

**Statistical Methods for Scientists and Engineers**  
**Prof. Somesh Kumar**  
**Department of Mathematics**  
**Indian Institute of Technology - Kharagpur**

**Lecture - 34**  
**Nonparametric Methods - VII**

Now, we move on to the next problem that is called the goodness of fit test, in this problem we want to test whether the distribution is a particular distribution or not. So basically this is the problem of modelling of distributions.

**(Refer Slide Time: 00:45)**

Lecture 34

© CET  
I.I.T. KGP

Goodness of Fit Tests      Chi-square Test (Karl Pearson)

Let  $X_1, \dots, X_n$  be a random sample from  $F(x)$

We want to test

$H_0: F(x) = F_0(x)$        $X \sim F(x)$

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

We consider  $K$  categories  $C_1, \dots, C_k$

$P_F(X \in C_i) = \theta_i, i=1, \dots, k.$

Suppose  $P_{F_0}(X \in C_i) = \theta_i^0, i=1, \dots, k$

So our hypothesis testing problem can be transformed to

$H_0: \theta_i = \theta_i^0, i=1, \dots, k$

$H_1: \text{at least one inequality in the above statement}$

So roughly speaking we have a sample, so let us consider say  $X_1, X_2, \dots, X_n$  this is a random sample from say if  $F(x)$  is alright, and we want to test whether  $F(x) = F_0(x)$  or not. If you want to test  $F(x) = F_0(x)$  against  $F(x) = F_1(x)$  where  $F_1$  is different from  $F_0$ , then we have the most powerful test using the Neyman–Pearson lemma. However, it is not that, here we want to null hypothesis we are specifying completely, but alternatively we are not able to specify.

Therefore, we cannot apply the Neyman–Pearson lemma here, so what we do? We consider say we classify the data into  $K$  categories say  $C_1, C_2, \dots, C_k$ , and we calculate probability of  $X$  belonging to  $C_i$  where  $X$  is of course  $F(x)$ , what is the probability of  $X$  belonging to  $C_i$  let us denote by  $\theta_i$  okay, when the distribution is  $F$  for  $i=1$  to  $k$ . Now suppose this probability that  $X$  belonging to  $C_i$  that is say  $\theta_i^0$ , when  $F$  is taken to be  $F_0$ .

So actually we make use of this fact that means under the null hypothesis the probability of each category is specified, so we actually frame it as a multinomial testing problem, so our hypothesis testing problem can be transformed to  $H_0: \theta_i = \theta_{i0}, i = 1 \text{ to } k$ , against at least 1 inequality in the above statement. So you can see that in this Kolmogorov this Chi square test for goodness of fit, I am going to discuss it is one of the oldest nonparametric test it was developed by Karl Pearson.

**(Refer Slide Time: 04:32)**

Let  $f_i =$  no. of  $X_i$ 's belonging to category  $C_i$   
 $i = 1, \dots, k$

Then  $(f_1, \dots, f_k) \sim$  multinomial  $(n, \theta_1, \dots, \theta_k)$

We develop a likelihood ratio test

$$\Omega = \{(\theta_1, \dots, \theta_k) : \theta_i \geq 0, \sum_{i=1}^k \theta_i = 1\}$$

$$L(\theta_1, \dots, \theta_k) = \left( \prod_{i=1}^k \theta_i^{f_i} \right) \frac{n!}{f_1! \dots f_k!}$$

To maximize  $L$  over  $\Omega$ , we maximize

$$\ln L = \sum_{i=1}^k f_i \log \theta_i - \lambda (\sum \theta_i - 1)$$

$$\frac{\partial \ln L}{\partial \theta_i} = \frac{f_i}{\theta_i} - \lambda = 0 \Rightarrow \theta_i = \frac{f_i}{\lambda}, \quad i = 1, \dots, k$$

$$1 = \sum \theta_i = \frac{\sum f_i}{\lambda} = \frac{n}{\lambda} \Rightarrow \lambda = n$$

$$\hat{\theta}_i = \frac{f_i}{n} \rightarrow \text{maximizes } L(\underline{\theta}).$$

In this test actually cleverly the problem of full testing has been transformed to checking  $K$  categories, so we can actually consider let  $f_i$  be = number of  $X_i$ 's which are belonging to category  $C_i$  for  $i = 1$  to  $k$ , then we can say that this  $f_1, f_2, f_k$  this is having a multinomial distribution, total number of observations is  $n$  and the probabilities of categories  $C_1, C_2, C_k$  they are  $\theta_1, \theta_2, \theta_k$  respectively.

The simplest things that we can do is we consider a likelihood ratio test for this problem, so for likelihood test we know that we would like the likelihood function, we develop a likelihood ratio test, so here the full parametric space that is  $\theta_1, \theta_2, \theta_k$ , where  $\theta_i$ 's are  $\geq 0$  and  $\sum \theta_i = 1$ ,  $\theta_i$  to the power of  $f_i$   $n$  factorial /  $f_1$  factorial and so on  $f_k$  factorial  $i = 1$  to  $k$ . So basically this part is constant, so maximization problem is reduced to this part only.

So to maximize  $L$  over the parametric space  $\omega$ , we maximize  $\ln L$  so we take a log here,  $\ln L = \sum_{i=1}^k \theta_i \ln \left( \frac{f_i}{n} \right)$ , subject to the condition that  $\sum_{i=1}^k \theta_i = 1$ , so introduce Lagrange's multiplier, let us call this term as  $-\lambda$  Lagrange's multiplier let we call it  $\lambda$  here or we can give some other notations say  $M$  here. So that means we differentiate with respect to each  $\theta_i$ , then I will get  $f_i/\theta_i - \lambda = 0$  that gives me  $\theta_i = f_i/\lambda$  for  $i=1$  to  $k$ .

This gives me the value of  $\lambda$ , because I can apply the condition here  $\sum_{i=1}^k \theta_i = \sum_{i=1}^k f_i/\lambda = 1$ , this  $= n/\lambda$  this means  $\lambda = n$ , so  $\theta_i = f_i/n$  here. So this is the maximizing  $L$  theta.

**(Refer Slide Time: 08:14)**

So  $\sup_{\theta \in \Omega} L(\theta) = \hat{L}(\Omega_1) = \frac{n!}{f_1! \dots f_k!} \prod_{i=1}^k \left( \frac{f_i}{n} \right)^{f_i}$

$\hat{L}(\Omega_0) = \sup_{\theta \in \Omega_0} L(\theta) = \frac{n!}{f_1! \dots f_k!} \prod_{i=1}^k (\theta_i^0)^{f_i}$

The likelihood ratio is

$$\Lambda = \frac{\hat{L}(\Omega_0)}{\hat{L}(\Omega_1)} = \prod_{i=1}^k \left( \frac{n \theta_i^0}{f_i} \right)^{f_i}$$

$$-2 \log_e \Lambda = -2 \sum_{i=1}^k f_i (\log n \theta_i^0 - \log f_i)$$

$$= -2 \sum_{i=1}^k f_i \left[ \log f_i + (n \theta_i^0 - f_i) \cdot \frac{1}{f_i} + \frac{(n \theta_i^0 - f_i)^2}{2!} \left( -\frac{1}{f_i^2} \right) + \frac{(n \theta_i^0 - f_i)^3}{3!} (-1) (-2) \cdot \frac{1}{f_i^3} + \dots - \log f_i \right]$$

So what is the maximum value then, so the supremum value of  $L$  theta for theta belonging to  $\omega = \hat{L}(\Omega_1) = n! / (f_1! \dots f_k!)$  and so on  $f_k$  factorial product of  $f_i/n$  to the power  $f_i$   $i=1$  to  $k$ , because I have substituted the values of theta I is a  $f_i/n$  here, and  $\hat{L}(\Omega_0)$  that means supremum of  $L$  theta for theta belonging to  $\Omega_0$  that  $= n! / (f_1! \dots f_k!)$  product of  $\theta_i^0$  to the power  $f_i$   $i=1$  to  $k$ .

So if we consider the likelihood ratio that is  $L(\Omega_0)/\hat{L}(\Omega_1)$ , so this term will get cancelled out you will left with product  $n \theta_i^0 / f_i$  to the power  $f_i$ , let me call it say  $\lambda$  then  $-2 \log$  of  $\lambda$  that  $=$ , so this what we are doing to develop the distribution here, that is  $-2 \sum_{i=1}^k f_i (\log n \theta_i^0 - \log f_i)$ . Now this term  $\log$  of  $n \theta_i^0$  I expand around  $f_i$ , so that  $=$

twice  $\sum_{i=1}^k n \theta_i \log \left( \frac{n \theta_i - f_i}{n \theta_i} \right)$  then the derivative of this that is becoming 1 by that at  $f_i$  so it is simply becoming  $f_i/n$ .

Then second term will give me  $n \theta_i^2 / 2$  factorial second derivative will give me  $-1/f_i + n \theta_i^3 / f_i^2$ , this will be square here and so on - log of  $f_i$  that is this term here So we can now adjust the terms here, this term will get cancelled out, this term will give me  $n \theta_i^2 / f_i$  this  $f_i$  will get cancelled out, and in other terms here  $f_i$  will come here.

**(Refer Slide Time: 12:00)**

Handwritten mathematical derivation on a blue background:

$$-2 \log \lambda = -2 \left[ \sum_{i=1}^k (n \theta_i - f_i) - \frac{1}{2} \sum_{i=1}^k \frac{(n \theta_i - f_i)^2}{f_i} + \frac{1}{3} \sum_{i=1}^k \frac{(n \theta_i - f_i)^3}{f_i^2} - \dots \right]$$

$$= \sum_{i=1}^k \frac{(n \theta_i - f_i)^2}{f_i} - \frac{2}{3} \sum_{i=1}^k \frac{(n \theta_i - f_i)^3}{f_i^2} + \dots$$

Since  $\frac{f_i}{n} \xrightarrow{p} \theta_i$  under  $H_0$  i.e.  $f_i \xrightarrow{p} n \theta_i = e_i$  under  $H_0$

$$-2 \log \lambda \approx \sum_{i=1}^k \frac{(n \theta_i - f_i)^2}{f_i} = \sum_{i=1}^k \frac{(e_i - f_i)^2}{f_i} = Q$$

$$= \sum_{i=1}^k \frac{e_i^2 + f_i^2 - 2e_i f_i}{f_i}$$

$$= \sum_{i=1}^k \frac{e_i^2}{f_i} - n$$

Asymptotic dist<sup>n</sup> of  $Q$  is  $\chi_{k-1}^2$   
 So we reject  $H_0$  if  $Q \geq \chi_{k-1, \alpha}^2$  (at significance level  $\alpha$ )

So let us write this there is  $-2 \log \lambda = -2 \sum_{i=1}^k n \theta_i \log \left( \frac{n \theta_i - f_i}{n \theta_i} \right) - \frac{1}{2} \sum_{i=1}^k n \theta_i^2 / f_i + 1/3 \sum_{i=1}^k n \theta_i^3 / f_i^2$  and of course summation will be there and so on, this term is simply so  $\sum_{i=1}^k n \theta_i$  so  $\sum_{i=1}^k n \theta_i$  is also  $= n$ , this is  $n$  then  $\sum_{i=1}^k f_i$  is  $n$  so this first term will become 0, second term is giving me  $\sum_{i=1}^k n \theta_i^2 / f_i - 2/3 \sum_{i=1}^k n \theta_i^3 / f_i^2$  and so on.

Now  $f_i/n$  converges to  $\theta_i$  in probability under  $H_0$  that is you can say  $f_i$  converges to  $n \theta_i$  that we denote by  $e_i$  under  $H_0$  in probability. So therefore, I can say that  $-2 \log \lambda$  is asymptotically  $= \sum_{i=1}^k n \theta_i^2 / f_i$  for  $i=1$  to  $k$ , this is written as  $e_i - f_i / f_i$  square  $i=1$  to  $k$  that  $= Q$  here, so this is converging so therefore, the higher order terms we neglect here and we are writing only this,  $f_i$ 's are called observed frequencies,  $e_i$ 's are called expected frequency of the  $i$ th class.

And we have an alternative formula for this, this is also see if I expand it this will becoming  $e_i^2/f_i + f_i - 2e_i f_i/f_i$  that =  $\sum e_i^2/f_i + \sum f_i - 2 \sum e_i f_i$  both are  $n$  here so this will becoming  $-2 \sum_{i=1}^k (e_i^2/f_i - e_i f_i)$ . So this is the test statistic which is coming from the likelihood ratio, in the likelihood ratio we know that we accept the null hypothesis, if  $L(\omega) \geq L(\omega_0)$  we reject the null hypothesis if the denominator is  $\log$ .

So basically we have taken  $-2$ , so  $\lambda$  should be small for closer to  $H_0$ , and for rejection  $\lambda$  should be small, so this  $-2$  I have taken  $-2 \log$  of this so outcome will be reverse that means for large values of  $-2 \log \lambda$  will be rejecting. Another interpretation you can make out from here this is actually the difference between the observed and expected frequencies square, so if the 2 distribution is not  $f_0$ , then there will be large difference here that means these differences will propagate and this term will become large.

So basically this gives an indication the value of  $-2 \log \lambda$  whether the null hypothesis is true or not, now that gives a rough indication, but to get a real picture of this we will need the distribution of that, for that we see that this is multinomial. Therefore, asymptotic distributions of this simply will become the some of the chi-square, because we considered 2 then binomial convergence to normal, so here it is converging to  $k-1$  dimensional thing.

And therefore when we are taking some of the square it will convert to chi square on  $k-1$  degrees of freedom, so asymptotic distribution of this quantity let me call it  $W$  or  $Q$  we have called it that is chi square  $k-1$ , so we reject  $H_0$  if  $Q \geq \chi^2_{k-1, \alpha}$  at significance level  $\alpha$ . This test is widely used in all the applications for modelling of the statistical distributions and it is extremely useful.

However, since it is asymptotic certain assumptions are there, for example the expected frequency of each cell must be  $>5$  for a good approximation, if that is not so then this test is not very good.

**(Refer Slide Time: 18:10)**

Kolmogorov-Smirnov One-Sample Statistic

$H_0: F(x) = F_0(x) \forall x$   
 $H_1: F(x) \neq F_0(x) \text{ for some } x.$

$X_1, \dots, X_n$  is a random sample from  $F$ .

$F_n(x) \rightarrow$  EDF of  $(X_1, \dots, X_n)$



$D_n = \sup_x |F_n(x) - F_0(x)| \rightarrow$  Kolmogorov-Smirnov statistic  $D_n$

We further define

$D_n^+ = \sup_x (F_n(x) - F_0(x)),$  and

$D_n^- = \sup_x (F_0(x) - F_n(x))$

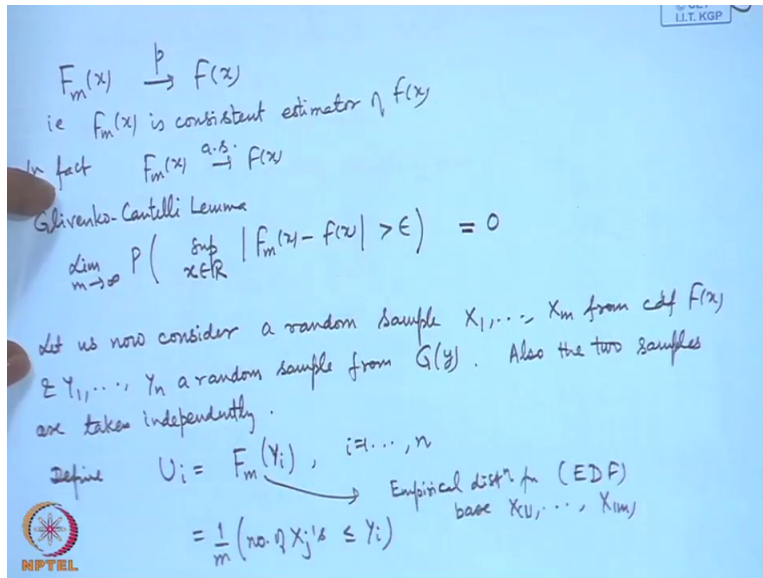
$D_n = \max(D_n^+, D_n^-)$

Now another test for the goodness of fit was developed by Kolmogorov and Smirnov, so this is called Kolmogorov-Smirnov one sample statistic, as before we are writing down our hypothesis testing problem as  $F(x) = F_0(x)$  for all  $x$  or  $F(x) \neq F_0(x)$ . We have at our disposal a random sample from this population, we define the sample distribution function that is  $F_n(x)$  that is empirical distribution function of  $X_1, X_2, \dots, X_n$  that is the ordered statistics from this.

We define the maximum absolute difference between the empirical distribution function and the assumed the distribution function, so here actually you take  $F_0$ . So what is the idea for this? The idea for this is the result about the empirical distribution function, which we gave earlier that was that it is strongly consistent not only strongly consistent.

**(Refer Slide Time: 19:58)**



We had actually proved that the limit of the probability supremum  $F_m(x) - F(x) > \epsilon$  actually goes to 0, so because of this, this is a very good indicator of the actual discrepancy between the assumed model and the sample that means based on the sample actually we calculating the empirical distribution function, so if there is too much discrepancy then this statement or this value will be large. So based on this idea this Kolmogorov-Smirnov they have defined the statistic called the  $D_n$ .

Now suddenly as in the previous chi-square goodness-of-fit, we need to discuss the distribution of  $D_n$ , if we are looking at  $F_n$  then certainly we knew the distribution, but since we are considering the maximum here so then this problem become slightly different. So we further define, so this one is actually the Kolmogorov-Smirnov statistic we call it  $D_n$  okay, so we further define 2 quantities called  $D_n^+$  that is= supremum of  $F_n(x) - F(x)$ .

And  $D_n^-$  that is= supremum of  $F(x) - F_n(x)$ , so here I will put reverse that is  $D_n$  is actually = the maximum of  $D_n^+$  and  $D_n^-$  that means actually I am taking the maximum positive difference and maximum negative difference. Now we will try to analyze this distribution of  $D_n^+$  and  $D_n^-$  separately.

**(Refer Slide Time: 22:21)**

© CET  
I.I.T. KGP

$$\begin{aligned}
 D_n^+ &= \sup_x (F_n(x) - F(x)) \\
 &= \max_{i=0, \dots, n} \sup_{X_{(i)} \leq x < X_{(i+1)}} (F_n(x) - F(x)) \\
 &= \max_{i=0, \dots, n} \left( \sup_{X_{(i)} \leq x < X_{(i+1)}} \left( \frac{i}{n} - F(x) \right) \right) \\
 &= \max_{0 \leq i \leq n} \left( \frac{i}{n} - F(X_{(i)}) \right) \\
 &= \max \left\{ 0, \max_{i=1, \dots, n} \left( \frac{i}{n} - U_{(i)} \right) \right\} \text{ from } U[0,1]
 \end{aligned}$$

Thus we have shown that  $D_n^+$  is distribution free

$X_{(0)} < X_{(1)} < \dots < X_{(n)} < X_{(n+1)}$   
 $\downarrow \qquad \qquad \qquad \qquad \qquad \qquad \downarrow$   
 $-\infty \qquad \qquad \qquad \qquad \qquad \qquad +\infty$

$X_{(1)} \quad X_{(2)} \quad X_{(3)} \quad \dots \quad X_{(n-1)} \quad X_{(n)}$

So  $D_n^+$  that is the supremum of  $F_n(x) - F(x)$  overall  $x$ , we have ordered  $X_1 \leq X_2 \leq \dots \leq X_n$ , and I take  $X_0$  and  $X_{n+1}$  also or  $X_{-\infty}$  and  $X_{+\infty}$ , so this is I am taking to  $-\infty$ , this I am taking to  $+\infty$ . So we can then express it as maximum of supremum value between  $F_n(x) - F(x)$ , let me put this into one side, this is for  $i = 0, 1, 2, \dots, n$  that means I am considering that supremum value in the intervals, so this is  $X_1$  that is  $-\infty$  to  $X_1$ , then between  $X_1$  to  $X_2$ , then between  $X_2$  and  $X_3$ , then between  $X_{n-1}$  and  $X_n$  and then  $X_n$  to  $+\infty$ .

So I have divided this problem into looking at the difference into each of the intervals, but the advantage of this approach is that actually the value of the empirical distribution function in this interval is  $i/n$ , so basically I am looking at that what is the difference of  $F(x)$  from  $i/n$ , when  $X$  is in the interval  $X_i$  to  $X_{i+1}$  and this we are doing over all  $i$ 's. Now capital  $F$  is an increasing function because it is a cdf, so when I look at the supremum here.

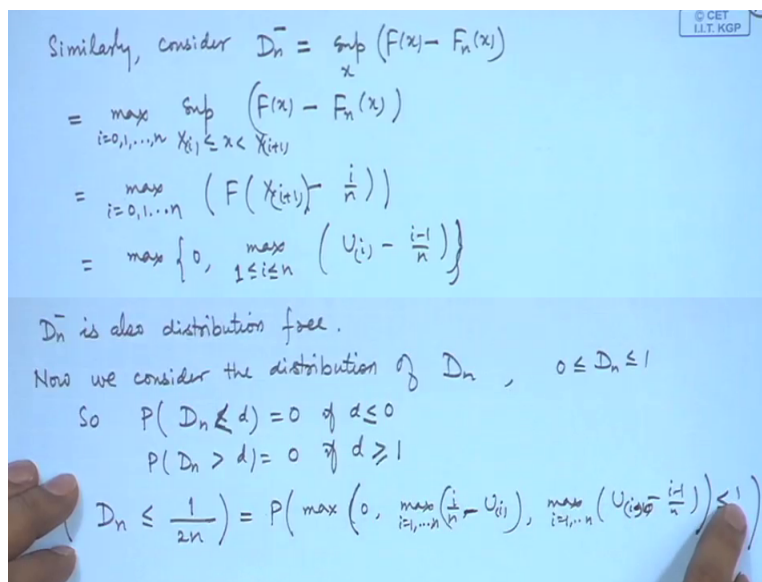
Now it is supremum over  $X$ , so this will become a fixed quantity, so this becomes actually the minimum value of  $F$  the minimum value of  $F$  in this interval is attained at  $X_i$ , so this is becoming  $= \max_{i=0, 1, 2, \dots, n} (i/n - F(X_i))$  and now this  $F(X_i)$ 's are actually  $U_i$ 's that we have already seen, so this is  $i/n - U_i$ 's, where  $U_i$ 's are the order statistics from the uniform  $0,1$  okay this is from uniform  $0, 1$  here, and we are looking at the maximum for  $i=1$  to  $n$ , and corresponding to  $0$  then this is actually  $0$ .



So this is maximum of this and this, now this is very interesting here, I started with some sample here okay, now based on that sample I have considered the difference between  $F_n(x) - F(x)$  that is an empirical distribution function  $F_n(x)$ , but this quantity if you look at this quantity has become free from the original distribution, because this is nothing but from the uniform  $0, 1$ . Thus we have shown that  $D_n^+$  is distribution free.

As I mentioned earlier in the beginning of this particular section on nonparametric methods that here we develop those methods which are free from the distribution original distribution assumption, so that means that whatever the distribution originally it does not affect our distribution that means distribution assumption is not required except of course we consider continuity etc. here.

**(Refer Slide Time: 27:01)**



Similarly, consider  $D_n^- = \sup_x (F(x) - F_n(x))$

$$= \max_{i=0,1,\dots,n} \sup_{X_{(i)} \leq x < X_{(i+1)}} (F(x) - F_n(x))$$

$$= \max_{i=0,1,\dots,n} (F(X_{(i+1)}) - \frac{i}{n})$$

$$= \max \left\{ 0, \max_{1 \leq i \leq n} \left( U_{(i)} - \frac{i-1}{n} \right) \right\}$$

$D_n^-$  is also distribution free.

Now we consider the distribution of  $D_n^-$ ,  $0 \leq D_n^- \leq 1$

So  $P(D_n^- \leq d) = 0$  if  $d \leq 0$

$P(D_n^- > d) = 0$  if  $d \geq 1$

$$P(D_n^- \leq \frac{1}{2n}) = P\left(\max\left(0, \max_{i=1,\dots,n} \left(\frac{i}{n} - U_{(i)}\right), \max_{i=1,\dots,n} \left(U_{(i)} - \frac{i-1}{n}\right)\right) \leq \frac{1}{2n}\right)$$

Now in a similar way I can consider  $D_n^-$  here, so let us consider  $D_n^-$  now what is  $D_n^-$ ,  $D_n^- =$  supremum of  $F(x) - F_n(x)$  that means I am taking the negative value here, so this = maximum of supremum  $F(x) - F_n(x)$  for  $X_{(i)} \leq x < X_{(i+1)}$   $i=0, 1, 2, n$  that is = maximum of  $i=0, 1, 2, n$  that is =  $F(X_{(i+1)} - i/n)$ , that is = maximum of  $0, \max_{1 \leq i \leq n} (U_{(i)} - i/n)$ , so I have shifted by 1 here because I am taking 1 to  $n$ , so I can in place of  $i+1$  I can write  $i$  then this becomes  $i-1$ .

So once again as in  $D_n^+$  here also you see so  $D_n^-$  is also distribution free, now first thing is that we are able to derive the form of  $D_n^+$  and  $D_n^-$  in terms of the order statistics from uniform  $0,1$ .

The distribution of  $U$  bracketed  $i$  is known that we have derived as the beta distribution, now here the form that is coming out is the maximum that means when we are considering several dependent distributions or dependently distributed random variables.

Then what is the distribution of the maximum of that, and then again maximum of the to2, let me express it here. So now we consider the distribution of  $D_n$ , now one thing that you note these values are between 0 to 1, and what are these values here  $i-1/n$  these values are also between 0 to 1, so these all values actually always lie between 0 to 1. If you look at this also this is  $i/n$  these values lie between 0 to 1, so this values will also lie between 0 to 1 only okay.

So the entire thing is that  $D_n$  lies between 0 to 1, so that means when we are considering distribution of this  $\leq$  say some  $d$  then it is  $=0$  if  $d$  is  $<0$ , and probability of  $D_n >$  say  $d$  that is  $=0$  if  $d$  is  $>=1$  okay. Now let us consider  $D_n < d$  between 0 to 1, so we put a particular form here, why that particular form? It will clear when we derive the expression here. Let us consider say probability of  $D_n \leq 1/2n$ .

So there is a reason that why I am considering  $1/2n$ , the reason is that if you look at these values here they are of the form  $1/n$  etc. in each interval if I look at this, so the differences will be of this nature and let me firstly derive this here. So  $D_n$  is nothing but probability of maximum of 0 maximum of  $i/n - U_i$  where  $i$  is from 1 to  $n$ , and maximum of  $U_i - i-1/n$  and again here  $i=1$  to  $n$ , so what we are saying is that this is  $\leq 1/2n$ , so maximum of this and this, and this  $\leq 1/2n$ .

**(Refer Slide Time: 32:43)**

$$\begin{aligned}
 &= P\left(\max_{i=1, \dots, n} \left(\frac{i}{n} - U_{(i)}\right) \leq \frac{1}{2n}, \max_{1 \leq i \leq n} \left(U_{(i)} - \frac{i-1}{n}\right) \leq \frac{1}{2n}\right) \\
 &= P\left(\frac{i}{n} - U_{(i)} \leq \frac{1}{2n}, U_{(i)} - \frac{i-1}{n} \leq \frac{1}{2n}, i=1, \dots, n\right) \\
 &= P\left(\frac{i}{n} - \frac{1}{2n} \leq U_{(i)} \leq \frac{i}{n} - \frac{1}{2n}, i=1, \dots, n\right) \\
 &= P\left(U_{(i)} = \frac{i}{n} - \frac{1}{2n}, i=1, \dots, n\right) \\
 &= 0 \text{ as } U_{(i)} \text{'s are continuous r.v.'s.}
 \end{aligned}$$

So we need to consider  $P\left(D_n < \frac{1}{2n} + \epsilon\right)$  Proceeding as above we get

$$= P\left(\frac{2i-1}{n} - \epsilon < U_{(i)} < \frac{2i-1}{n} + \epsilon, i=1, \dots, n\right)$$

We splitted, so  $0 \leq \frac{1}{2n}$  is always true, so this probability then can be expressed as probability of maximum of  $\frac{i}{n} - U_i \leq \frac{1}{2n}$  for  $i=1$  to  $n$ , and maximum of  $U_i - \frac{i-1}{n} \leq \frac{1}{2n}$  this is  $1 \leq i \leq n$ , this is = probability of now these are already ordered here, so this we can considered as  $\frac{i}{n} - U_i \leq \frac{1}{2n}$  and  $U_i - \frac{i-1}{n} \leq \frac{1}{2n}$  this is true for  $i=1$  to  $n$ , these 2 statements I can combine. Then this is nothing but  $U_i \in [\frac{i}{n} - \frac{1}{2n}, \frac{i}{n} - \frac{1}{2n}]$  and on this side I have I now if you take it to the other side then it is becoming  $\frac{i}{n} - \frac{1}{2n}$   $i=1$  to  $n$ .

So this is interesting both the sides are the same so this is actually becoming= probability of  $U_i = \frac{i}{n} - \frac{1}{2n}$  so that was the reason that I mentioned that why I am considering  $D_n \leq \frac{1}{2n}$ , because for this particular part this is giving extremely simple expression that is the probability of  $U_i$ , so certainly  $U_i$ 's are continuously random variable therefore, this probability will be=0. So what we are finding here that the probability of  $D_n \leq \frac{1}{2n}$ .

That means  $D_n$  has to start from  $\frac{1}{2n}$  from the original definition it is not clear what is the starting point, so we said  $D_n$  lies between 0 to 1, but now we see that even  $D_n \leq \frac{1}{2n}$  it is giving you probability 0. So we consider then probability of  $D_n < \frac{1}{2n} + \text{something}$  okay as going as before what will happen here? Here I will get  $\frac{1}{2n} + v$ , here I will get  $\frac{1}{2n} + v$ , so here  $\frac{1}{2n} + v$ , here I will get  $\frac{1}{2n} + v$ , then if I am having this term here I will get  $-v$  here and here I will get  $+v$ .

So proceeding as above we get this= probability of  $2i-1/n-v < U_i < 2i-1/n+v$   $i=1$  to  $n$ , now this is the joint probability for the random variables  $U_1, U_2, U_n$  which are the order statistics from uniform  $0,1$ , the joint pdf of this is known, so it is nothing but the  $n$ -fold integral over this region like for  $U_1$  you will have from  $1/n-v$  to  $1/n+v$ , for  $U_2$  it will be maximum of so actually then this region because you also have even  $U_1 < U_2 < U_n$  and they are lie between  $0$  to  $n$ .

**(Refer Slide Time: 37:16)**

This is an  $n$ -fold integral over the joint pdf of  $U_{(1)}, \dots, U_{(n)}$ ,  
 i.e.  $n!$ ,  $0 < u_1 < u_2 < \dots < u_n < 1$   
 over the given region.  
 Owen, Birnbaum & Smirnov have tabulated upper  $100\%$  points  
 $D_{n,\alpha}$  of dist<sup>n</sup> of  $D_n$  i.e.  $P(D_n > D_{n,\alpha}) = \alpha$   
 $P(D_n > D_{n,\alpha}) = \alpha$   
 $\lim_{n \rightarrow \infty} P(D_n < \frac{c}{\sqrt{n}}) = 1 - 2 \sum_{i=1}^{\infty} (-1)^{i-1} e^{-2i^2 c^2}$   $0 < c < 1$   
 $P(D_n^+ < c) = P(\max_{i=1, \dots, n} (\frac{i}{n} - U_{(i)}) < c)$   
 $= P(\frac{i}{n} - U_{(i)} < c, i=1, \dots, n) = P(U_{(i)} > \frac{i}{n} - c, i=1, \dots, n)$

So this is nothing but, this is an  $n$ -fold integral over the joint pdf of  $U_1, U_2, U_n$  that is  $n$  factorial for  $0 < u_1 < u_2 < u_n < 1$  over the given region okay, so if I have to take say  $n=2$  or  $n=3$ , then these things can be evaluated. Then Owen, Birnbaum and Smirnov, they have tabulated the values of the upper 100% alpha points of the distribution of the  $D_n$ , let us call it  $D_n, \alpha$  of the distribution of  $D_n$  that is probability of  $D_n > D_n, \alpha = \alpha$  for various values of  $n$  and  $\alpha$ . We can actually consider probability of  $D_n > D_n, \alpha$  that= under  $H_0 = \alpha$ .

So he also considered some asymptotic distribution also that is probability of  $D_n < v/\sqrt{n}$  if we consider as limit  $n$  tends to infinity, it was shown that it is  $1 - 2 \sum_{i=1}^{\infty} (-1)^{i-1} e^{-2i^2 v}$ . And if I consider say a number  $c$  between  $0$  and  $1$ , then probability of  $D_n < c$  that is= probability of maximum  $i/n - U_i < c$  for  $i=1$  to  $n$  that is= probability of  $i/n - U_i < c$   $i=1$  to  $n$  that is= probability of  $U_i > i/n - c$  for  $i=1$  to  $n$ . Once again you see that this can be evaluated in terms of the joint distribution of the  $U_1, U_2, U_n$ .

**(Refer Slide Time: 40:13)**

© CET  
I.I.T. KGP

$$P(D_n^- < c) = P(U_{(i)} < c + \frac{i-1}{n}, i=1, \dots, n)$$

$$= P(1 - U_{(i)} > \frac{n-i+1}{n} - c, i=1, \dots, n)$$

If  $U_1, \dots, U_n$  are i.i.d.  $U[0,1]$   
 $1-U_1, \dots, 1-U_n$  are also i.i.d.  $U[0,1]$   
 $V_1, \dots, V_n$  i.i.d.  $U[0,1]$

$$V_{(i)} = U_{(n-i+1)}$$

$$P(V_{(i)} > \frac{n-i+1}{n} - c, i=1, \dots, n)$$

$$= P(U_{(n-i+1)} > \frac{n-i+1}{n} - c, i=1, \dots, n)$$

$$= P(U_{(i)} > \frac{i}{n} - c, i=1, \dots, n)$$

So  $D_n^+$  &  $D_n^-$  have the same dist.  
 We can directly use  $D_n^+$  &  $D_n^-$  for testing  
 $H_0: F(x) = F_0(x) \forall x$   
 $H_1: F(x) \neq F_0(x)$  for some  $x$ .

Similarly, if we consider  $D_{n-1} < c$  then this will be  $U_i < c + i-1/n$   $i=1$  to  $n$  that is = probability of  $1-U_i > n-i+1/n-c$ ,  $i=1$  to  $n$ , if we consider say  $U_1, U_2, U_N$  they are i i d uniform  $0,1$  then  $1-U_1, 1-U_2, 1-U_n$  are also i i d uniform  $0,1$  that means if I consider  $V_1, V_2, V_n$  which is  $1-U_i$  then they are i i d uniform  $0,1$ , so this is actually the same that means  $V_i = U$  of  $n-i+1$ , so this one then we can write as probability of  $V_{n-1} - V_i > n-i+1/n-c$ ,  $i=1$  to  $n$ .

Then this can be written as  $U_{n-i+1} > n-i+1/n-c$ ,  $i=1$  to  $n$ , then this is nothing but  $U_i > i/n-c$ ,  $i=1$  to  $n$ , now you compare this with this expression here, probability of  $D_{n+1} < c$  is probability of  $U_i > i/n-c$  and it is the same thing here also. So what we are getting that  $D_{n+}$  and  $D_{n-}$  have the same distribution, so one can directly use of the  $D_{n+}$  and  $D_{n-}$  for the testing problem. We can directly use  $D_{n+}$  and  $D_{n-}$  for testing  $H_0 F(x) = F_0(x)$ ,  $H_1 F(x) \neq F_0(x)$  for some  $x$ .

**(Refer Slide Time: 43:24)**

© CET  
I.I.T. KGP

$$\lim_{n \rightarrow \infty} P(D_n^+ < \frac{z}{\sqrt{n}}) = 1 - e^{-2z^2} = F(z) = 1 - \alpha \quad (*)$$

$$U = 4n D_n^{+2}$$

$$P(U \leq u) = P(4n D_n^{+2} \leq u)$$

$$= P(D_n^+ \leq \frac{\sqrt{u}}{2\sqrt{n}})$$

$$\rightarrow 1 - e^{-\frac{2u}{n}} = 1 - e^{-u/2}$$

So  $\lim_{n \rightarrow \infty} f_U(u) = \frac{1}{2} e^{-u/2}, u > 0$  is Negative exponential dist<sup>n</sup>.

From (\*) we can choose  $D_{n, \alpha}^+$

$$1 - \alpha = P(\sqrt{n} D_n^+ < z) = P(4n D_n^{+2} < 4z^2)$$

$4z^2 = \chi_{2, 1-\alpha}^2$

NPTEL

The asymptotic distribution of  $D_n^+$  etc. is also being worked out, if I look at  $D_n^+ < z/\sqrt{n}$  as  $n$  tends to infinity that is  $= 1 - e^{-2z^2}$  that is  $F(z) = 1 - \alpha$ , if I define  $U = 4n D_n^{+2}$  then probability of  $U \leq u$  then that is  $=$  probability of  $4n D_n^{+2} \leq u$  that is probability of  $D_n^+ \leq \sqrt{u}/2\sqrt{n}$ . So if I apply this formula the limit will become  $= 1 - e^{-u/2}$  that is  $1 - e^{-u/2}$ .

So the limiting pdf of  $u$  that  $= 1/2 e^{-u/2}$  that is negative exponential distribution which also can be said as the Chi square distribution on 2 degree of freedom, of course since it is a negative exponential distribution the percentage points of this can be easily calculated, and we can express the test in terms of this also, so asymptotic test, asymptotic confidence interval can be obtained in terms of this, if we call this as say star.

Then from star we can choose say  $D_{n, \alpha}^+$  such that  $1 - \alpha =$  probability  $\sqrt{n} D_n^+ < z =$  probability  $4n D_n^{+2} < 4z^2$ , that is  $4z^2 = \chi_{2, 1-\alpha}^2$  okay, so this can be easily calculated, one can easily find out the confidence interval for  $F_x$ .

**(Refer Slide Time: 45:58)**

We can also use  $D_n$  to find confidence interval for  $F(x)$

© CET  
I.I.T. KGP


$$P(D_n \leq D_{n,\alpha}) = 1 - \alpha$$

$$P\left(\sup_x |F_n(x) - F(x)| \leq D_{n,\alpha}\right) = 1 - \alpha$$

$$\Rightarrow P(|F_n(x) - F(x)| \leq D_{n,\alpha} \quad \forall x) = 1 - \alpha$$

$$\Rightarrow P\left(F_n(x) - D_{n,\alpha} \leq F(x) \leq F_n(x) + D_{n,\alpha} \quad \forall x\right) = 1 - \alpha.$$

So  $100(1-\alpha)\%$  confidence interval for  $F(x)$  is of the form  
 $\left(\max(0, F_n(x) - D_{n,\alpha}), \min(1, F_n(x) + D_{n,\alpha})\right)$



We can also use  $D_n$  to find confidence interval for  $f(x)$  that is probability  $D_n \leq D_{n,\alpha}$  that is  $1 - \alpha$ , so we write it as supremum of  $F_n(x) - F(x)$  this is equivalent to saying  $F_n(x) - F(x) \leq D_{n,\alpha}$  for all  $x$  that is  $1 - \alpha$ , this is equivalent to saying probability  $F_n(x) - D_{n,\alpha} \leq F(x) \leq F_n(x) + D_{n,\alpha}$  that is  $1 - \alpha$ . So  $100(1 - \alpha)\%$  confidence interval for  $F(x)$  is of the form maximum of  $0, F_n(x) - D_{n,\alpha}$  to minimum of  $1, F_n(x) + D_{n,\alpha}$ .

So we have seen that this test Kolmogorov-Smirnov test it is actually this is using more values compared to the Chi square test for goodness of fit was developed by Karl Pearson, in the Karl Pearson test essentially we reduced it to category problem that means we consider  $K$  classes out of the full distribution, and therefore, the test is more sensitive because what categories you are choosing, how many categories you are taking it will be dependent upon that.

Whereas this test is more robust, of course it is having sensitivity in the heavy tailed distribution but that is beside the point, there have been some modifications they have been proposed but essentially what we have seen is the distribution of the  $D_n$  is actually derivable here, so this is a much you can say improved thing compared to the chi-square test for goodness of fit, only thing is that the use of Kolmogorov-Smirnov is not that straight forward for the persons who have no idea about use of the statistics.

Because they need to understand the tabular version of the distribution of  $D_n$  that means how the percentage points are calculated, whereas for the Chi square test the percentage points are simply the percentage points of a chi-square distribution, so with a little knowledge of distribution theory one can actually apply the test. So it is a say compromise ease of applications is there in chi-square test, but the robustness is more in the Kolmogorov-Smirnov test.

**(Refer Slide Time: 49:40)**

Single Sample Location Problems

Let  $X_1, \dots, X_n$  be a random sample from a dist<sup>n</sup> with cdf  $F(x)$ .  
 Let  $\theta$  denote median of  $F(x)$ . Let  $F$  be strictly increasing and continuous at  $x = \theta$ .  
 We want to test  $H_0: \theta = \theta_0$   
 vs  $H_1: \theta > \theta_0$ ,  $H_2: \theta < \theta_0$ ,  $H_3: \theta \neq \theta_0$

If we shift our observations  $X_i$  to  $X_i - \theta_0$ , then median of new dist<sup>n</sup> will become zero. So without loss of generality, we take  $\theta_0 = 0$ .

We define 
$$Y_i = \begin{cases} 1 & \text{if } X_i > 0 \\ 0 & \text{if } X_i \leq 0 \end{cases}$$

$S = \sum_{i=1}^n Y_i \rightarrow$  Sign Test Statistic  
 $\downarrow$  no. of positive  $X_i$ 's.

Under  $H_0$ ,  $S \sim \text{Bin}(n, \frac{1}{2})$

In general  $P(X_i > 0) = p$   
 under  $H_0: p = \frac{1}{2}$

Next we consider single sample location problems, actually in the beginning I have given you some applications of that 2 sample functions that  $F_m, Y_i$  that means which are based on  $F$  and  $G$  you have 2 samples and based on that some restrict is for the location etc. are given. Now I am getting into use of all this ranks here and to derive the few test for the location problem, so first we let us consider 1 sample problem.

So let us consider  $X_1, X_2, X_n$  be a random sample, suppose you have the cdf  $f_x$ , let  $\theta$  denote median of  $F_x$ , and we assume let  $F$  be strictly increasing and continuous at  $x=0$  that means we are assuming median to be the unique. So we want to test  $\theta = \theta_0$  against say  $\theta > \theta_0$ , or  $\theta < \theta_0$ , or  $\theta \neq \theta_0$ , so these 3 types of alternative hypothesis we will be considering, and you can compare it with the parametric testing problem.

In the parametric testing problem for one sample we were testing whether the mean value is = something < something > something != something,  $\mu = \mu_0$  etc. we have done that testing

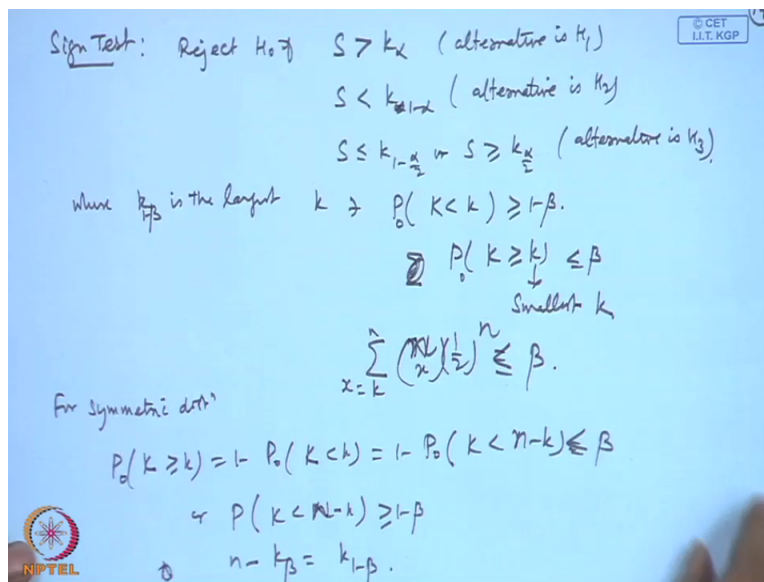


problems under the assumptions of the normality. So here there is no distribution assumptions made except that it is a continuous distribution and strictly increasing and continuous at theta that means the median is uniquely defined.

So now we want to test about some value theta 0, so whatever be the value theta 0 since this is known, if we shift our observations Xi to Xi-theta 0 then median of new distribution will become 0, so actually we do that thing. So without loss of generality, we take theta 0 to be 0, so that the problem becomes slightly simpler. We define say psi i= 1 if Xi is >0, it is 0 if Xi is <=0. And we define S= sigma psi i, i=1, this is called sign test statistic.

Because this is giving you the number of positive Xi's, how many Xi's are positive did exactly telling that thing. Actually under H0 if the null hypothesis is true then the 0 will be the median then under H0 S will follow binomial n, 1/2. See in general you will have probability of Xi>0 some p, but under H0 we will p=1/2, so under H0 the distribution of S is binomial n, 1/2 so we can actually devise a simple heuristic test based on the sign test.

**(Refer Slide Time: 54:41)**



So sign test, reject H0 if S is > let us say some k alpha, this is if alternative is H1, and if S is < k alpha that is k 1-alpha if alternative is H2, and S <= k 1-alpha/2 or S >= k alpha/2 if alternative is H3. Where k beta is the largest k such that probability of k 1-beta is the largest k such that probability K < k is >= 1-beta, we have to take this largest etc. because the distribution is assumed

to be discrete, so we may not actually achieved equality here that is probability of  $K > k$  alpha need not be alpha.

So that is why we choose the largest such cut off point, so basically we are saying that this condition is equivalent to that sigma, basically we are saying probability  $K \geq k$  is  $\leq \alpha$ , so then this small k is the smallest k that is sigma  $N C x 1/2$  to the power n,  $x=k$  to n that is  $\leq \alpha$ , so this can be easily calculated from the tables of the binomial distribution. If the distribution is symmetric, then we have this  $K \geq k$  that is  $= 1 - \text{probability } K < k$  that is  $1 - P_0 K < n - k$  that is  $\leq \alpha$ . Then we can say that probability of  $K < n - k$  is  $\geq 1 - \alpha$ , so you will have  $n - k = k - 1 - \alpha$ .

(Refer Slide Time: 57:37)

Power Functions

$$\beta(\theta) = P_{\theta}(K \geq k_{\alpha}) = \sum_{x=k_{\alpha}}^n \binom{n}{x} (1 - f_{\theta}(0))^x (f_{\theta}(0))^{n-x}$$

where  $k_{\alpha}$  is the smallest  $k$  such that  $P_{\theta=0}(K \geq k) \leq \alpha$   
or  $\sum_{x=k}^n \frac{1}{2^n} \binom{n}{x} \leq \alpha$

Example  $F(x) \in N(\theta, \sigma^2)$   
 $F_{\theta}(0) = P_{\theta}(X \leq 0) = \Phi\left(-\frac{\theta}{\sigma}\right) = 1 - \Phi\left(\frac{\theta}{\sigma}\right)$   
 $f_{\theta}(0) = \frac{1}{2}$ ,  $\alpha = 0.0384$ ,  $n = 16$ ,  $\sigma^2 = 1$   
We can see from binomial tables that  $k_{\alpha} = 12$ .  
at  $\theta = 1.04 = \sum_{x=12}^{16} \binom{16}{x} (0.8508)^x (0.1492)^{16-x} \approx 0.9211$   
Power =  $P(T > 1.77) = 0.9510$   
 $t_{0.0284} = 1.77$

We can also calculate the power function here, power function of the sign test, probability of  $K \geq k$  alpha when theta is the true value true medium then it is  $= x = k$  alpha to n,  $n C x 1 - F_{\theta}(0)$  to the power  $x * F_{\theta}(0)$  to the power  $n - x$ , where k alpha is the smallest k such that probability  $K \geq k$  under  $\theta = 0 \leq \alpha$  or we can say  $1/2$  to the power n  $C x$  is  $\leq \alpha$ , for  $x = k$  to n. Let us take us an example here.

Suppose, I consider  $F(x)$  to be normal  $\theta$  sigma square, so  $F_{\theta}(0)$  that is probability  $X \leq 0$  that is  $\Phi(-\theta/\sigma)$ , that is  $1 - \Phi(\theta/\sigma)$ , so  $F_{\theta}(0) = 1/2$ , let us take say  $\alpha = 0.0384$   $n = 16$ ,  $\sigma^2 = 1$ . Then we can see from the binomial tables that  $k_{\alpha} = 12$ , so if we

consider the power function at say  $\theta = 1.04$ , then it is  $\sigma = 16 \times 0.8508$  to the power  $x$ ,  $0.1492$  to the power  $16-x$ ,  $x=12$  to  $16$ , then that is approximately  $0.9211$ .

If we consider the corresponding  $t$  test that is  $4 \times \bar{S} > t_{0.0384}$  that is  $1.77$ , then power of this test is = probability of  $T > 1.77$  that is  $0.9918$ , so certainly we can say that the power of the sign test is  $<$  the power of the usuality test that we already know, but this is under the assumption of the normality, if we actually have known all this about the normality then this test will be measurable, it will fail because this will simply give a wrong thing.

Another thing is that asymptotically also we can use the sign test, because we are saying  $k$  follows the  $S$  follows the binomial distribution, so this binomial distribution asymptotically becomes a normal distribution. So one can actually use this also, so in both the cases the results can be obtained and the cutoff point, the critical point and the power of the test can be easily calculated.

In the following lectures, I will be describing some other test statistics which are based on the order statistics, in the sign test actually the order statistics are not used only the sign of the term is important. So therefore, in that sense you can say extremely simplistic test for the median of the distribution. Next, we will define certain test which will be based on the actual values or the actual measurements, so that I will be starting like Wilcoxon signed-rank test statistic, Mann Whitney and so on.