Nonparametric Statistical Inference Professor. Niladri Chatterjee Department of Mathematics Indian Institute of Technology, Delhi Lecture No. 07

Welcome students to MOOCs series of lectures on Nonparametric Statistical Inference, this is lecture number 7. As I said in the last class that in this lecture we shall study tests for randomness. Now random process or random number or random graph, random indexing are certain terms which are nowadays very important because of machine learning.

(Refer Slide Time: 00:56)



But from a statistician's perspective most important thing is random sample. Question is what is random? Cambridge dictionary gives this is the definition of randomness that is happening, done or chosen by chance rather than according to a plan.

(Refer Slide Time: 01:12)



So, when we talk about drawing inference about a population or about a distribution we often take a random sample. What does it mean? It means that sample is taken without any bias or predesigned rule of selection or to support some purpose etcetera. Therefore, when we have a data we need to test the randomness of the data. How to do that? In order to check whether something is random it is very important to maintain the order in which the samples are chosen.

If we consider the data points only as a set then we lose the flavor of randomness because the order in which the data points have been chosen that is very important as I am going to illustrate now.

(Refer Slide Time: 02:12)



Suppose, you are taking 12 balls from a basket with a large number of red and blue balls. Now we have taken 6 red balls and 6 blue balls. Suppose, the order of selection is this first 6 are blue balls and the second set of 6 are red balls. Will you consider it to be a random selection? No. What about this if we choose it alternatively blue, red, blue, red like that then do you consider this to be a random selection?

Again we understand that it is not random. Let us consider a third example. If the pattern of selection is like this: first 3 Rs then 6 Bs then again 3 Rs we understand that there is a specific way of selecting the balls from the basket. Therefore, this also is not random. However, if we consider this selection red, blue, blue, red, red then blue then again red, red then blue, blue, blue, blue and red, then it appears that the selection of red and blue ball from the basket is rather random. Question is how do we test because all the four selection orders are having 6 blue and 6 red balls.

(Refer Slide Time: 03:55)



So, that is the basic question how to test whether a particular sequence is random or not.

(Refer Slide Time: 04:04)



The basic theory developed in test of randomness is based on the concept of runs.

(Refer Slide Time: 04:15)



So, what is a run? Given an ordered sequence or two or more types of symbols, a run is defined to be succession of one type of symbol that is succeeded and preceded by another type of symbol or no symbol at all. So, let us consider the 4 selection orders that you have just seen. With respect to the first one, this is one run of blue ball because there is no red ball in between. It is a continuous sequence of observations of 6 blue balls similarly a continuous sequence of observation of 6 red balls.

Therefore, how many runs are there? There are 2 runs. This is succeeded by this run and this sequence of R balls is preceded by the sequence of blue balls. Consider the second one how many runs are there? There are 12 runs as I have explained here this is 1 run blue it is intercepted by a red again that is intercepted by a blue then red so like that there are 12 number of runs in this sequence.

With respect to the third one we can say that this is one run of R, this is one run of blue and this is one run of R again that gives us 3 runs and this one is very similarly first run, second run, third, fourth, fifth, sixth and seventh there are seven different runs in this sequence. What does it tell us?

(Refer Slide Time: 06:08)



So, it is clear from the above example that if the number of runs is too large or too small, then it is unlikely that the points are chosen randomly. However, what is too large and what is too small that depends upon the composition of the data. For example, we have in the previous illustrations 12 balls of which 6 B's and 6 R's that means 6 are blues and 6 are reds. Therefore maximum possible run can be 12 of course, the minimum will be 2, one sequence of red and one sequence of blue.

However, suppose there are 10 B's and 2 reds. then the maximum number of runs can be 5. One sequence of blue, one sequence of red may be one red, one sequence of blue, one red and another sequence of blue. This is the maximum number of runs that can be achieved from 10 blues and 2 Rs. Thus, the maximum number of runs is 5 although the minimum will still remain 2. Therefore, when we are talking about if the number of runs is too large we have to be very careful about the composition of the data as well.

(Refer Slide Time: 07:52)



Now the tests based on runs are applicable to either quantitative or qualitative data. It is not necessary that the data has to be numeric. It can be symbolic as we have already seen sequence of R's and B's in our previous illustrations. However, in the case of numeric data typically the dichotomy is effected by comparing each number with a focal point which can be median or mean. What do we mean by that?

Dichotomy means that when there is a sequence of numbers how do we split it into two different sets like R and B. In sign test we have seen plus and minus so how do we do that? So, that is done by comparing with median or mean. For example, suppose we have taken 12 real numbers between 50 and 100. So this is one sequence 95.6, 88.3, 75.8 like that 59.2. Another sequence is 59.2, 88.3 and up to 60.6, which one of them is random and why?

To test that what we do we first compute the mean or median. In this case median is coming out to be 72.8 because this is sort of the average of these two then what we do as we did in the case of sign test we consider how many of these observations are bigger than the median and how many of them are smaller than the median. The same we do with respect to this sequence of observations as well and we find the following.

(Refer Slide Time: 09:50)

```
By converting the data like what we did for Sign Test, we get the Following:

Sample1: +++++- - \# of Runs = 2

Sample2: -++-++-- \# of Runs = 9
```

With respect to the first one we get that there are 6 plus followed by 6 minus. Therefore, number of runs is 2. However, with respect to the second sample we get minus, plus, plus another minus plus then 1 minus plus minus, minus, plus, plus minus that means the 6 plus and 6 minuses are very nicely intermingled among themselves and we get that the number of runs is equal to 9. Now if I ask you which one of them is random we can easily see that this is definitely random and this is definitely not.

So, while working on numeric data this is one we resolve the dichotomy and thereby we can solve the randomness problem and we can check if the data is actually random or not in the following way. However, if we have a numeric data there is another important test which is called runs up and down test which we will do later, but before that let us just study the common what is called run test. (Refer Slide Time: 11:18)



So, the run test are based on the total number of runs in the sequence.

(Refer Slide Time: 11:23)

```
* Assume an ordered sequence of N elements of two types, A & B
with n<sub>1</sub> of them being A and n<sub>2</sub> are B, where n<sub>1</sub> + n<sub>2</sub> = N.
n<sub>2</sub> = 6
N = 2
* In order to derive a test of randomness we need to obtain the
distribution of R = r<sub>1</sub> + r<sub>2</sub>, where r<sub>1</sub> is the number of runs of A
and r<sub>2</sub> is the number of runs of B.
* Total number of distinct arrangements of N elements having
n<sub>1</sub> elements of type A and n<sub>2</sub> elements of type B, is i.e.
(n<sub>1</sub>+n<sub>2</sub>)
(n<sub>1</sub>+n<sub>2</sub>)
(n<sub>1</sub>+n<sub>2</sub>)
(n<sub>1</sub>+n<sub>2</sub>)
(n<sub>1</sub>+n<sub>2</sub>)
```

Assume an ordered sequence of N elements of two types A and B with n1 of them being A and n2 of them are B such that n1 + n2 = N it is the total number of data points. With respect to the previous example n1 was 6, n2 was 6 and N was 12. In order to derive a test of randomness we need to obtain the distribution of R what is R, R is the total number of runs which is r1 + r2 where r1 is the number of runs of A and r2 is the number of runs of B.

Now what is the total number of possible arrangements of N elements when there are n1 of them are of type A and n2 of them are type B. The number of possible arrangements is

 $^{n1+n2}C_{n1}$ or equivalently $^{n1+n2}C_{n2}$. Why? Some people often asked because arrangement is a permutation problem, but this is a combination. So, how do we interpret a permutation problem using a combination?

The answer is simple because n1 plus n2 many elements can be arranged among themselves in n1 plus n2 factorial. However, n1 of them are same and n2 of them are same. Therefore, we need to divide it by n1 factorial n2 factorial and this is nothing, but n1 plus n2 c n1 which we have written here. Ok

(Refer Slide Time: 13:25)



Now the number of distinct ways of distributing n like objects into r distinguishable cells with no cell empty is ${}^{n-1}C_{r-1}$ when $n \ge r$. This is a very simple combinatorics problem. So, suppose for example I have 7 balls and I want to put them in 3 boxes such that no box is empty. In how many ways one can do that? So, let us first understand that in order to put them in 3 boxes we need to draw some partitions there.

So, let us consider I draw these two partition therefore in the first box there are 2 balls, in the second box there are 2 balls, in the third box there are 3 balls. So, this we can write it as say 2, 2, 3. Another example suppose I have again 7 balls and I choose the following lines then what do I get? We get 1, 3, 3.

Therefore, each selection of two lines out of the 6 empty places between the 7 balls 1, 2, 3, 4, 5, 6, the number of ways of choosing 2 out of 6 that is going to give me 3 partitions each one

is non empty. In general, therefore when there are N objects to be put in R boxes when n is greater than equal to r number of ways of doing it is ${}^{n-1}C_{r-1}$.

(Refer Slide Time: 15:35)



Now we had n1 objects split into r1 runs and n2 objects split into r2 runs. Therefore, number of ways of doing it is n1 -1 C $_{r1-1}$. Number of ways of doing this is $^{n2-1}$ C $_{r2-1}$. Therefore, total number of combinations $^{n1-1}$ C $_{r1-1}$ multiplied by $^{n2-1}$ C $_{r2-1}$. Why you are multiplying? So, let us understand that suppose this is one sequence of blue balls.

Now in order to make these three to be runs in each of this we have to fill in some objects of the different kinds. So, if these are blue balls 1, 2, 3 so to make this run separate we have to insert some run of red balls and if these partitions of 2 runs can be done in k different ways so we can fill it in k different ways. Therefore, the number of ways of making 3 blue ball runs multiplied by the number of ways of making 2 red ball runs will give me 5 different runs.

That is why we see that this comes to be $^{n1-1}$ C $_{r1-1}$ x $^{n2-1}$ C $_{r2-1}$.

(Refer Slide Time: 17:45)



Also note that very interesting thing that r1 and r2 cannot be very arbitrary. Given r2 r1 can take values r2 + 1 or r2 - 1 or r1 can be equal to r2. Say for example, if these are the runs of r2 3 runs then there can be 3 possibilities for r1.; r1 can be 2, r1 can be 3 or r1 can be 4. So, if r2 is equal to 3 r1 can take either 2 or 3 or 4 it has no other possible choice because if there are more element of the first type then actually we will combine them into 1 run.

(Refer Slide Time: 19:03)



So, with these observations we get the following.

(Refer Slide Time: 19:07)



To obtain the probability that R takes a value small r where R is the total number of runs. Therefore, small r can be either even say 2d or odd that is 2d + 1. Now if the total number of runs is 2d that is even then both r1 and r2 have to be equal to d and the sequence may start with either A or B. Similarly, if R is equal to 2d plus 1 that is odd numbers then r1 can be d + 1 and r2 can be d. Or alternatively, r1 is equal to d and r2 is equal to d + 1.

In the first case the sequence starts with A and also ends with A. In the second case it starts with B and also ends with B. Therefore, probability that R is equal to small r is computed as follows.

(Refer Slide Time: 20:22)

The number of ways that
$$r_1$$
 can take a value d is $\binom{n_1-1}{d-1}$
and r_2 can take the value d is $\binom{n_2-1}{d-1}$.
 \therefore The number of ways that $R = 2d$ is $2 * \binom{n_1-1}{d-1} * \binom{n_2-1}{d-1}$
Hence, $P(R = 2d) = \frac{2*\binom{n_1-1}{d-1}*\binom{n_2-1}{d-1}}{\binom{n_1+n_2}{n_1}}$
In a similar way,
 $P(R = 2d + 1) = \frac{\binom{n_1-1}{d}*\binom{n_2-1}{d-1}+\binom{n_1-1}{d-1}*\binom{n_2-1}{d}}{\binom{n_1+n_2}{n_1}}$

Number of ways that r1 can take value of d is ${}^{n1-1}C_{d-1}$ just sometime back we studied this. Similarly, r2 can take value d in ${}^{n2-1}C_{d-1}$ ways. Therefore, the number of ways that R can be two times d is 2 x ${}^{n1-1}C_{d-1}$ x ${}^{n2-1}C_{d-1}$. Why the 2 is coming it is coming because it may start with A and end with B or it may start with B and end with A.

Therefore, probability R is equal to 2d is equal to

 $2 * {}^{n1} {}^{-1}C_{d-1} \times {}^{n1} {}^{-1}C_{d-1}$ divided by ${}^{n1+n2}C_{n1}$ as this is the total number of possible runs. In a similar way probability R is equal to 2d + 1 is equal to ${}^{n1-1}C_d \times {}^{n2-1}C_{d-1}$ that is this is the case when there are d + 1 runs of the first time and d runs for the second time plus

That means d runs of the first type and d + 1 runs of the second type divided by the total number of possibilities that is $^{n1 + n2} C_{n1}$.

(Refer Slide Time: 22:05)



Let us now consider for illustration a case with n1 = n2 = 5.

(Refer Slide Time: 22:14)

$$P(R = 2 = 2 * 1) = \frac{2 * \binom{n_1 - 1}{d - 1} * \binom{n_2 - 1}{d - 1}}{\binom{n_1 + n_2}{n_1}}$$

$$= \frac{2 * \binom{5 - 1}{1 - 1} * \binom{5 - 1}{1 - 1}}{\binom{10}{5}} = \frac{2 * \binom{4}{0} * \binom{4}{0}}{\binom{10}{5}} = \frac{2}{252} = 0.008 \xrightarrow{\text{ABABBBB}}{\text{B} - ... \text{B} - ... \text{A}}$$

$$P(R = 10 (= 2 * 5)) = \frac{2 * \binom{n_1 - 1}{d - 1} * \binom{n_2 - 1}{d - 1}}{\binom{n_1 + n_2}{n_1}} \xrightarrow{\text{ABAB} - ... \text{AB}}{\text{BABA} - ... \text{BA}}$$

$$= \frac{2 * \binom{5 - 1}{5 - 1} * \binom{5 - 1}{5 - 1}}{\binom{10}{5}} = \frac{2 * \binom{4}{4} * \binom{4}{4}}{\binom{10}{5}} = \frac{\binom{2}{252}}{252} = 0.008$$

Therefore, probability that R is equal to 2 that is number of runs is equal to 2 that is coming out to be from this formula 2 upon 252 which is 0.008. Now what is 2? 2 is equal to the number of arrangements that satisfy that total number of run is 2 and we know that it can be only in two ways that is A, A, A, A and A followed by B, B, B, B and B or it can be B followed by A.

There are only two possibilities to have total number of run is 2. In a similar way probability R is equal to 10 that is also coming out to be 2 upon 252 by applying the formula and here also we can justify the same that the number of run is 10 if we get AB, AB up to AB or BA, BA up to BA. So, this formula whatever it is giving is correct from practical justification.

(Refer Slide Time: 23:51)



What about the probability of R is equal to 6. So, the same formula by applying the value of d is equal to 3 we get it is giving us 72 upon 252 which is 0.2857. Question is how do we get 72 many arrangements to give us 6 runs? So, 5 A's can be arranged into 3 runs in 6 different ways. Similarly, 5 B's can be arranged into 3 runs in 6 different ways. Therefore, number of possible arrangements is equal to 2 times 6 into 6 is equal to 72. Thus, from here also we get the correct answer.

(Refer Slide Time: 25:43)

Justification: R = 6. Therefore $r_1 = r_2 = 3$ Now 3 runs of 5 A elements can be in the following six ways: (3 1 1), (2 2 1), (2 1 2), (1 3 1), (1 2 2), (1 1 3) In a similar way 3 runs of 5 B elements can be in six ways. Hence total possibility = 6 * 6 * 2 = 72

Now the question is what are the 6 different ways if we look at we have shown it to you that 5 elements can be put into 3 runs. The possibilities are 3, 1, 1 or 2, 2, 1 and similarly 1, 3, 1

and from here we get 1, 1, 3. So, these are the only 6 ways of splitting 5 elements into 3 runs. So, that is what we have applied there for both A and B and our ultimate answer came out to be 72.

(Refer Slide Time: 26:40)



In a similar way we can handle odd number of runs as well.

(Refer Slide Time: 26:44)

$$P(R = 7 (= 2 * 3 + 1)) = \frac{\binom{n1-1}{d} * \binom{n2-1}{d-1} + \binom{n1-1}{d-1} * \binom{n2-1}{d}}{\binom{n1+n2}{n1}}$$
$$= \frac{\binom{5-1}{3} * \binom{5-1}{3-1} + \binom{5-1}{3-1} * \binom{5-1}{3}}{\binom{10}{5}} = \frac{2 * \binom{4}{3} * \binom{4}{2}}{\binom{10}{5}}$$
$$= \underbrace{48}{252} = 0.190$$

For example consider R is equal to 7 therefore by applying the formula we get the number of such runs is 48 so that probability of 7 runs is 0.190.

(Refer Slide Time: 27:04)



Now the question is how do we know that there will be 48 runs? That comes from this calculation. Just now we have seen that probability of 3 runs of 5 elements is 6. Now getting 4 runs out of 5 elements there are 4 possibilities. We can have them 2 then 1, 1, 1 or 1, 2 then 1, 1, 1 and similarly from the other two you can understand the partition into 5 runs. Therefore, 6 from here and 4 from here.

Therefore, total number of runs is 2 into 6 into 4 is equal to 48. The 2 is coming because we may have 4 runs of A and 3 runs of B or alternatively 4 runs of B and 3 runs of A. So, that is why the factor 2 has come, but at the end this justifies why the formula gave us that there are 48 many different possibilities of arranging the elements such that total number of run is equal to 7.

(Refer Slide Time: 28:34)

	Thus we have : $P(R = 7) = 0.190$
	In similar way we have:
	$P(R = 8) = \frac{2 * \binom{5-1}{4-1} * \binom{5-1}{4-1}}{\binom{10}{5}} = \frac{2 * \binom{4}{1} * \binom{4}{1}}{\binom{10}{5}} = \frac{32}{252} = 0.127$
	$P(R = 9) = \frac{\binom{5-1}{4} * \binom{5-1}{4-1} + \binom{5-1}{4-1} * \binom{5-1}{4}}{\binom{10}{5}} = \frac{2 * \binom{4}{4} * \binom{4}{3}}{\binom{10}{5}} = \frac{8}{252} = 0.032$
F	$P(R = 10) = \frac{2 * \binom{5-1}{5-1} * \binom{5-1}{5-1}}{\binom{10}{5}} = \frac{2 * \binom{4}{4} * \binom{4}{4}}{\binom{10}{5}} = \frac{2}{252} = 0.008$

Therefore, we have probability R is equal to 7 is 0.190. In a similar way, we can calculate that probability R is equal to 8 is 0.127 probability R is equal to 9 is 0.032 probability R is equal to 10 that we have already obtained is 0.008.

(Refer Slide Time: 29:04)

Hence we have the following Right Tail probabilities: 1. $P(R \ge 10) = 0.008$ 2. $P(R \ge 9) = 0.008 + 0.032 = 0.040$ 3. $P(R \ge 8) = 0.008 + 0.032 + 0.127 = 0.167$ 4. $P(R \ge 7) = 0.008 + 0.032 + 0.127 + 0.190 = 0.357$ These values are needed for acceptance and rejection of Ho.

Now if we look at the right tail probabilities we can see that by adding the terms probability $R \ge 10$ is 0.008, $R \ge 9$ is 0.040, $R \ge 8$ is 0.167 and $R \ge 7$ is 0.357. Why did we calculate that because we need these values to accept or reject the null hypothesis.

(Refer Slide Time: 29:44)

							Righ	t-tail p	robal	bilitie	8					
	n1	n_2	R	Р	n_1	n_2	R	P	n_1	n_2	R	Р	n_1	n_2	R	Р
	2	2	4	.333	4	8	9	.071	5	11	11	.058	6	12	12	.075
	2	3	5	.100			8	.212			10	.154			11	.217
			4	.500			7	.467			9	.374			10	.395
1	2	4	5	.200	4	9	9	.098	5	12	11	.075	6	13	13	.034
	2	5	5	.286			8	.255			10	.181			12	.092
	2	6	5	.357	4	10	9	.126	5	12	9	.421			11	.257
1	2	7	5	.417			8	.294	5	13	11	.092			10	.439
	2	8	5	.467	4	11	9	.154			10	.208	6	14	13	.044
1	3	3	6	.100			8	.330			9	.465			12	.111
			5	.300	4	12	9	.181	5	14	11	.111			11	.295
	3	4	7	.029			8	.363			10	.234			10	.480
			6	.200	4	13	9	.208	5	15	11	.129	7	7	14	.001
			5	.457			8	.393			10	.258			13	.004
1	8	5	7	.071	4	14	9	.234	6	6	12	.002			12	.025
			6	.286			8	.421			11	.013			11	.078
1	3	6	7	.119	4	15	9	.258			10	.067			10	.209
			6	.357			8	.446			9	.175			9	.383
1	3	7	7	.167	4	16	9	.282			8	.392	7	8	15	.000
			6	.417			8	470	6	7	13	.001			14	.002
	3	8	7	.212	5	5	10	.008	χ,		12	.008			13	.012
			6	.467	12	2	9	.040	2		11	.084			12	.051
1	8	9	7	.255			8	.167	1		10	,121			11	.133
1	3	10	7	.294			7	.357	1		9	.267			10	.296

And for that we need to look at a table with right tail probabilities. So, here is the table and we are considering number of A and number of B is to be 5 and 5. So, if we look at we can see that the right tail probabilities for 10 is 0.008 for 9, 0.04 for 8 it is 0.167 and for 7 it is 0.357. Therefore, we see that we have obtained the same values through our calculations using the formula for obtaining the run probabilities.

(Refer Slide Time: 30:33)



The above table giving the right tail probabilities is taken from the book nonparametric statistical inference written by Gibbons and Chakraborti.

(Refer Slide Time: 30:45)

	Large sample Approximation.
	When either n_1 or n_2 are greater than 20 or both are greater than 10, one can approximate with Normal distribution with the following
	Mean and Variance: $m = \gamma^{0}$
	$E(R) = \frac{2*n1*n2}{n1+n2} + 1 \qquad \qquad$
	$Var(R) = \frac{2*n1*n2*(2*n1*n2-n1-n2)}{(n1+n2)^{2}*(n1+n2-1)}$
	The above test is known as Wald-Wolfowitz Run Test
()	

Now large sample approximations when either n1 or n2 are greater than 20 or both are greater than 10 then one can approximate with normal distribution with the following parameters. The mean is coming out to be

 $\frac{2*n1*n2}{n1+n2} + 1$

Note that we are using the notation m and n for the number of elements for the two populations here it is given n1 and n2.

So m is actually the n1 and n is actually the n2. Therefore, variance of R is given as

$$\frac{2*n1*n2*(2*n1*n2-n1-n2)}{(n1+n2)^2*(n1+n2-1)}$$

This is a complicated expression we are not going to derive it derive it here in this class. Let us accept this formula the test that we have described so far based on the runs is called Wald Wolfowitz run test. (Refer Slide Time: 32:08)



So, some simple mathematical observations on the number of runs.

(Refer Slide Time: 32:15)

0	
	Q1. What is the marginal distribution of \mathbb{R}_1 ? $\min_{max} 1$ Given n_1 and n_2 the random variable \mathbb{R}_1 can take integral values in the set $\{1, 2,, n_1\}$ To compute $\mathbb{P}(\mathbb{R}_1 = \mathbf{r}_1)$ when $\mathbf{r}_1 \in \{1, 2,, n_1\}$
MPTEL	

What is the marginal distribution of R1 where if you remember R1 is the number of runs of the first element. We know that the minimum value for R1 is 1 and the maximum value of R1 is n1 that is the number of elements. Therefore, given n1 and n2 the random variable R1 can take value integral values in the set 1, 2, 3 up to n1. To compute probability R1 is equal to small r1 we do the following.

(Refer Slide Time: 32:59)



If R1 of small r1 is equal to the joint distribution of r1, r1 in r1, r2 that means the number of runs of the first element and number of runs of the second element are same. We have already seen that this value can be same or +/-1 of R1. Therefore, we also add to it the probability of f_{R1R2} (r1, r1 -1) and similarly r1 and r1 + 1. We know all these values this we have already computed.

This is ${}^{2n1-1}$ C ${}_{r1-1}$ x by ${}^{n2-1}$ C ${}_{r1-1}$ similarly from this we get this and from here we get this. This we are getting by replacing the value of r2 with the corresponding parameter in the relevant expression. Now we need to simplify this. How we simplify? (Refer Slide Time: 34:14)



So, let us look at it. So $^{n1 + n2} C_{n1} x$ by f_{R1} (r1) that is this is the marginal distribution or marginal probability. Therefore, it is writing as

$$2 * \binom{n_1 - 1}{r_1 - 1} * \binom{n_2 - 1}{r_1 - 1} + \binom{n_1 - 1}{r_1 - 1} * \binom{n_2 - 1}{r_1 - 2} + \binom{n_1 - 1}{r_1 - 1} * \binom{n_2 - 1}{r_1}$$

Now let us take ${}^{n2-1}C_{r1-1}$.common and we write two times this as summation of this. Therefore we are getting ${}^{n2-1}C_{r1-1}$. Another one we write it here ${}^{n2-1}C_{r1-1}$ ${}^{n2-1}C_{r1-2}$ + ${}^{n2-1}C_{r1}$. Now summation of ${}^{n2-1}C_{r1-1}$ and ${}^{n2-1}C_{r1-2}$ = ${}^{n2}C_{r1-1}$. So, we know that ${}^{n}C_{r}$ + ${}^{n}C_{r-1}$ = ${}^{n+1}C_{r}$. So, ${}^{n2-1}C_{r1-1}$ + ${}^{n2-1}C_{r1-1}$ by applying this formula is ${}^{n2}C_{r1-1}$.

In a similar way this is giving us $^{n2}\mathrm{C}_{r1}$. Therefore the whole thing comes out to be

 $^{n1-1}C_{r1-1}$. Again we are applying the same formula we get $^{n2+1}C_{r1}$.

(Refer Slide Time: 36:34)

$$\therefore f_{R_1}(r_1) = \frac{\binom{n_1 - 1}{r_1 - 1}\binom{n_2 + 1}{r_1}}{\binom{n_1 + n_2}{n_1}}$$

Thus we get the marginal probability of R_1 taking the values r_1 , where possible values for r_1 are $1, 2, ..., n_1$.
In a similar way, marginal probability of R_2 can be computed for $r_2 = 1, 2, ..., n_2$
We accept the following theorem without proof.

Therefore

$$f_{R_1}(r_1) = \frac{\binom{n_1 - 1}{r_1 - 1} \binom{n_2 + 1}{r_1}}{\binom{n_1 + n_2}{n_1}}$$

Thus, we get the marginal probability of r1 taking the values small r1 in the set 1, 2, 3 up to n1.

In a similar way, the marginal probability of r^2 can be computed for r^2 is equal to 1, 2 up to n^2 and it does not take much to understand that it is going to be

$$f_{R2} \left(\ r2 \ \right) \ = \ \left(\ ^{n1+1}C \ _{r2} \ \textbf{X} \ ^{n2} \ ^{-1}C \ _{r2} \ _{-1} \ \right) \ / \ ^{n1+n2}C_{n1}.$$

Now we accept the following theorem without proof.

(Refer Slide Time: 37:51)



The probability of distribution of R that is the total number of runs with n1 objects of type A and n2 objects of type B in a random sample is given by this complicated formula when r is even and this is when r is odd. It looks complicated, but we are familiar with this ^{n1 -1} C r by 2 minus 1 because it is even say 2d. Therefore, this is actually the d that we are talking about and here also this is the d that we are talking about. In a similar way when r is odd it is 2d plus 1 so one has to be d other has to be d plus 1 and by replacing 2d plus 1 by r we can get this formula. I want you to verify that because all the necessary mathematics we have already studied.

(Refer Slide Time: 38:49)

Observation Suppose we define $n_1 + n_2$ random variables: $R_{11}, R_{12}, \cdots, R_{1n_1}$ and $R_{21}, R_{22}, \cdots, R_{2n_2}$ Where, R_{ii} denotes the number of runs of type *i*, (where i = 1, 2) of length j, $j = 1, 2, \cdots, n_i$ Note that Note that $\begin{array}{c} \sum_{j=1}^{n_1} jr_{1j} = n_1, \quad \sum_{j=1}^{n_2} jr_{2j} = n_2 \quad z_1, z_2, \cdots \\ R_{11} = \# r_{11} \\ R_{12} = \# r_{12} \\ R_{11} = \# r_{12} \\ R_{12} = \# r_{12} \\ R_{12} = \# r_{12} \\ R_{12} = \# r_{12} \\ R_{13} = \# r_{13} \\ R_{13} = r_{$

Also there is another interesting observation with respect to the runs is that if we define n1 plus n2 random variables R11, R12, R1n1 and R21, R22, R2n2 where Rij denotes the number

of runs of type i where i is equal to 1, 2 of length j that is j is equal to this should be j 1, 2, 3 up to ni. What does it mean? That means the number of possible values for R1 is equal to 1, 2, 3 up to n1.

So, R11 is number of runs of length 1 R12 is equal to number of runs of length 2 up to R1 n1 is equal to number of runs of length n1. This is with respect to element of first type that is A. Similarly, for B what we will get R21 is equal to number of runs of length 1 up to R2 n2 is equal to number of runs of length n2 that is how we get R11. R12 up to R1 n1. R21. R22. R2n2 so many variables that is n1 plus n2 too many variables.

But let us note that if R1n1 it cannot take any value more than 1 because if the length of the run is n1 that means that all A observations are in that run therefore there can be only at most one run of length n1. Similarly, there can be at most one run of length n2 with respect to the elements of type B. Now note that sigma j r1j j is equal to 1 to n1 is equal to n1 and j is equal to 1 to n2 summation j r2j is equal to n2. Why? Since r1j is the number of runs of length j therefore basically we are looking at number of runs of length j multiplied by the j that is the length of that run.

Therefore, when you sum over all the possible values we get the total number of elements which is n1 and in this case which is n2. Also note that sigma over j 1 to n1 r1j is the total number of runs of the first element and similarly this expression gives the total number of runs of the second element.

(Refer Slide Time: 42:40)



Let me now describe a variation of the run test which is called runs up and down method this is applicable only for numerical data.

(Refer Slide Time: 42:54)



Suppose, we have a time series data and one may like to check if there is any trend in that data. For example, the number of cars passing through a corridor every day, it can be number of infected cases of Corona Virus on a daily basis. So, if we take the data for different time instance which in this case are the days. Suppose, the data is increasing then we can find a trend like this which is an upward trend or we may find a downward trend which may look like this.

And we are often interested in testing if there is any trend in the data or it is absolutely random. If we just want to test a randomness we can split the data into two parts. Some of them are plus and some of them are minus depending upon whether the sample value is above or below the median. This we have seen.

(Refer Slide Time: 44:13)

Cor	isider the following Tab	ole										
	Time	1	2	3	4	5	6 ~					
	Sample Value	8	13	1	3	4	7 🗸					
	Symbol	+	+/	-	-	-	+~					
But	(San in this way some infor ent in identifying patter	nple n matic rn in f	medi on is l the ti	an i ost me -	5.5 which orde) n can l red ob	be used to some oservation.					
Ru	Runs Up and Down uses this information for more detailed analysis.											

For example, consider the table if we have 6 time instance and the observed values are like this therefore the sample median is coming to be 5.5 and therefore what we are looking at that these 3 are above the sample median and these 3 are below the sample median. However, in this type of test we lose the information about the trend because we have done similar thing with respect to sign test. Runs up and down test preserves the information about the trend.

(Refer Slide Time: 44:58)

1	Here, given n time-ordered n – 1 element sequence	d sample n-1 of +	<i>x</i> ₁ , and	x ₂ ,. 1 –	wh	x _n ose	we con j th ele	struct a ment is,	n	
	$\begin{cases} + if x_{j+1} - \\ - if x_{j+1} - \end{cases}$ Thus for the above dat	$x_j > 0$ $x_j < 0$ a we get	the	fol	llow	vinį	, table	:		
	m	•	•			-				
	Time	1	2	3	4	5	0			
	Sample Valu	ue 8	13	1.	3	4	7			
~	Sequence of Symbols (D	; ;) ↓	J	*	*/	+/				
(*)		ux								

How it works. So, here what we do given n time ordered sample x1, x2, xn we construct an n minus 1 element sequence which you call D n minus 1 or plus and minus, but whose jth element is plus if xj plus 1 minus xj is greater than 0 and it is minus if xj plus 1 minus xj is

less than 0. Therefore, from the same data now we get a sequence of 5 symbols, but it is different from the one we get for the sign test.

Now why this is plus because at time instance 1 the value is 8 at time instance 2 the value is 13 therefore there is an increment hence it is plus from 2 to 3 the value has come down from 13 to 1 therefore it is minus. From 3 to 4 value goes from 1 to 3 therefore it is plus here it is 3 to 4 plus and here it is 4 to 7 therefore another plus. Therefore, we get a different sequence here.

(Refer Slide Time: 46:16)

Under the null hypothesis of randomness H_0 the + and are expected to be nicely intermingled, so that there will not be a longer run of any of the two symbols. Thus here focus is NOT on the number of runs. Rather here we look at the length of the runs. A long run of + (or -) implies that the data is showing some clear trend. / Hence H_0 will be rejected if there are r runs of length t or more, where r and t are decided on the basis of α .

Therefore, under the null hypothesis of randomness the plus and minus are expected to be nicely intermingled so that there will be not be longer run of any of the two symbols because if there is a longer run of plus, plus, plus like that what we understand that for the entire time period the value is increasing that means there is an upward trend and therefore the value is not random.

Therefore, here we look at the length of the runs not the total number of runs that we have just seen with respect to Wald–Wolfowitz runs test. If there is a long run of plus or minus it implies that the data is showing some clear trend. Therefore, H naught will be rejected if there are r runs of some length t or more where r and t are decided on the basis of alpha that means if there are large number of runs of certain length or more than that then we are going to reject the randomness of the data. (Refer Slide Time: 47:41)

\odot	The Analysis
	The method considers the number of runs of length j ,
	$j = 1, 2, \cdots, n - 1$
	Let R_j be the random variable denoting the no. of runs,
	either up or down of length j in the sequence $D_{(n-1)}$
	+++ +
	For illustration, let $n = 3$.
	Therefore, there can be runs of length 1 or length 2.
	Now three numbers can be arranged among themselves
	in 3! = 6 possible ways each one is equally likely.
	For each of them D_2 will contain runs of + and/or -

So, the method considers the number of runs of length j, j is equal to 1, 2, 3 up to n minus 1 because the length of a run can be 1 can be 2 or maximum it can be n minus 1 because if the entire data is showing an upward trend there will be one run of length n minus 1 all plus. For illustration let n be equal to 3. Therefore, there can be runs of length 1 or 2. Now the three numbers can be arranged among themselves in 6 possible ways each one of them is equally likely.

Therefore, since n is equal to 3 you are looking at D2 and D2 will contain runs of plus or minus or maybe mixture of that one.



(Refer Slide Time: 48:38)

So, let us illustrate suppose we have three elements a1, a2, a3. A1 is the value obtained at the first instance, a2 is the value obtained at the second instance and a3 is the value obtained at the third instance. Now if a1 is less than a2 less than a3 then D2 is going to be plus because a2 is greater than a1 and another plus because a3 is greater than a2. Therefore we get a run of length 2.

Therefore, what we are saying r1 is equal to 0 r2 is equal to 1 so what is r1? R1 is equal to number of runs of length 1 R2 is equal to number of runs of length 2 and therefore we get 0, 1. If a1 is less than a3 less than a2 therefore from a1 to a2 we get a plus from a2 to a3 we get a minus therefore we get 2 runs of length 1 therefore r1 is equal to 2 and r2 is equal to 0. We know that there will be 6 possible permutations. For each of them we calculated the value of r1, r2 and we are getting that two instances of 0, 1 and 4 instances of 2, 0 as it is clear from here.

(Refer Slide Time: 50:34)



Corresponding to each arrangement we can compute the value of r2 and r1. Hence, r2 and r1 represents the number of runs of length 2 and 1 respectively as I have indicated before and therefore f3 0, 1 this is the notation which is actually the probability that we get a 0, 1 sequence that probability is 1 by 3 and f3 of 2, 0 that means we get 2 runs of length 1 and 0 runs of length 2 that probability is going to be 2 by 3 as we have seen just now.

(Refer Slide Time: 51:20)



Now let us generalize it if there are n samples we can think of fn that is the probability of r1, r2, rn minus 1 where r1 is the number of runs of length 1, r2 is the number of runs of length 2 and rn minus 1 is the number of runs of length n minus 1.

(Refer Slide Time: 51:44)



It is obvious that rn minus 1 cannot be more than 1 because if there is run of n minus 1 many elements whether plus or minus that means it has taken care of all the transitions and therefore there cannot be any other element in the sequence therefore there can be only one run of length n minus 1, but if rn minus 1 is equal to 1 then all other values are 0. However, if rj is greater than 0 for any j not equal to n minus 1 then rn minus 1 is equal to 0.

Because then it can have a run of length n minus 1 if r1 is equal to n minus 1 where all others are going to be 0 because we have got n minus 1 many runs of length 1 that means it has to be plus, minus, plus, minus or something like minus, plus, minus, plus like that. R1 can never be equal to n minus 2 which is pretty obvious and sigma over jrj where j is the length of the run and rj is the number of runs of length j. So, that gives us the total number of symbols in the run and that is going to n minus 1.

(Refer Slide Time: 53:15)



However, the basic question is how to compute this probability fn, r1, r2, rn minus 1. So, in order to compute these probabilities we need to count the number of possible arrangements of n minus 1 sign of plus or minus and which leads to definition of n minus 1 random variables where R1 is the random variable denoting the number of runs of length 1. Similarly R2 and R2 Rn minus 1 and we are looking at the probabilities of these variables taking the value small r1 small r2 up to small rn minus 1.

So, let un r1, r2, rn minus 1 be the number of possible arrangements so that we can get the runs of these numbers. Therefore, fn r1, r2, rn minus 1 is going to be un r1, r2, rn minus 1 divided by n factorial because that is the number of way of arranging n many elements.

(Refer Slide Time: 54:34)



For illustration with respect to the earlier example u3 0, 1 was 2 u3 2, 0 was 4 therefore f3 0, 1 was coming out to be 1 by 3 f3 0, 2 was coming out to be 2 by 3.

(Refer Slide Time: 54:54)

S ₃	D_2	r_1	r_2
$(a_1 < a_2 < a_3)$	(+,+)	0	1
$(a_1 < a_3 < a_2)$	(+,-)	2	0
$(a_2 < a_1 < a_3)$	(-,+)	2	0
$(a_2 < a_3 < a_1)$	(+,-)	2	0
$(a_3 < a_1 < a_2)$	(-,+)	2	0
$(a_3 < a_2 < a_1)$	(-,-)	0	1

As we have seen from this table sometime back.

(Refer Slide Time: 54:59)



Now question is how do we compute all these un's basically there is a recurrence relation which helps us to get un from un minus 1.

(Refer Slide Time: 55:12)



So, suppose we have n minus 1 observations a1, a2, an minus 1 and extra observation an is added. Now addition of this observation can either split an existing run or increase the length of a run or it can introduce a new run of length 1. For illustration, consider 4 data points a1, a2, a3, a4 and such that it is an increasing order therefore in D3 we will have plus, plus, plus because from here to here plus, here to here plus and from here to here there is another plus.

(Refer Slide Time: 56:00)

	No	ĥs	in the Cargura						
		S ₅	D ₄	<i>r</i> ₁	r_2	r_3	r_4		element
	1	$(a_1, a_2, a_3, a_4, a_5)$	+++++++++++++++++++++++++++++++++++++++	0	0	0	1		
	2	$(a_5, a_1, a_2, a_3, a_4)$	Ott	1	0	1	0		
	3	$(a_1, a_5, a_2, a_3, a_4)$	tott	2.	1.	0.	0.		
	4	$(a_1, a_2, a_5, a_3, a_4)$	++ -+	2	1	0	0		
	5	$(a_1, a_2, a_3, a_5, a_4)$	+++⊖	1	0	1	0		
			-						
	So,								
	In c								
(*)	In c	ases 3, 4, an existing ru	n of length 3 is	split.					

Now, suppose the new element a5 is such that it is the largest of all then what will happen we will get another plus here on the top of that 3 plus and therefore it is increasing the length of this run of plus. Instead if the fifth element comes at the very first instance then here we will get a minus because a1 is less than a5, but after that there is going to be plus, plus, plus. Similarly, if a5 comes at the second instance then we will get plus, minus, plus, plus.

And similarly for all the 5 cases note that here we have assumed that a5 is the largest element. Therefore, introduction of a new element will change the pattern of the runs of plus and minus. So, as we have written in case of 2 and 5 the insertion introduced a new run of length one. As you can see here we get a minus, here we get a minus, but in case of 3 and 4 an existing run is split into two parts.

Because here it was plus, plus, plus that is split into plus, minus, plus, plus and accordingly you can see that the values are changed as we have indicated now it is taking as 2 runs of length one, one run of length 2 and no runs of length 3 and 4. So, the entire calculation is run here let me not elaborate that.

(Refer Slide Time: 57:54)

Note that



Let me go further. Now if the initial arrangement of a1, a2, a3 and a4 is different instead of that they are increasing and then the effect of introducing a5 will also be different as the arrangement of the plus and minus will depend upon the relative order of a1, a2, a3, a4 and also a5. Hence, we need to look for a generalization in fact we can identify for mutually exclusive and exhaustive cases as mentioned below.

(Refer Slide Time: 58:32)

So, suppose this n minus 1 denote the number of arrangements with n minus 1 many elements of plus and minus then an additional run of length 1 can be added in the arrangement of Sn because of the introduction of the new element that we have already seen or a run of length i

minus 1 in Sn minus 1 can be changed into a run of length i in Sn that means that an existing run is now increased by length 1 or it can be split into two parts.

Here we have divided that also into two different segments. Suppose h is 2i that means it is an even number and the new element comes at the middle so it may split it into run of length i followed by a run of length 1 followed by a run of length i that means suppose we had something like this and the new element converts it into plus, plus, minus, plus, plus. So, this is the case that is taken care of in 0.3,.

(Refer Slide Time: 59:52)

4) A run of length h = i + j in S_{n-1} , where $h > i > j, 3 \le h \le n - 2$ can be split up into a) A run of length *i*, followed by a run of length 1, followed by a run of length j. b) A run of length j, followed by a run of length 1, followed $\frac{1}{2}$ = by a run of length i. ++++ The overall recurrence relation between u_{n-1} and u_n can be obtained from this formula:

But more generally if the length is h which is i plus j i and j are different numbers then it can split it into a run of length i followed by a run of length 1, followed by a run of length j that means the split is not equal for example we can have plus, plus, plus, plus divided into three parts plus followed by a minus followed by three plus. So here, i is equal to 1 j is equal to 3. On the other hand it can also split it like plus, plus, plus, plus it goes to plus, plus, plus, minus, plus.

Therefore, i is equal to 3 j is equal to 1. The overall recurrence relation between un minus 1 and un can be obtained from the following formula.

(Refer Slide Time: 60:59)

$$\begin{split} u_n(r_{n-1},r_{n-2},\cdots,r_h,\cdots,r_i,\cdots,r_j,\cdots,r_1) \\ &= 2u_{n-1}(r_{n-2},\cdots,r_1-1) \\ &+ \sum_{i=2}^{n-1} (r_{i-1}+1) \; u_{n-1}(r_{n-2},\cdots,r_i-1,r_{i-1}+1,\cdots,r_1) \\ &+ \sum_{\substack{i=2\\(h=2i)}}^{n-2} (r_h+1)u_{n-1}(r_{n-2},\cdots,r_h+1,\cdots,r_i-2,\cdots,r_1-1) \\ &+ 2\sum_{i=2}^{n-3} \sum_{j=1}^{i-1} (r_h+1)u_{n-1}(r_{n-2},\cdots,r_h+1,\cdots,r_i-1,\cdots,r_j-1,\cdots,r_1-1) \end{split}$$

As you can understand this is pretty complicated we are not deriving it here neither I am asking you to check the correctness of this things. It is only for you to conceptualize how runs up and down works.

(Refer Slide Time: 61:19)



Now Levene and Wolfowitz have computed the means and variance of number of runs of length t or more. For n less than equal to 14 the exact probabilities of getting at least r runs of length t or more are also available and they are tabulated and from there the critical regions can be constructed and H naught can be tested against suitable alternatives.

(Refer Slide Time: 61:54)



However, when the number of values is large we can go for normal approximation for the same and the calculation is the following.

(Refer Slide Time: 62:02)



The mean is coming out to be 2n minus 1 upon 3 and the variance is coming out to be 16n minus 29 upon 90. As you can understand that these are little bit complicated so I am not going to derive the formulae, but if n is greater than 25 and if we incorporate the continuity correction then we can get the left tail and right tail critical regions by using normal approximation as follows that.

(Refer Slide Time: 62:34)



For the left tail it is R plus 0.5 minus the mean divided by the standard deviation has to be less than equal to minus z alpha where Z α is the critical value. For normal distribution similarly for right hand tail it is going to be R- 0.5 minus 2n - 1 which is the mean divided by the standard deviation. So, this is for one sided and if you want to do it two sided then replace Z α by Z α /2.

Okay friends I stop here today. I think you have understood the intuition behind using runs for testing whether a given data is random or not. I stop here and in the next class I shall study some more test based on two samples. In particular I will look at Kolmogorov Smirnov two sample test and median test for testing the equality of distribution. Okay friends. Thank you.