critical value for  $df = 5$  and  $\alpha = 0.05$  is 11.07

Obtained  $Q = 4.96$

Thus,  $H_0$  is not rejected at 5% level of significance.

*$\therefore$  We accept  $H_0$  at 5% level of significance.*



Therefore, what I am writing here is that the degrees of freedom is 5, alpha is equal to 0.05, the critical value is 11.07. The obtained value is 4.96, therefore we cannot reject at 5 percent level of significance; therefore, we can accept the hypothesis at 5 percent level of significance.

(Refer Slide Time: 27:40)



A manufacturer produces an item and pack them in bags of 100. He claims that on the average in a bag 6% items will be defective.

To check the above claim a customer opens a bag and tests all the items. ✓

He finds that 10 items are defective.

Should he accept that on the average 6% will be defective?



Let me now give you a practical example, suppose a manufacturer produces an item and pack them in bags of 100. He claims that on the average in a bag 6% items will be defective. Now, to check the above claim a customer opens a bag and tests all the items; he finds that 10 of the items out of 100 are defective. Question comes whether he should accept that on the average 6% of the items are going to be defective?

That means that in another bag there maybe lesser number of defectives; like that when we keep on taking different bags of samples. We may end up with at the end that on the average, the number of defectives is going to be about 6% the total number of products.

(Refer Slide Time: 28:46)

Note:

The above problem is typical problem on decision making where the producer and customer have opposite interests:

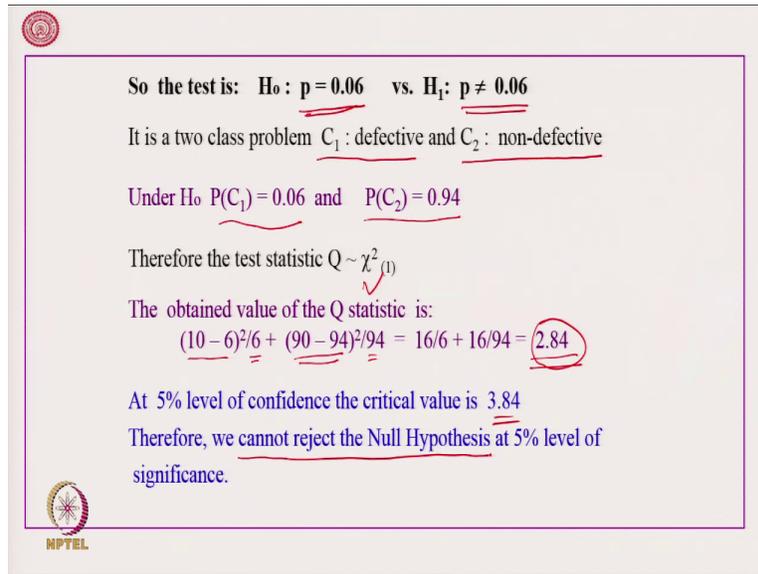
- Producer wants the hypothesis will be accepted – so that customer will go for placing purchase orders.
- Customer will be careful that he does not buy defective products. So will not be willing to accept consignments with proportion of defectives at a higher level.

Typically, the testing procedure is decided in advance before actual examination of the items starts.

Now, this is a very very typical problem on decision making; where the producer and the customer have opposite interests. A producer wants the hypothesis will be accepted, so that the customer will go for placing purchase orders. Because, if the purchaser is convinced that the number of defectives is small; then he is likely to go for purchasing the items from the producer. However, the customer or the purchaser will be careful that he does not buy defective products; so, will not be willing to accept consignments with proportion of defectives at a higher level.

Therefore, if the test fails that the number of defectives on the average is greater than 6% then customer may not like to purchase these items from the producer. So, typically how the test will be carried out that is decided in advance before the actual examination of the item starts. So, we are not going into the detail of this statistical procedures.

(Refer Slide Time: 30:00)



So the test is:  $H_0: p = 0.06$  vs.  $H_1: p \neq 0.06$

It is a two class problem  $C_1$ : defective and  $C_2$ : non-defective

Under  $H_0$   $P(C_1) = 0.06$  and  $P(C_2) = 0.94$

Therefore the test statistic  $Q \sim \chi^2_{(1)}$

The obtained value of the Q statistic is:

$$\frac{(10 - 6)^2}{6} + \frac{(90 - 94)^2}{94} = \frac{16}{6} + \frac{16}{94} = 2.84$$

At 5% level of confidence the critical value is 3.84

Therefore, we cannot reject the Null Hypothesis at 5% level of significance.

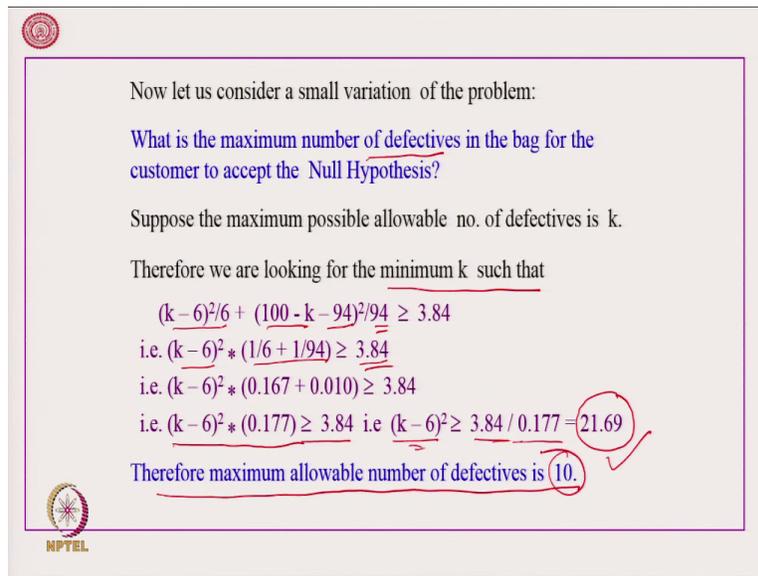


So, let us look at how we test for  $H_0: p = 0.06$  versus  $H_1: p \neq 0.06$ . So, we are looking at two sided tests. Therefore, it is a two class problems, where  $C_1$  is the class of defectives and  $C_2$  is the class of non-defectives. Therefore, under  $H_0$  probability of class 1 is equal to 0.06, and probability of class 2 is equal to 0.94. Therefore, the test statistic  $Q$  is chi square with 1 degrees of freedom. Why 1 degree? Because there are two classes. And we know that there are  $k$  classes then the chi square will have  $k-1$  degrees of freedom.

The obtained value of the  $Q$  statistic is therefore  $10$  minus  $6$  whole square divided by  $6$ , plus  $90$  minus  $94$  whole square divided by  $94$ ; which is coming out to be  $2.84$ . Therefore, now the problem is whether we should accept the null hypothesis or reject, for that we have to refer to the chi square table. And we see the following that for 1 degrees of freedom at 5% level of significance; the critical value is  $3.84$ .

We have taken the table from this link, so that you can also check how to see a chi square table. But, going back to the problem, we can see that our obtained value is somewhat smaller than the tabulated value of  $C\alpha$ ; when  $\alpha = 0.05$ . Therefore, we cannot reject the null hypothesis at 5% level of significance.

(Refer Slide Time: 32:10)



Now let us consider a small variation of the problem:

What is the maximum number of defectives in the bag for the customer to accept the Null Hypothesis?

Suppose the maximum possible allowable no. of defectives is  $k$ .

Therefore we are looking for the minimum  $k$  such that

$$(k-6)^2/6 + (100-k-94)^2/94 \geq 3.84$$

i.e.  $(k-6)^2 * (1/6 + 1/94) \geq 3.84$

i.e.  $(k-6)^2 * (0.167 + 0.010) \geq 3.84$

i.e.  $(k-6)^2 * (0.177) \geq 3.84$  i.e.  $(k-6)^2 \geq 3.84 / 0.177 = 21.69$

Therefore maximum allowable number of defectives is 10.

NPTEL

Now, let us consider a small variation of the problem. Suppose we want to ask, what is the maximum number of defectives in the bag, for the customer to accept the null hypothesis? That means that if the number of defective is less than that value say  $k$ , then the customer is going to accept the null hypothesis; otherwise he is going to reject the null hypothesis. So, what is going to be the maximum number of defectives that is going to be the question?

So, we are therefore looking at the minimum value of  $k$ , such that

$$(k-6)^2/6 + (100-k-94)^2/94$$

That has to be greater than equal to 3.84, if we have to reject the null hypothesis. That is

$$(k-6)^2 * (1/6 + 1/94)$$

has to be greater than equal to 3.84. Now, that boils down to  $k$  minus 6 whole square multiplied by 0.1 double seven has to be greater than equal to 3.84.

That is

$$(k-6)^2 \geq 3.84 / 0.177 :$$

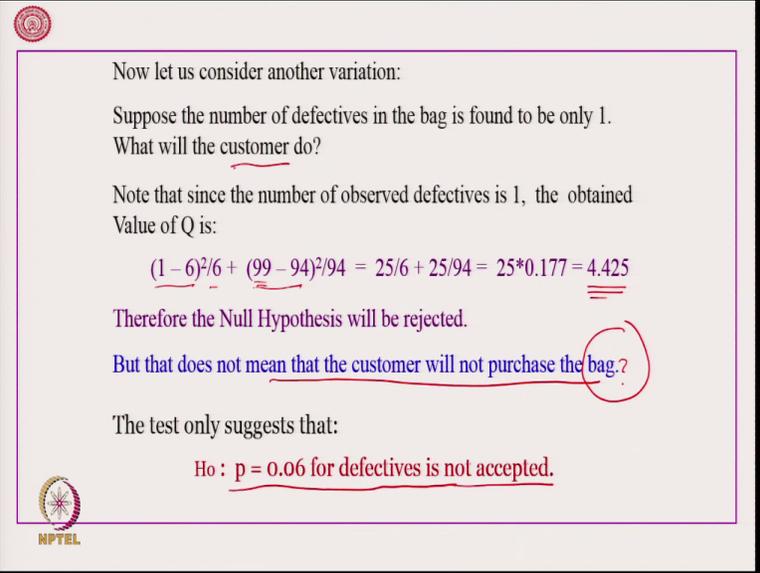
which is coming out to be 21.69. Therefore, we shall reject the null hypothesis, if  $k$  is such that

$$(k - 6)^2 \geq 3.84 / 0.177 = 21.69$$

Therefore we can easily see that the value of  $k$  can be 11. Hence, we can see that the maximum allowable number of defective is 10. Because, in that case  $k$  minus 6 whole square is coming out to be 4 square that is 16, which is less than 221.69 and we are going to accept the null hypothesis.

But, if  $k$  is equal to 11, then this is coming out to be 25 and therefore we are going to reject the null hypothesis; therefore, the maximum allowable number of defective is 10.

(Refer Slide Time: 34:43)



Now let us consider another variation:  
Suppose the number of defectives in the bag is found to be only 1.  
What will the customer do?

Note that since the number of observed defectives is 1, the obtained Value of Q is:

$$\frac{(1 - 6)^2}{6} + \frac{(99 - 94)^2}{94} = \frac{25}{6} + \frac{25}{94} = 25 * 0.177 = 4.425$$

Therefore the Null Hypothesis will be rejected.

But that does not mean that the customer will not purchase the bag?

The test only suggests that:

Ho :  $p = 0.06$  for defectives is not accepted.

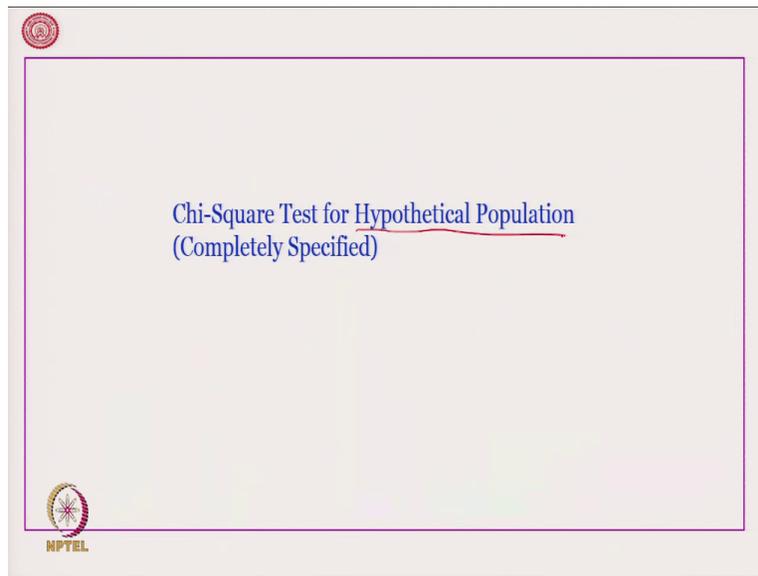
The slide includes a copyright symbol in the top left and the NPTEL logo in the bottom left. Handwritten annotations include a red circle around the question 'But that does not mean that the customer will not purchase the bag?' and a red underline under the null hypothesis statement.

Now, let us consider another variation: Suppose the number of defectives in the bag is found to be only 1. Question is what will the customer do? Note that since the number of observed defective is 1; the obtained value of  $Q$  is going to be 1 minus 6 whole square by 6. Then 99 minus 94 whole square divided by 94, this is coming out to be 99; because since there is only one defective, 99 of the other elements or items are non-defective. This value is coming out to be 4.425, which is greater than the  $\alpha$  that is 3.84 which we have seen sometime back.

What does it mean? Does it mean that the customer will not purchase the bag because we are rejecting the null hypothesis? No. The customer will go for purchasing the bag because number of defective has been found to be very small. Then what does it mean that we are going to reject

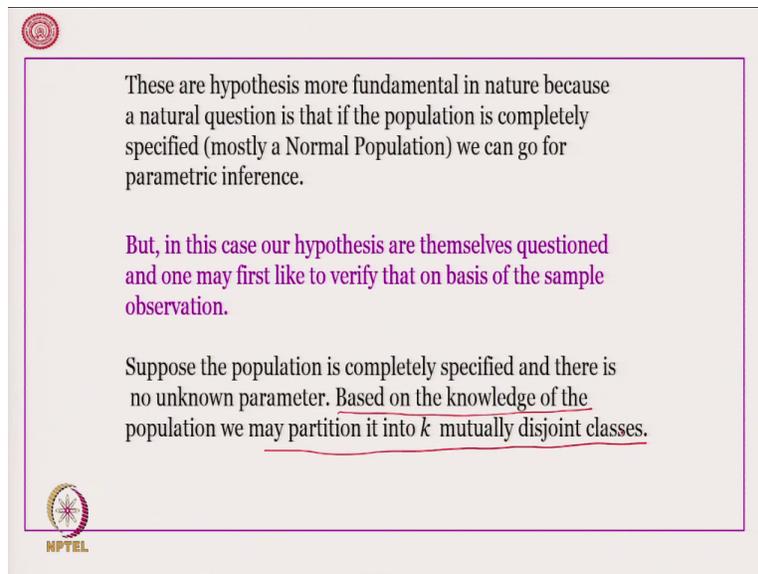
the null hypothesis; it means that we are rejecting that  $p$  is equal to 0.06. That is the proportion of defectives is not accepted, because the test suggest that the proportion of defective is quite different from 6 percent.

(Refer Slide Time: 36:17)



Now, we can use Chi-Square test in a slightly different way as follows. What I am saying Chi-Square test for hypothetical population which is completely specified.

(Refer Slide Time: 36:34)



So, these are hypothesis more fundamental in nature because a natural question is that, if the population is completely specified; mostly is a normal population, then we can go for parametric inference. So, what does it mean? I stated at the very beginning that parametric inference is more robust or powerful than nonparametric; but, parametric inference works under normality assumption.

So, if we are given a sample and we can test, that it is actually coming from a normal population, then we are really happy because we can use all our parametric inference techniques. So, that is why sometimes we try to check if the given sample is coming from a standard distribution. And that is why I am saying that the hypothesis that it is coming from a normal population itself can be questioned, and we need to make sure that is correct.

That is that assumption is correct and therefore we like to verify that on the basis of a sample observation. So, how to do that? Suppose the population is completely specified and there is no unknown parameter. We know that if there is no unknown parameter, then the density is completely known to us; and therefore the distribution is known to us. But, based on that knowledge what we do? We try to partition the entire sample space into  $k$  mutually disjoint classes.

(Refer Slide Time: 39:04)

Let according to the underlying assumption i.e.  $H_0$  the proportion of the population that falls under the  $j^{\text{th}}$  class be  $p_j^0$ .

Then the number of observation falling in the  $j^{\text{th}}$  class is expected to be  $np_j^0$ .

But actual number of observations in the  $j^{\text{th}}$  class may be somewhat different.  $n_j$ .

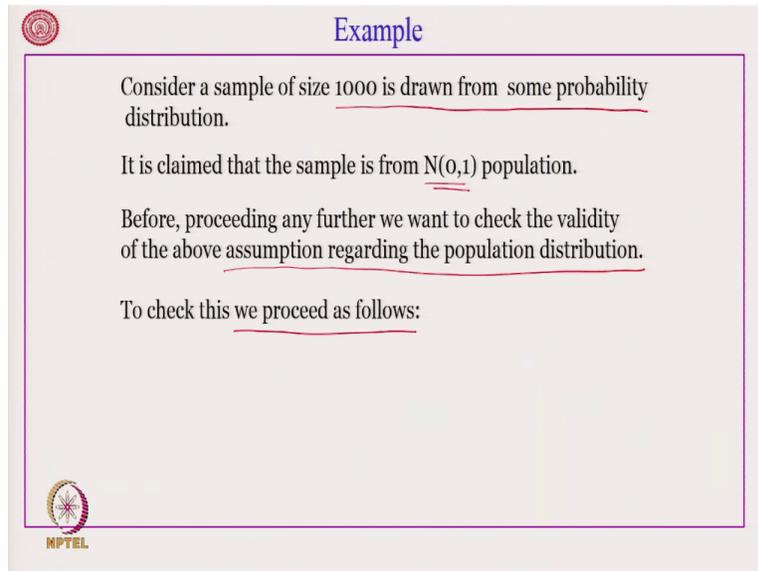
The sum of the square of deviation normalized by the expected value gives the desired Chi-sq statistic.

The slide includes the NPTEL logo in the bottom left corner and navigation icons at the bottom.

So, let according to the underlying assumption that is  $H_0$  is the proportion of population that falls under the  $j^{\text{th}}$  class. Then the number of observations falling in the  $j^{\text{th}}$  class is expected to be  $np_j^0$ . Because this is the population proportion,  $n$  is the sample therefore, expected number of sample in the  $j^{\text{th}}$  class is  $np_j^0$ . But, actual number of observations in the  $j^{\text{th}}$  class maybe somewhat different as we have already seen.

Thus, there may be a gap between the actual observations from the  $j$ th class, which if we call  $n_j$ ; between  $n_j$  and  $np_j^0$ . The sum of the square of the deviation normalized by the expected value gives the desired Chi-Square statistic. This we have already explained.

(Refer Slide Time: 40:02)



The slide is titled "Example" in blue text at the top center. It contains four lines of text, each underlined in red. The text reads: "Consider a sample of size 1000 is drawn from some probability distribution." followed by "It is claimed that the sample is from  $N(0,1)$  population." then "Before, proceeding any further we want to check the validity of the above assumption regarding the population distribution." and finally "To check this we proceed as follows:" The slide also features a small circular logo in the top left and bottom left corners, and the text "NPTEL" in the bottom left corner.

So, to consider a sample of size 1000 that is drawn from a probability distribution. We claimed that the sample is coming from  $N(0, 1)$  population; but we need to verify it. Therefore, before proceeding any further we want to check the validity of the above assumption regarding the population distribution. So, how to do that that is the question; so we proceed as follows.

(Refer Slide Time: 40:33)

We decide to partition the data into 7 classes.

Then from the Normal Table we get the cumulative distribution of the classes.

Suppose we partition as follows:

$(-\infty, -2.5), [-2.5, -1.5), [-1.5, -0.5), [-0.5, +0.5),$   
 $[+0.5, +1.5), [+1.5, +2.5), [+2.5, \infty)$



What we are doing? We decide to partition the data into 7 classes; then from normal table we get the cumulative distribution of the classes. Suppose we partition it as follows:-  $-\infty$  to  $-2.5$ ,  $-2.5$  to  $-1.5$ ,  $-1.5$  to  $-0.5$ ,  $-0.5$  to  $+0.5$ ; similarly,  $+0.5$  to  $1.5$ ,  $1.5$  to  $2.5$  and  $2.5$  to  $\infty$ . Thus, there are 7 classes.

(Refer Slide Time: 41:13)

Class #	Class Range	Expected Class Probability ( $p_j^e$ )	No. of Observations ( $O_j$ )	Expected Class Frequency ( $E_j$ )	$(O_j - E_j)^2$	$\frac{(O_j - E_j)^2}{E_j}$
1	$x < -2.5$	0.0062	10	6.2	14.44	2.32
2	$-2.5 \leq x < -1.5$	0.0606	80	60.6	376.36	6.21
3	$-1.5 \leq x < -0.5$	0.2417	220	241.7	470.89	1.95
4	$-0.5 \leq x < +0.5$	0.3829	400	382.9	292.41	0.76
5	$+0.5 \leq x < +1.5$	0.2417	195	241.7	2180.89	9.02
6	$+1.5 \leq x < +2.5$	0.0606	60	60.6	0.36	0.01
7	$+2.5 \leq x$	0.0062	35	6.2	829.44	133.78
Total			1000			154.05





z	0.0	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
1.5	0.9332	0.9345	0.9357	0.9370	0.9382	0.9394	0.9406	0.9418	0.9429	0.9441
1.6	0.9452	0.9463	0.9474	0.9484	0.9495	0.9505	0.9515	0.9525	0.9535	0.9545
1.7	0.9554	0.9564	0.9573	0.9582	0.9591	0.9599	0.9608	0.9616	0.9625	0.9633
1.8	0.9641	0.9649	0.9656	0.9664	0.9671	0.9678	0.9686	0.9693	0.9699	0.9706
1.9	0.9713	0.9719	0.9726	0.9732	0.9738	0.9744	0.9750	0.9756	0.9761	0.9767
2.0	0.9772	0.9778	0.9783	0.9788	0.9793	0.9798	0.9803	0.9808	0.9812	0.9817
2.1	0.9821	0.9826	0.9830	0.9834	0.9838	0.9842	0.9846	0.9850	0.9854	0.9857
2.2	0.9861	0.9864	0.9868	0.9871	0.9875	0.9878	0.9881	0.9884	0.9887	0.9890
2.3	0.9893	0.9896	0.9898	0.9901	0.9904	0.9906	0.9909	0.9911	0.9913	0.9916
2.4	0.9918	0.9920	0.9922	0.9925	0.9927	0.9929	0.9931	0.9932	0.9934	0.9936
2.5	0.9938	0.9940	0.9941	0.9943	0.9945	0.9946	0.9948	0.9949	0.9951	0.9952
2.6	0.9953	0.9955	0.9956	0.9957	0.9959	0.9960	0.9961	0.9962	0.9963	0.9964
2.7	0.9965	0.9966	0.9967	0.9968	0.9969	0.9970	0.9971	0.9972	0.9973	0.9974
2.8	0.9974	0.9975	0.9976	0.9977	0.9977	0.9978	0.9979	0.9979	0.9980	0.9981
2.9	0.9981	0.9982	0.9982	0.9983	0.9984	0.9984	0.9985	0.9985	0.9986	0.9986
3.0	0.9987	0.9987	0.9987	0.9988	0.9988	0.9989	0.9989	0.9989	0.9990	0.9990



So, these are the 7 class;  $x$  is less than minus 2.5,  $x$  is between minus 2.5 to minus 1.5 up to  $x$  is greater than plus 2.5. Now, the question is what is the expected class probability  $P_j$ ? So,  $x$  less than minus 2.5 is 0.0062. How to get that? We get it from the normal table. So, this is minus 2.5, this is 0. So, this is the cumulative probability of a standard normal distribution taking value less than minus 2.5 that is 0.0062. Therefore, that we have obtained. Now, let me do it one or two more for your understanding. What is the probability that a standard normal random variable is between minus 2.5 and minus 1.5.

So, again we go back less than minus 1.5 is 0.0668 0.0668 and less than minus 2.5 is 0.0062. Therefore, we get the difference is 0.0606. So, the probability is 0.0606 and that is what we have written here. In a similar way we can compute for all I have given the entire normal table for your help. You can compute from here the class probabilities for all the 7 classes which I am showing here, like this in this column. Suppose the actual observations belonging to this 7 classes are as follows: 10, 80, 220, 400, 195, 60 and 35.

Therefore expected class frequency how do we get? We multiply these values by 1000; therefore we get 6.2, 60.6, 241.7 like that up to 6.2, and that should sum up to 1000. So, this is the expected class frequency, this is the observed values. Therefore we need to compute the square of their difference, which is coming out to be like this 14.44, 376.36 up to 829.44. Therefore, for each one of them we calculate this normalized thing by dividing by  $E_j$ , and these are the values that we obtained.

And the sum of this values

$$\sum_j \frac{(O_j - E_j)^2}{E_j}$$

summing over j is equal to 1 to 7, we get 154.05. Is it too high? We do not know but looking at it; I can itself see that this cannot be from normal 0, 1. Because x greater than 2.5 that should number should be very small; but here we find that number is pretty large which is 35. And if we look at this sum of square, we can see that that is giving us a huge value, making the sum to be too high. But, still it can be from normal 0, 1.

(Refer Slide Time: 45:15)

Therefore value of the Statistic is 154.05

The critical value for  $\chi^2_{(6)}$  is

Cr value for df = 6 at 0.01 level is 16.812.  
Hence reject the  $H_0$

df	$\alpha = .995$	$\alpha = .990$	$\alpha = .975$	$\alpha = .950$	$\alpha = .900$	$\alpha = .850$	$\alpha = .800$	$\alpha = .750$	$\alpha = .700$	$\alpha = .650$	$\alpha = .600$	$\alpha = .550$	$\alpha = .500$	$\alpha = .450$	$\alpha = .400$	$\alpha = .350$	$\alpha = .300$	$\alpha = .250$	$\alpha = .200$	$\alpha = .150$	$\alpha = .100$	$\alpha = .050$	$\alpha = .025$	$\alpha = .010$	$\alpha = .005$	df
1	0.0000393	0.000157	0.000982	0.00393	3.841	5.024	6.635	7.879	1																	
2	0.0100	0.0201	0.0506	0.103	5.991	7.378	9.210	10.597	2																	
3	0.0717	0.115	0.216	0.352	7.815	9.348	11.345	12.838	3																	
4	0.207	0.297	0.484	0.711	9.488	11.143	13.277	14.860	4																	
5	0.412	0.554	0.831	1.145	11.070	12.833	15.086	16.750	5																	
6	0.676	0.872	1.237	1.635	12.592	14.449	16.812	18.548	6																	

So, let us check with the chi square value; so, again we look at the chi square table. This is going to be chi square with 6 degrees of freedom; because there are 7 classes. And therefore k is equal to 7 and therefore k minus 1 is equal to 6 is the degrees of freedom; therefore, we look at chi square with 6 and for 0.05, the critical value is 12.592. Therefore, so first we look at what is the critical value at 1% level of confidence that is 16.812 and at 5 % level of confidence, it is 12.592.

Therefore, since for degrees of freedom at 6 at 1 percent level; the value is 16.812. Therefore, we reject the null hypothesis at 1% level of significance; therefore quite naturally it has to be

rejected at 5% level of significance as well. Or, in other words we are not going to accept that this data is coming from a standard normal distribution.

(Refer Slide Time: 46:46)



If the  $n$  sample observations are values of a continuous r.v. one may expect that the final decision may differ with different values of  $k$  and the boundaries chosen to demarcate the class.

Hence the question is that can all the data points be taken into consideration to test any such hypothesis.

In case of continuous distribution comparisons can be made between observed and expected cumulative relative frequencies for each of the observed values.

Several goodness of fit tests statistics are functions of the deviations between the observed cumulative distribution and corresponding cumulative probabilities expected under  $H_0$ .



NPTEL



However, one can think of that if  $n$  is the sample observations and they are coming from a continuous random variable. One can expect that the final decision may differ with the different values of  $k$  and the boundaries chosen to demarcate the class. I have just given you an example with only one way of choosing the class boundaries, and the number of classes is equal to 7. One can think of making different class boundaries, different number of classes and test. I will leave it as an exercise for you to check with a few different classes and computing the corresponding chi square statistic.

Therefore, one important problem is that here we have considered elements belonging to a particular class. Therefore, the points are losing their individuality, can we make a test where we look all data points into consideration to test such a hypothesis. In case of continuous distribution comparisons can be made between observed and expected cumulative relative frequencies, for each of the observed values.

So, instead of putting them into several classes, we look at each point in a sorted order, in an increasing order and check up to that point what is the observed frequency, and what could have been the actual frequency, if you are looking at certain underlying distribution. And then we try to see how much difference it is actually generating; so let us go and illustrate. So, several goodness of fit tests statistics are functions of the deviations between the observed cumulative distribution and the corresponding cumulative probabilities expected under  $H_0$ .

(Refer Slide Time: 49:00)

The test can be used for an underlying distribution which can be continuous.

We design an Empirical distribution function  $F_n$  based on the  $n$  samples taken from the population.

$$F_n(x) = \begin{cases} 0 & \text{if } x < X_{(1)} \\ \frac{i}{n} & \text{if } X_{(i)} \leq x < X_{(i+1)}, i = 1, 2, \dots, n-1 \\ 1 & \text{if } x \geq X_{(n)} \end{cases}$$

Here,  $X_{(i)}$  is the  $i^{\text{th}}$  order statistics.

$F_n(x)$  is the proportion of sample values  $\leq x$

$\therefore nF_n(x)$  is the number of samples values  $\leq x$

Possible values for  $nF_n(x)$  are: 0, 1, 2, ... n

*Handwritten notes:*  $X_{(1)} < X_{(2)} < \dots < X_{(n)}$  and  $X_{(i)} < X_{(i+1)}$

So, we shall study one major algorithm in this regard which is called Kolmogorov-Smirnov one sample test. The test goes as follows: we design an Empirical distribution function  $F_n$  based on the  $n$  samples taken from the population. So, we have the ordered as I have said that we have to take them in sorted order; so,  $x_1$  less than  $x_2$  less than  $x_n$  is the ordered sample. What does it mean? Below this point there was no observation; that is  $F_n(x)$  is equal to 0, if  $x$  is less than  $x_1$ . And  $x_n$  has to be 1, if  $x$  is greater than  $x_n$ .

So, for any  $x$  more than this we know that the entire sample is less than that. Therefore, the cumulative Empirical distribution function  $F_n$  will take the value 1. In between for  $x_i$  and  $x_{i+1}$ , if these are the observed statistic; any  $x$  between this interval means that there are  $i$  many observations out of  $n$  less than this. Therefore, what we are getting the cumulative probability to be  $i/n$ . So, that is how we design the Empirical distribution function  $F_n$ , where  $x_i$  is the  $i^{\text{th}}$  order statistic.

$F_n(x)$  is the proportion of sample values less than equal to  $x$ ,  $nF_n(x)$  is the number of sample value less than equal to  $x$ , possible values for  $nF_n(x)$  are 0, 1, 2 up to  $n$ . As you can understand that if I multiply this by  $n$ , we get that number of sample values less than equal to  $x$ .

(Refer Slide Time: 51:02)

**Distribution of  $nF_n(x)$**

The distribution of  $nF_n$  is  $\text{Bin}(n, \theta)$  where  $\theta = F(x)$ , i.e.

Distribution function of  $X$   $\Rightarrow$  *Distrib under  $H_0$  i.e. cdf =  $F(x)$*

$$P(nF_n(x) = k) = \binom{n}{k} F(x)^k (1 - F(x))^{n-k}, k = 0, 1, 2, \dots, n$$

Thus

$$P\left(F_n(x) = \frac{k}{n}\right) = \binom{n}{k} F(x)^k (1 - F(x))^{n-k}, k = 0, 1, 2, \dots, n$$

Hence we have the following:

$$E(nF_n(x)) = n\theta = nF(x)$$

$$\text{Var}(nF_n(x)) = n\theta(1 - \theta) = nF(x)(1 - F(x))$$

*Binomial Dist<sup>n</sup>*



What is the distribution of  $nF_n(x)$ ? The distribution of  $nF_n$  is binomial with  $n$  theta; where theta is equal to  $F_x$ . What is  $F_x$ ?  $F_x$  is the distribution under null hypothesis; that is the cdf is equal to  $F_x$ . Therefore, under the null hypothesis, the distribution function  $F_x$  means that probability and observed value is less than equal to  $x$  is  $F_x$ . Therefore we are considering it to be a binomial distribution with  $F_x$  is the probability of falling in this side; and  $1 - f_x$  is the probability that it will fall on the other side of  $x$ .

Therefore, probability  $nF_n(x)$  is equal to  $k$  is equal to out of  $n$  we are choosing  $k$  that is  ${}^n C_k$ ; they are getting the probability  $F_x$  to the power  $k$  because they are less than equal to  $x$ .  $1 - F_x$  to the power  $n - k$  because they are falling here; where  $k$  is equal to  $0, 1, 2$  up to  $n$ . Therefore, probability  $F_n(x)$  is equal to  $k$  by  $n$ , I am just dividing both side by  $n$ ; is equal to  ${}^n C_k F_x$  to the power  $k$ ,  $1 - F_x$  to the power  $n - k$ .

So, if that is a binomial distribution, then expected value of  $nF_n(x)$  is equal to  $n\theta$  is equal to  $nF_x$ ; where  $F_x$  is the distribution. We are comparing with and its variance is going to be  $nF_x(1 - F_x)$ . This is coming from binomial distribution.

(Refer Slide Time: 53:05)



In One sample Kolmogorov-Smirnov test we check if the Unknown population distribution  $F(x)$  is actually a known distribution  $F_0(x)$ .

This is written as:

$$H_0: F(x) = F_0(x) \quad \checkmark$$

vs.

$$H_1: F(x) \neq F_0(x) .$$

The test statistic for this purpose is the maximum point-wise deviation between  $F_n(x)$  and  $F_0(x), \forall x \in \mathbb{R}$



What does it imply? It implies that the expected value of  $F_n(x)$  is equal to  $F(x)$ , and variance of  $F_n(x)$  is equal to  $F(x)(1 - F(x))$  by  $n$ . In one sample Kolmogorov-Smirnov test, we check if the unknown population distribution  $F(x)$  is actually a known distribution  $F_0(x)$ ; that is you are checking  $F(x) = F_0(x)$  versus  $F(x) \neq F_0(x)$ . This test statistic for this purpose is the maximum point wise deviation between  $F_n(x)$  and  $F_0(x)$ ; for all  $x$  belonging to  $\mathbb{R}$ .

(Refer Slide Time: 53:52)



Hence the test statistic is :

$$D_n = \max_{-\infty < x < \infty} |F_n(x) - F_0(x)|$$

The statistic is then compared with critical values tabulated for different values of  $n$ , from the Kolmogorov -Smirnov one sample statistic table as the one given in the link

<http://www.real-statistics.com/statistics-tables/kolmogorov-smirnov-table/>



Therefore, the test statistic is going to be  $D_n$  is equal to maximum of  $F_n(x) - F_0(x)$ ; for  $x$  between minus infinity to infinity. This statistic is then compared with critical values of tabulated

for different values of  $n$ , from the Kolmogorov-Smirnov one sample statistic table as the one given in the following link. So, I am giving you here a link to check Kolmogorov-Smirnov cut off values or critical values for one sample test.

(Refer Slide Time: 54:31)

**Example**

Suppose 10 values are chosen from the interval  $(0,1)$ . We want to test if the chosen values justify that they are picked from uniform distribution over the interval  $(0,1)$  i.e.  $U[0,1]$

Let the values in the sorted order be:

0.08 0.16 0.21 0.29 0.35 0.43 0.58 0.61 0.68 0.85 ✓

Note that in this case

$F_0(x) = x, x \in (0,1), F_0(x) = 0, x \leq 0$  and  $F_0(x) = 1, x \geq 1$

$U(0,1)$



So, let me illustrate it with an example, suppose we have taken ten observations from uniform 0, 1 distribution. Let these values be like this 0.08, 0.16, 0.21, up to 0.85. Therefore, what  $F_0(x)$  is uniform 0, 1; therefore  $F_0(x)$  is equal to  $x$ , for  $x$  belonging to  $(0, 1)$ . Because uniform 0, 1 distribution that means it is ranging only from 0 to 1 that means it takes values only from 0 to 1.

What is  $F_0(x) = 0$ , for  $x \leq 0$ ; and  $F_0(x) = 1$ , for  $x \geq 1$ . Therefore,  $F_0(x)$  will have a following shape it is 0 up to 0 and then it goes uniformly like this. At 1, 1 and then it takes the, it becomes horizontal. So, that is how this uniform distribution will look like.

(Refer Slide Time: 55:47)

Hence we compute sample cumulative frequency at each point  $x$  and compare that value with the  $F_0(x)$  and supremum of the differences will be obtained.

The maximum of these will be the K-S statistics

Hence, we compute the sample cumulative frequency at each point  $x$ , and compare that value with the  $F_0(x)$  and we consider the supremum of the differences and that is what is going to give me the Kolmogorov-Smirnov of statistic.

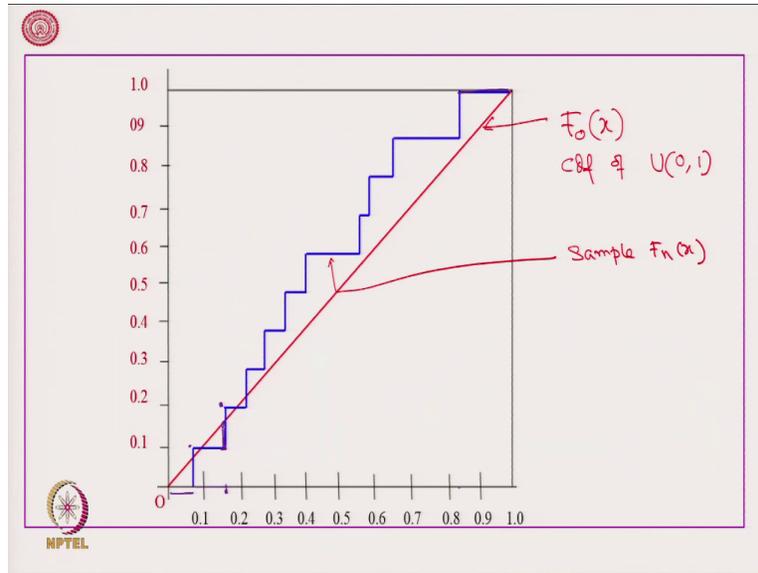
(Refer Slide Time: 56:05)

	$F_n(x)$	$\sup  F_n(x) - F_0(x) $
$x < 0.08$ ✓	0 ✓	0.08
$0.08 \leq x < 0.16$ ✓	$\frac{1}{10} = 0.1$ ✓	0.06
$0.16 \leq x < 0.21$ ✓	$\frac{2}{10} = 0.2$ ✓	0.04
$0.21 \leq x < 0.29$ ✓	$\frac{3}{10} = 0.3$	0.09
$0.29 \leq x < 0.35$	$\frac{4}{10} = 0.4$	0.11
$0.35 \leq x < 0.43$	$\frac{5}{10} = 0.5$	0.15
$0.43 \leq x < 0.58$	$\frac{6}{10} = 0.6$	0.17
$0.58 \leq x < 0.61$	$\frac{7}{10} = 0.7$	0.12
$0.61 \leq x < 0.68$	$\frac{8}{10} = 0.8$	0.19
$0.68 \leq x < 0.85$ ✓	$\frac{9}{10} = 0.9$ ✓	0.22
$x \geq 0.85$ ✓	1 ✓	0.15

So, what we are doing we had these values, so  $x$  less than 0.08; we did not get any observation. So, therefore that is 0 for between 0.8 to less than 0.16; we have got only one observation 0.08. Therefore,  $F_n(x)$  is equal to 0.1, point 16 to just less than point 21; we have got two observations point eight and point one six therefore,  $F_n(x)$  is equal to point 2. Similarly, for up to point 85, we

got 9 observations and at tenth observation in the ordered sample is point 85. Therefore, for  $x$  greater than equal to point 85, we get the value 1. Now, what is the value of the supremum?

(Refer Slide Time: 57:03)



To do that let us first look at a graph. So, this is the graph that we have drawn for our understanding. So, this red line is the  $F_0(x)$  that is the cdf of uniform 0, 1. The blue line which is the sample  $F_n(x)$ . So, what is happening till point 8 it is 0, at point 8, we get one observation; therefore it goes up to point 1. Then the next observation is at point 16, therefore it gets a jump up to point 2 and like that the last one is at point 85 and that gives a jump up to this point that is 1.

Therefore, for each of this interval we can always look at what is the maximum of the difference for any  $x$  belonging to that interval. Say for example, in this interval the maximum difference coming out between the red line and the blue line is this much; that is the supremum.

(Refer Slide Time: 58:24)

	$F_n(x)$	$\sup  F_n(x) - F_0(x) $
$x < 0.08$	0	0.08
$0.08 \leq x < 0.16$	$\frac{1}{10} = 0.1$	0.06
$0.16 \leq x < 0.21$	$\frac{2}{10} = 0.2$	0.04
$0.21 \leq x < 0.29$	$\frac{3}{10} = 0.3$	0.09
$0.29 \leq x < 0.35$	$\frac{4}{10} = 0.4$	0.11
$0.35 \leq x < 0.43$	$\frac{5}{10} = 0.5$	0.15
$0.43 \leq x < 0.58$	$\frac{6}{10} = 0.6$	0.17
$0.58 \leq x < 0.61$	$\frac{7}{10} = 0.7$	0.12
$0.61 \leq x < 0.68$	$\frac{8}{10} = 0.8$	0.19
$0.68 \leq x < 0.85$	$\frac{9}{10} = 0.9$	0.22
$x \geq 0.85$	1	0.15

Compare this value 0.22 with the critical value taken from the K-S Table

Therefore, we have computed those values when  $F_n(x)$  is 0 and  $x$  is less than 0.8; then the supremum is coming out to be 0.08. Because  $F_0(x)$  at 0.08 is equal to the same value  $F_0(x)$  is equal to  $x$ . Similarly, for point 08 to point 16,  $F_n(x)$  is point 1; therefore the supremum is coming out to be the difference between point 1 and point 16 that is point 016. Point 16 to point 21, the  $F_n(x)$  is 0.2; therefore the maximum difference is coming out is from point 16 and point 2 that is 0.04. Point 21 to point 29, the value is point 3; therefore the supremum of the difference is coming out to be point 21.

0.3 – 0.21 that is 0.09, like that we have computed all of them and we have got the maximum value to be 0.22. Now, if we look at the graph, we will find that the maximum value is coming here; this is the maximum difference which is coming out to be 0.22. So, the statistic for us is going to be that 0.22, therefore we need to compare this value with the critical value; that we can take from the Kolmogorov-Smirnov of table.

(Refer Slide Time: 60:08)

<http://www.real-statistics.com/statistics-tables/kolmogorov-smirnov-table/>

$n \backslash \alpha$	0.001	0.01	0.02	0.05	0.1	0.15	0.2
1		0.99500	0.99000	0.97500	0.95000	0.92500	0.90000
2	0.97764	0.92930	0.90000	0.84189	0.77639	0.72614	0.68377
3	0.92063	0.82900	0.78456	0.70760	0.63604	0.58582	0.56481
4	0.85046	0.73421	0.68887	0.62394	0.56522	0.52476	0.49265
5	0.78137	0.66855	0.62718	0.56327	0.50945	0.47439	0.44697
6	0.72479	0.61660	0.57741	0.51926	0.46799	0.43526	0.41035
7	0.67930	0.57580	0.53844	0.48343	0.43607	0.40497	0.38145
8	0.64098	0.54180	0.50654	0.45427	0.40962	0.38062	0.35828
9	0.60846	0.51330	0.47960	0.43001	0.38746	0.36006	0.33907
10	0.58042	0.48895	0.45662	0.40923	0.36866	0.34250	0.32257

Since the critical value at 5% level of significance is 0.409 we accept  $H_0$  at 5% level of significance. ✓

Since the critical value at 10% level of significance is 0.369 we do not reject  $H_0$  at 10% level of significance. ✓



So, let us look at the table, we have value of  $n$  is equal to 10 at 5 percent level of significance; it is 0.409, we are looking at only 3 decimal spaces. At point 1, we are getting 0.369; so we are comparing our obtained value with these values. What we have obtained is 0.22 therefore, we accept the  $H_0$  at 5% level of significance; and also we cannot reject it at 10% level of significance. That means we are going to accept it also at 10 percent level of significance.

(Refer Slide Time: 60:52)

A natural extension of the above test is: to check if two Samples  $X = \{x_1, x_2, \dots, x_m\}$  and  $Y = \{y_1, y_2, \dots, y_n\}$  are from the same distribution.  $\Rightarrow$

We shall consider that in Lecture 8.

Before that in Lecture 7 we shall see some other two sample Problems Such as Randomness.



Now, a natural extension of the above test is to check, if we take 2 different samples  $x_1, x_2, \dots, x_n$ ; and  $y_1, y_2, \dots, y_n$  are they coming from the same distribution. In that case what is going to happen?

We are not going to compare the sample frequency distribution with a standard hypothetical distribution. Rather we shall look at their sample cumulative functions for both  $x$  and  $y$  and we will check how much they differ. So, if the maximum value is greater than some threshold, then we are going to reject that they are coming from the same distribution; otherwise we are going to accept that.

So, that we will consider in lecture 8, but before that in the next lecture we shall see some other two sample tests. In particular our focus is going to be on Randomness testing; that is whether the sample taken by us is actually random. Okay friends I stop here today, in the next class we shall do these problems. Till then thank you.