**Statistics for Experimentalists**
**Prof. Kannan. A**
**Department of Chemical Engineering**
**Indian Institute of Technology – Madras**

**Lecture – 26**
**Analysis of Experiments involving Single Factor - Part B**

Welcome back. In today's class or rather classes we will be completing our discussion on single factor experimentation. Most likely we will not be doing experiments by changing only one factor or only one variable and keeping all other variables constant. This is simplified representation and it will involve a lot of definitions and notation which will help us in our further discussions on 2 factors and multiple factors experimentation.

**(Refer Slide Time: 01:09)**

**References:**

Mathews, P., Design of Experiments with Minitab. New Delhi: New Age International, 2010.

Montgomery, D. C., G.C. Runger, Applied Statistics and Probability for Engineers. 5th ed. New Delhi: Wiley-India, 2011.

So, let us see the reference books in addition to Montgomery and Runger's Applied Statistics and Probability for Engineers book. We also have the book by P. Mathews on Design of Experiments with MINITAB.

**(Refer Slide Time: 01:28)**

## Contents

❖ Review of ANOVA

❖ Further Analysis of Treatment Means

❖ Confidence Intervals

❖ Fischer's LSD

❖ Blocking and Randomization

So in today's classes we will be looking at the review of the analysis of variance. We will be further analyzing the treatment means. We will be looking at confidence intervals. Fisher's least square difference LHD blocking and randomization.

**(Refer Slide Time: 01:52)**



## Analysis of Variance Table (ANOVA)

| Source of variation | Sum of Squares | Degrees of Freedom | Mean Square | $F_o$ |
|---|---|---|---|---|
| Treatments | $SS_{treatments}$ | $a-1$ | $MS_{treatments}$ | $MS_{treatments}/MS_{error}$ |
| Error | $SS_{Error}$ | $a(n-1)$ | $MS_{error}$ | |
| Total | $SS_T$ | $an-1$ | | |

We already looked at the analysis of variance table previously. It is divided into a set of 4 columns. We have the source of variation in which we have the sum of squares contribution from treatments. We know by now what treatments are it is different levels of the factor which is applied during the experimental work. It may be fertilizer A, fertilizer B, fertilizer C and their effect on the crop output is monitored.
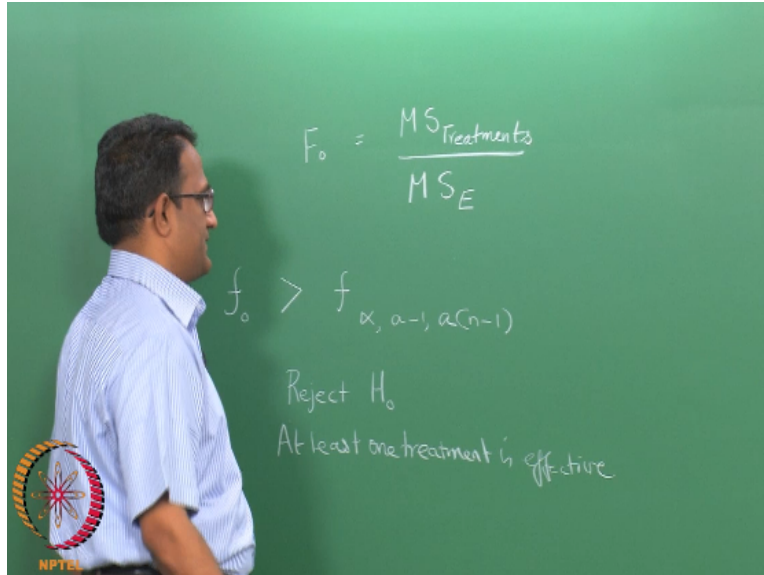
Then there is a contribution from the experimental error to capture the experimental error we do repeats of the experiments. And then we have the degrees of freedom. If you are having a treatments, a levels of the factors there are a-1 degrees of freedom only because out of these a-1 treatment levels are independent because we are getting the treatment mean from the different treatments themselves so we have only a-1 independent quantities.

The degrees of freedom for the error are for a treatments you have n-1 independent repeats because the repeats are also added up to get the average. So, we divide the sum of squares by the degrees of freedom to get the means square treatments. We divide the sum of squares of the error by degrees of freedom for the error to get means square error. The ratio of these 2 mean squares give us the F value.

And then we can do the F test to see whether the F statistics is laying in the critical region of in the acceptance region. If alpha is the chosen level of significance, alpha is usually .05 but you may want to choose other values of alpha as well not 1., not 2.1 etcetera and if you find F alpha a-1 numerator degrees of freedom a*n-1 denominator degrees of freedom is such that this F value is exceeding that F alpha a-1, a*n-1 then you reject the null hypothesis.

The F value if laying in the rejection region if this F not value is > F alpha a-1, a*n-1. I just write it on the board in case you did not follow my statement.

**(Refer Slide Time: 05:16)**

Finding this value of F alpha mean square and then you also find the value f alpha and once this calculation are done you get a value of f not. And if this value if greater then f alpha what is the degrees of freedom for the treatments? You have a treatments, you have n repeats so you have n-1, a treatments a *n-1. So, if this f not is greater than f alpha a-1, a*n-1. The null hypothesis is rejected at least one treatment is effective over and above the random noise that has taken care of that.

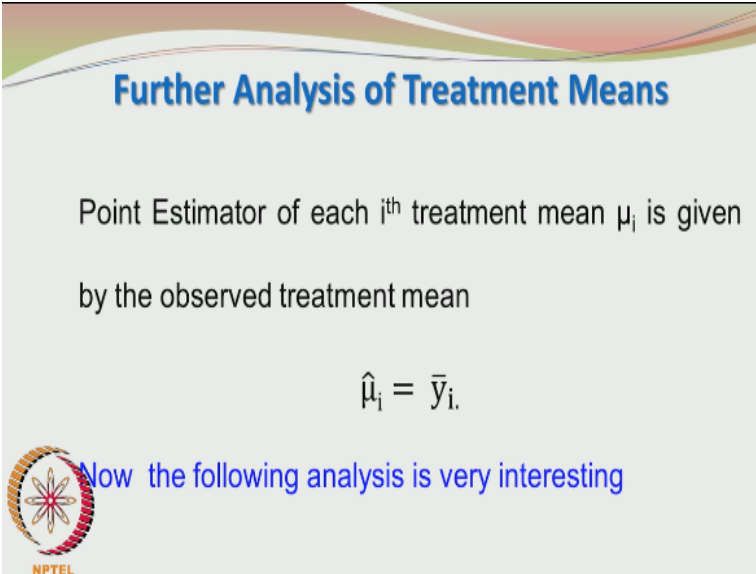**(Refer Slide Time: 06:36)**



And error variance actually this is the central point of our experimental work. The entire design of experiments and analysis will collapse if we are not able to resolve the error variance issue properly. We do a pooled standard deviation. What is this pooled standard deviation mean? You

might have heard about pooled car service to save on petrol people travel by one vehicle to their offices from their home.

So, what we do by pooling is we calculate the mean square error from all the repeats. Even though we are comparing for a particular treatment the errors are pooled from all the repeats, for all the treatments. You pool the repeat values across all the treatments. So, the mean square error is used as a representation of the error variance and more data points you have the better would be your estimate of the error variance sigma square.

So, we are not talking a particular treatment repeat set and then using that as the error we are taking all the repeats and finding the mean square error. So the degrees of freedom is a*n-1.

**(Refer Slide Time: 08:23)**

## Further Analysis of Treatment Means

Point Estimator of each $i^{th}$ treatment mean $\mu_i$ is given

by the observed treatment mean

$$\hat{\mu}_i = \bar{y}_{i.}$$
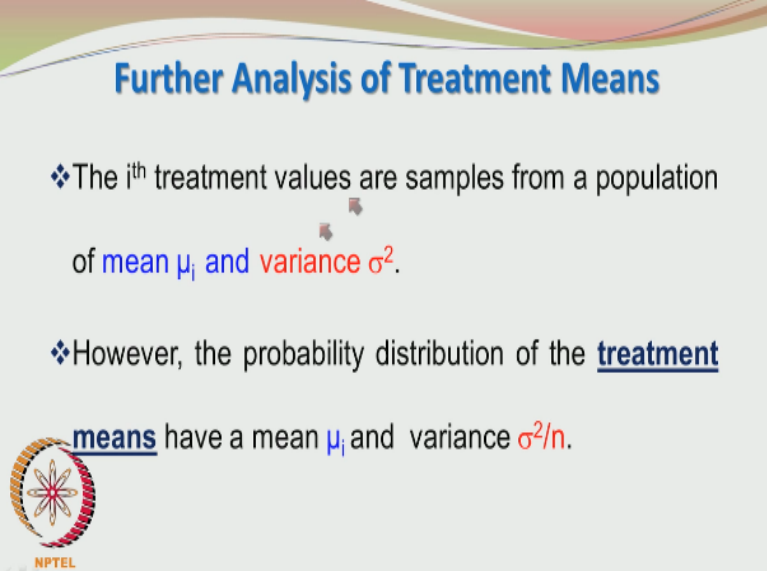
Now the following analysis is very interesting

NPTEL

Now this becomes quite interesting because whatever we studied for the (()) (08:27) distribution, the confidence intervals etcetera will find immediate application here. Now, we know that the treatment mean is given by mu i for the i'th treatment. So the point estimator of each of the i'th treatment mu i is given by the observed treatment mean. So, we do not know the exact value of mu y so we carry out the experiment, collect sum data.

And then we find the treatment mean y bar i dot. This bar represents the averaging, i represent the i'th treatment and dot represent the totalling across all the j values. j running from 1 to n, the j

representing the repeated measurements. This we have also seen in the previous class. So the estimator is indicated by this hat so the estimator for mu i the i'th treatment mean is nothing but the experimental observed treatment mean y bar I dot.

This again we have studied in our point estimation may be about 10 lectures previously. So essentially whenever you are doing experiments you are sampling the population space.

**(Refer Slide Time: 10:23)**



So, the i'th treatment values are samples from a population of mean mu i and variance sigma square. So, the probability distribution of the treatment means which we also should be familiar by now is having the parameters mean mu i and variance sigma square/n. So, we are having a population of mean, mu i and variance sigma square. We are taking sample out of it and taking average and these averages would also be having a probability distribution.

Everytime we do sampling from a population and calculate the respective sample means we may not get the same value each sample may give a different value of the mean. Each sample may have its own mean that is what I am trying to say. So, different samples will have different sample means and hence we may describe a distribution of the sample means. And it so happens that the center value or the average value of this distribution of sample means is the population mean itself.

However, the spread of the means is reduced when compared to the spread of the population values. So, the sample mean distribution spread is actually dampened by the sample size. So, the sample mean distribution will have a variance sigma squared by n and we always want the spread to be small and hence we want the n to be large.

**(Refer Slide Time: 12:43)**

## T-Test

Define the T random variable as follows

$$T_i = \frac{\bar{y}_{i.} - \mu_i}{\sigma/\sqrt{n}}$$

$$= \frac{\bar{y}_{i.} - \mu_i}{\sqrt{MS_E/n}}$$

We can define a T random variable as Ti for the i'th treatment as y bar i.-mu i/sigma/root n. This nomenclature seems to be a bit funny. The first one is the presence of the subscript i that is not serious we are just showing that as the i'th treatment mean and we always use to look at x bar but that is again also not a problem. We are now having y bar i dot which stands for the treatment mean from the experimentation so that is no problem.

Mu i is the unknown treatment mean on which we are speculating or putting up on hypothesis and the sigma normally we say that we use the T distribution when the population standard deviation sigma is not known which is also true in this case. But as far as the nomenclature goes we are putting it as sigma square and we are using the pooled error variance. In other words, we are using instead of sigma square we are using the mean square error.
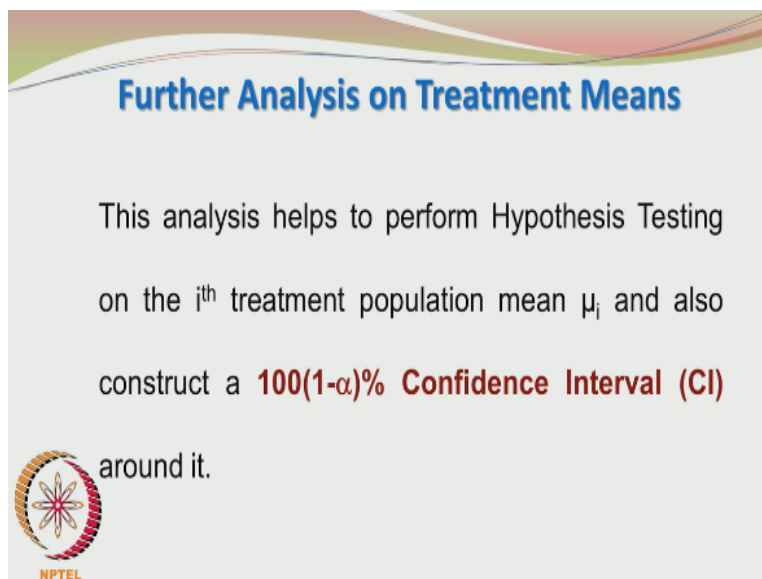
So any how we are not using the sigma as it is because it is not known and usually the repeats are also small in numbers no point in doing something like 30 repeats or 40 repeats that is not practical nor is it wise. So we need to curtail upon the number of repeats and we also have the

classic situation that the sigma is not known and if we can reasonably assume that the parent population is normal.

We can describe the resulting distribution of the same means in terms of a T distribution and the T distribution will involve the parameter sample variance in this case our variance is the mean square error. Again what are the degrees of freedom which we associate with the T distribution these are questions which may create some niggling doubts when we are actually analyzing data. What should be the degrees of freedom I should use?

The answer is very simple as far as this T test is concerned we always associate the degrees of freedom with the degree of freedom with the standard deviation we are using in the T test. In this case we are using the standard deviation as square root of mean square error. And the mean square error is associated with a*n-1 degrees of freedom. Hence the degrees of freedom in the T distribution is also a*n-1.

**(Refer Slide Time: 15:47)**



Now if you can do a T test based on a set of data you can also construct the confidence interval. So let us see how we can construct a 100*1-alpha% confidence interval around the treatment mean mu i. Please note that the treatment mean mu i is not known all we have a set of data corresponding to the i'th treatment we have the set of repeated measurements and we are going to construct a confidence interval using that.

Very simple the confidence interval, the 2 sided confidence interval is bounding the value of the treatment mean and that is given by y bar i. – t alpha/2 a*n-1 square root of mean square error/n <= mu i <= y bar i. + t alpha/2*a*n-1, a correction please t alpha/2, a*n-1 square root of MSe/n. This is the average value for the i'th level of the factor.

The i'th treatment mean from the experimental data t alpha/2, a*n-1 you may recall as the percentage point corresponding to the upper tail of the t distribution a*n-1 is the degree of treatment associated with the mean square error which is here. n is the sample size or the number of experiments we have repeated we are assuming that the treatments have the same number of repeats and this is the population treatment mean which is not known.

So, whatever we studied a few classes back finds immediate application and we are also now in a better position to understand what this confidence interval really means.

**(Refer Slide Time: 18:06)**

**Difference Between Treatment Means**

A T-test may also be performed on difference in individual means with the following hypotheses

$$H_0: \mu_i - \mu_j = 0$$

$$H_1: \mu_i - \mu_j \neq 0$$

So if you can do a T-test on a single treatment mean why can we do it for difference between the i'th treatment mean and the j'th treatment mean? So, we construct the null hypothesis as mu i – mu j = 0 and H1 is mu i – mu j != 0. So can you put these 2 statements in your own words? If you can pause a bit and see what should be the problem statement corresponding to these 2 hypothesis.

So, let us move along. H not mu y- mu j = 0 also implies mu i = mu j which means that there is absolutely no difference between the i'th and the j'th treatments or treatment mean. Then we look at the alternate hypothesis H1 mu i- mu j != 0. So, there is a difference between the 2 treatments. The difference may be negative or positive but the differences are difference and there is significant deviation.

In fact, between the i'th treatment and the j'th treatment which means that either the i'th treatment is better or j'th treatment is better depending up on what your output performance criteria is?

**(Refer Slide Time: 19:46)**

**T-Test for Difference Between Treatment Means**

$$T_0 = \frac{(\bar{y}_{i.} - \bar{y}_{j.}) - 0}{\sqrt{\dfrac{\sigma^2}{n} + \dfrac{\sigma^2}{n}}} = \frac{(\bar{y}_{i.} - \bar{y}_{j.})}{\sqrt{\dfrac{2MS_E}{n}}}$$
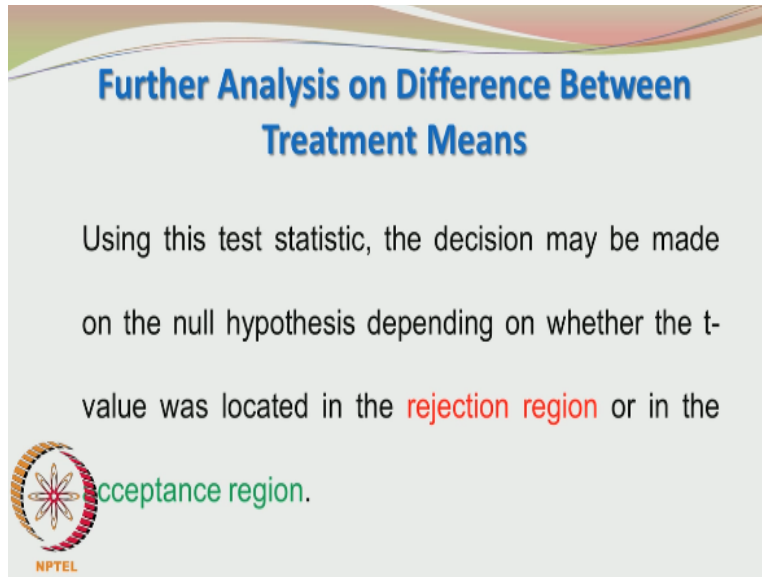
*Assuming equal number of repeats 'n' in each treatment*

Now, if we want to construct the T-test based on this we have the T statistics as y bar i dot –y bar j dot - 0 divided by the standard deviation and that would be sigma square by n + sigma square/n. We have used the same number of repeats and we also know that the variance of the sample mean would be sigma square/n. We have seen this derivation also previously and we know that variance of the difference between 2 random sample means.

Would be sigma square/n + sigma square by n irrespective of whether you put + or – the variance of a linear combination of the 2 sample means either x1 bar + x2 bar are x1 bar – x2 bar would be sigma square/n + sigma square/1. We are assuming in this case a couple of important things that the 2 random sample means are from identical distributions that is why they are both having the same variances.

Another important thing is that the 2 random samples are independent of each other. So the covariant between the 2 would vanish and hence you are having sigma square/n + sigma square/n. So, we have again seen these concepts previously. I request you to please revise these concepts at this particular stage before we move on. So, we have 2 sigma square/n which is written as 2*mean square error/n.

I like these tests, these tests are very elegant and once we know their origin they make a lot of sense.
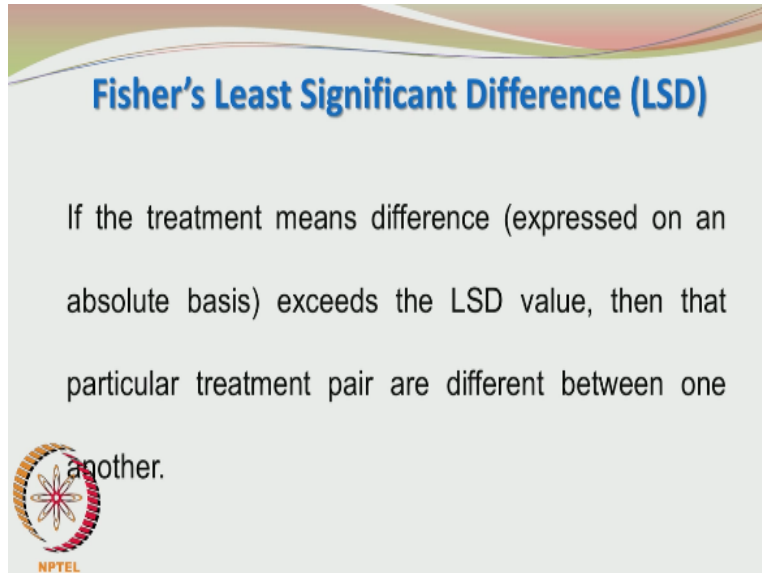
**(Refer Slide Time: 22:13)**



Now using this test statistic, the decision may be made on the null hypothesis depending on whether the t-value was located in the rejection region or in the acceptance region. Again not much more has to be said we will do a few problems to demonstrate this.
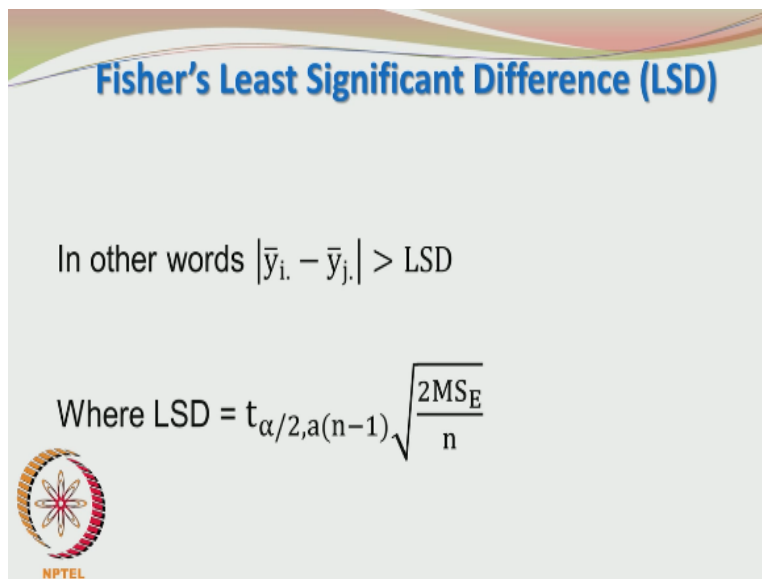
**(Refer Slide Time: 22:36)**



Next, we come to the Fisher's least significant difference. The analysis of variance tests only told that okay on of the treatment at least is having an effect so we cannot reject or we have to reject the null hypothesis. But it does not really tell you which of the treatment means is different and if among 10 treatments for example 3 of them are significant and others are pretty much comparable which of the 3 are significant?

So the analysis of variance test does not really tell us that so for this purpose we can use several methods but one of them which is commonly encountered is the Fisher's least significant difference. The short form for that would be LSD. So let us take the treatment mean values and find the difference between the 2 and express them on an absolute bases and if this difference exceeds the LSD value then that particular treatment pair are different between one another.

So, we have to resort to the experimental data. We compare the treatment means using the experimental data find the difference take the absolute value of the difference, calculate the LSD value see if the difference between the treatment mean values are higher or lower then this LSD value. If it is higher than the LSD value, then the 2 treatments are significantly different from one another.

So very interesting, the next question which follows naturally is how do you define this LSD value?

**(Refer Slide Time: 24:35)**



Fisher's Least Significant Difference (LSD)

In other words $\left| \bar{y}_{i.} - \bar{y}_{j.} \right| > LSD$

Where LSD $= t_{\alpha/2, a(n-1)} \sqrt{\dfrac{2MS_E}{n}}$

And LSD is given by t alpha/2, a*n-1 and the square root of 2MSE/n. LSD= t alpha/2, a*n-1 square root of 2MSE/n.

**(Refer Slide Time: 24:54)**

**Fixed Effects Analysis**

❖ The present analysis is termed as the fixed effects approach. Here, the factors, for e.g. the type of fertilizer, takes discrete specific values only.

❖ There is no point in interpolating or extrapolating the responses to other treatments that were not included in the analysis.

So we have to understand one important limitation if you want to put that pay on this method of analysis. This is called as the fixed effects analysis. We may be doing experiments with discrete entities like fertilizer you have fertilizer A, fertilizer B, fertilizer C. As far as the former is concerned there are 3 different brands of fertilizers. There is no fertilizer intermediate to A and B.
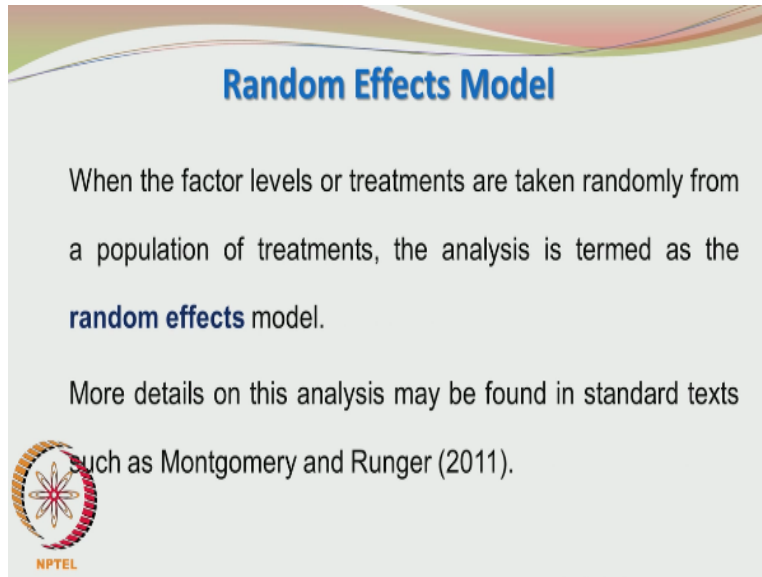
Well an engineer or a scientist may argue I will blend the 2 fertilizers in whatever proportion I want and then I can create a continuous distribution of fertilizers but in may practical physical situations it may not be possible even though the fertilizers can be blended to give any blend of fertilizer. But if you are only going to buy the fertilizers of the shelf you are not going to get these blends.

You will probably get only type A, type B, type C and so on. And another example for chemical engineers would be a type of catalyst you can have catalyst A, catalyst B, catalyst C so you do not have anything in between. Sometimes mechanical engineers may be working with machine A, machine B, machine C. Obviously, you cannot have anything in between you can neither interpolate nor extrapolate.

You will only have to see the effect of these 3 machines so those effects are fixed. I can give plenty more examples but I think you got the idea so we have these fixed effects analysis and

that corresponds to the discussion we were having so far. So the main point is we cannot have any interpolation or any extrapolation possible because our effects are factor levels or treatments are fixed to whatever we have considered.

**(Refer Slide Time: 27:11)**
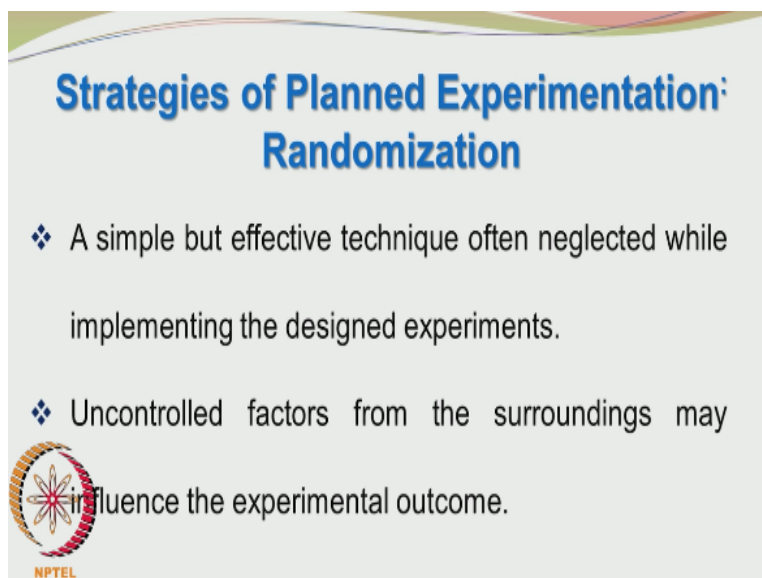


**Random Effects Model**

When the factor levels or treatments are taken randomly from a population of treatments, the analysis is termed as the **random effects** model.

More details on this analysis may be found in standard texts such as Montgomery and Runger (2011).

There is also another model called as the random effects model whereas the limitations are no longer present and more detailed discussion can be found from Montgomery and Runger's book for this. Discussion is quite easy based on whatever we have learnt you should be able to understand that. Let us now move on the basic idea behind this discussion is to provide the terminology and the concepts.

**(Refer Slide Time: 27:48)**



**Strategies of Planned Experimentation: Randomization**

❖ A simple but effective technique often neglected while implementing the designed experiments.

❖ Uncontrolled factors from the surroundings may influence the experimental outcome.

Now, how do you do the experimentation? I am going to tell you 2 interesting ways to carry them out and these are not usually done the main limitation for people who do experiments is the feeling, I would not put it as fear. The feeling of uncertainty and feeling of lack of time so they forget to take some simple precautions when doing the experiments. However, if you can practise yourself to carry out the experiments in the ways I am going to recommend.

You will find it lot more useful and effective. The initial effort may be slightly more but it will definity be helpful to you there are 2 aspects one is the randomization and the other is the blocking. Randomization makes sure or helps you to make sure that the effects not accounted for (()) (28:57) experiments are truly random. There are no systematic factors influencing your experiment.

So, when you do randomization you are doing the experiments in such a way one particular treatment is not selectively affected systematically by the uncontrolled factor. All the treatments are affected by the uncontrolled factors so the effect of the uncontrolled factors spread across all the levels of your experimentation and hence you do not have any systematic factors not accounted for influencing your experiment.

In simple terms whenever you do the sequence of runs you do it in a random fashion. Do not say I will do on Monday the lowest level, on Tuesday I will do the medium level and on Wednesday I will do the highest level. So, please do not carry out your experiments you do them in a randomized fashion.

**(Refer Slide Time: 30:10)**

**Randomization**

❖ The errors are assumed to be **random,** independent, normally distributed and having the same variance.

❖ By **randomization,** we ensure that the effects of unaccounted factors are randomly and uniformly dispersed among the different experiments.

So by randomization we ensure that the effects of unaccounted factors are randomly, and uniformly dispersed among the different experiments.

**(Refer Slide Time: 30:22)**



**Randomization**

This way, no experiment or subset of experiments are alone subjected to these unaccounted extraneous influences.

In this way you are minimizing the effects of unaccounted extraneous influences.

**(Refer Slide Time: 30:31)**

## Strategies of Planned Experimentation: Randomization

Randomization is implemented by running the designed experiments in a random fashion and the allocation of the experimental material to the different runs is also done in a random fashion

NPTEL

So, randomization is implemented by running the designed experiments in a random fashion and the allocation of the experimental material to the different runs is also done in a random fashion. I found one interesting example in one of the books in which they were trying out different brands of petrol on motors and the allocation of the petrol to the motors was run in a randomized fashion.

So, if there are 15 motors and 3 brands of petrol you want to test you do not give petrol 1 to the first five motors which are located in one building and then you give petrol to the another set of motors located in another company. So, you do not do it that way. You assign the petrol brand randomly to all the 15 motors. So, we will be looking at some more examples.

**(Refer Slide Time: 31:54)**

So, the effects of uncontrolled variables will be spread across all the experiments in a random manner that by preventing any systematic, additive or negative effect on the responses.

**(Refer Slide Time: 32:18)**



How do you do the randomization? You can if you are knowledgeable about spread sheeting you can put the sequence of you runs in the standard order first in one column, the second column you can generate random numbers corresponding to each entry and what will happen is you will have a random number tagged down to each of your experimental run then you can sort in the ascending order of these random numbers.

Or even the descending order of these random numbers and all the values would be aligned and then you will find that that runs are randomized. This is one technique of doing it. So, let us now move on to another important concept in experimentation which very rarely is implemented. Blocking is pretty much none existence in most experimental analysis but this is again a very important concept and it increases the sensitivity of your experimentation.

So that is something very good if you want your experiments has to be more sensitive so that you can capture the differences between the different treatments that is very good. So, let us see how to do blocking. Well, you carry out your different treatments on us specimen and then of course we want to do repeats. So, rather than doing the repeats on the same specimen we may do repeats on different specimens.

So, on each specimen we are carrying out different treatments and on different specimens we are carrying out the repeats of our experiments. So, the natural question is the specimens we are using for each repeat are they identical? Usually these specimens will have some minor difference between one another and so we cannot consider all these specimens to be exactly identical in all aspects.

There is an element of uncertainty by the difference in the specimen we are using in our experiment. So, we can use locking to take this factor into account. So, the repeats on experiments may not be carried out on a single specimen but on different ones. Each specimen is subject to different treatments and these treatments are repeated on different specimens. So, please understand this each specimen is subject to different treatments.

And these treatments are repeated on different specimens.
**(Refer Slide Time: 35:20)**

**Concept of Blocking**

If we choose to ignore the difference between the specimens it is as if we are repeating experiments with identical specimens using different treatments, but this is rarely possible.

The question is how to account for the difference for the difference between the specimens on which the treatments were carried out? Can we block out the effect of the types of specimen on which the tests were carried out. So, understand these statements before we move further. What will happen if you ignore the difference between the specimens then you are assuming that these specimens are identical in all aspects and that is usually not correct.

One simple example I can give is you are having 5 different sections in a school and 3 teachers are teaching in each section so if you want to compare between the 3 teachers then you not only go by the student performance in the first class but you will look at the performances in all of the five sections but each section may be different in the type of students they have. For example, section 2 may have unusually bright set of students.

Whereas section 5 maybe having an unusually large set of mischievous students who are not interested in studies yet but they are likely do well later on anyway. So, in this case we are having blocks as the sections the treatments are of course the teachers or the teaching methodology which is being adopted and the blocks would be the different sections. In order to identify the difference between the 3 teacher.

Or the teaching methodologies you have to take into consideration the difference between the sections in which they are teaching. If you do that you are making a test more sensitive. On the

other hand, if you do not account for this you are implicitly assuming that the sections are identical in all aspects which is not true. So, you may have to take blocking into account.
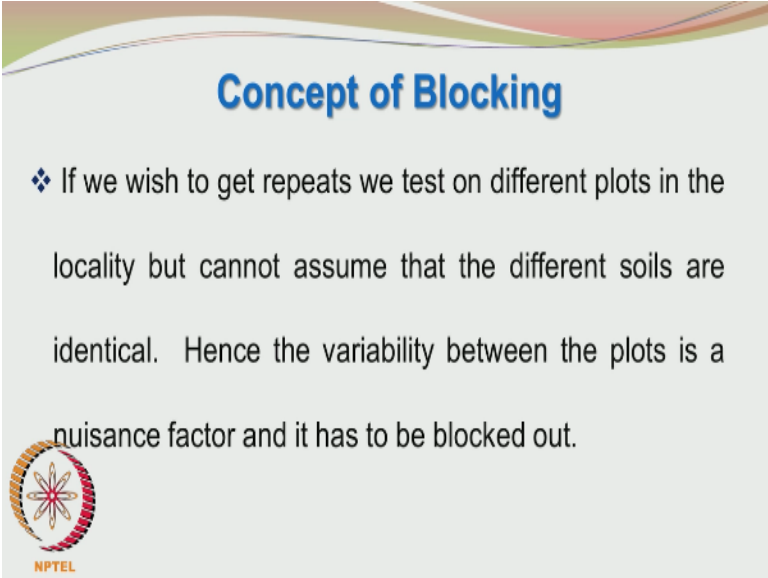
**(Refer Slide Time: 37:56)**



Let us look at the fertilizer example I seems to partial to this example involving fertilizers any way let us carry on with it. Let us try different brand of fertilizers to see which gives better crop yield. So the farmer would be interested in applying different brands of these fertilizers on a given plot of soil and then observe the yield.
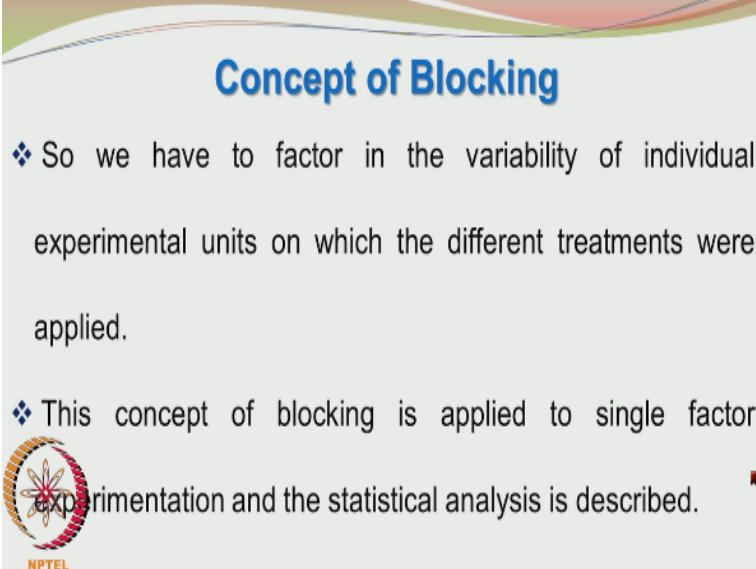
**(Refer Slide Time: 38:35)**



So by just testing on 1 plot we cannot really generalize the conclusion. So, we have to repeat the test on different plots in the locality. But the different plots in the locality may not be identical.

One plot may be close to the river and may be quite fertile where as another block may be quite far way where the watering is not as good or the water level are not as good so it may not be identical with the plot close to the river.

So we have to account for the different plots variability. It is actually a nuisance factor and we have to block out the effect of the type of plot. We are not really worried about the type of land. We want to see the most suitable fertilizer for our crop.

**(Refer Slide Time: 39:48)**



So please note that we are doing our experiments involving a single factor. In the fertilizer experiment again we have the different brands of fertilizers as the different levels of this factor. And what is this plot of lands? The plots of land seem to be looking like an additional variable. In fact, it is not so. It is a nuisance factor not a controllable factor and we are locking it out. So, it looks like a 2 variable experiment but in fact it is not so.

**(Refer Slide Time: 40:35)**

## Concept of Randomization within blocks

Now that we have demarcated the different experimental units into blocks, the treatments within each block is carried out in a random fashion.

Once you have identified the different blocks. Within each block you can carry out the experiments in a random fashion. So it does not mean now that you have divided your experimentation into different blocks you can forget about randomization. Whether you are doig blocking or not you have do randomization. So, even after dividing your experimental programs into different blocks within each block randomization should be implemented.

**(Refer Slide Time: 41:16)**

## Concept of Randomization within blocks

Why should the different treatments be carried out in a random fashion?
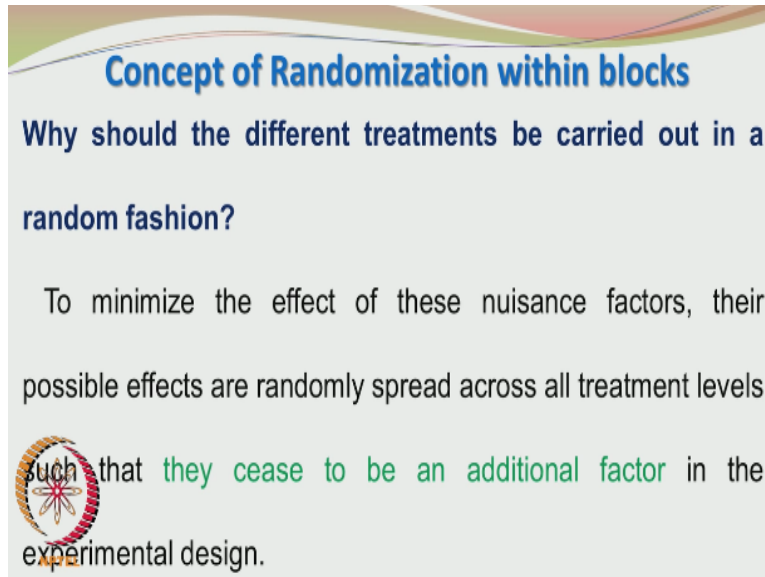
If we carry out the treatments in each block in a definite sequence, then the nuisance factors may interfere with the same treatment levels in each block that is being tested.

So, why should the different treatment be carried out within each block in a random fashion? If we carry out the treatments in each block in a definite sequence when the nuisance factor may interfere with the same treatment levels in each block that is being tested. So, I will just give a brief example to drive home this point.
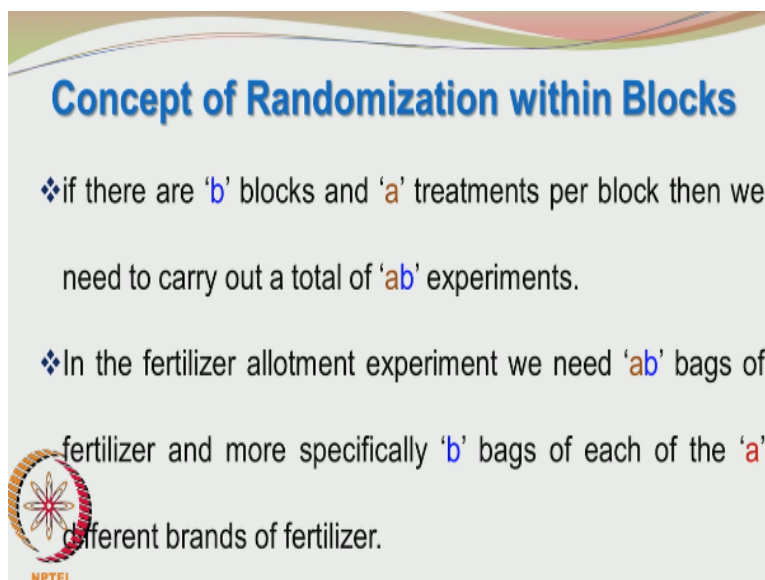
Why should the different treatments be carried out in a random fashion? We do that in order to minimize the effect of the nuisance factors. The effects of the nuisance factors are randomly spread across all treatment levels such that the nuisance factor no longer are acting as additional factor in the experimental design. So, we are making sure that there is no hidden factor in our experimental design.

Whenever we buy some appliance we have to be very careful that there are no hidden charges. Whenever we do experiments, we also have to be careful that there no hidden factors.
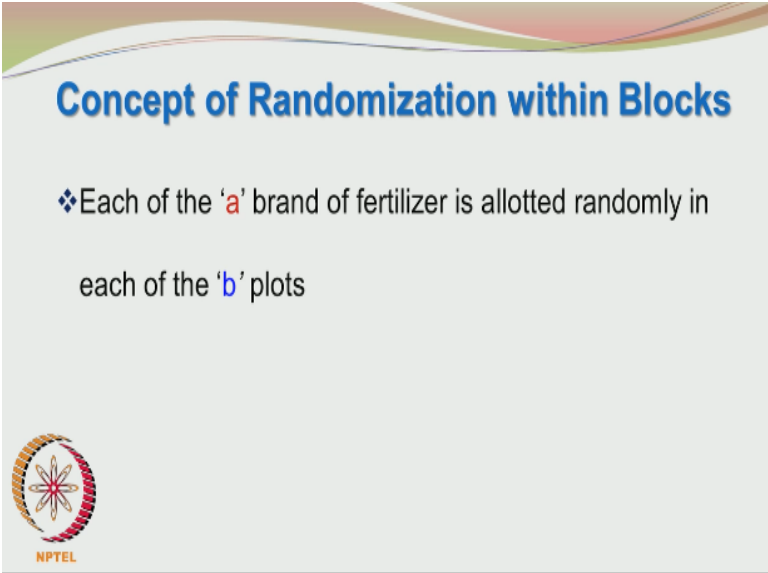
**(Refer Slide Time: 42:46)**

So, let us do the analysis involving blocks. Let us say that there are 'b' blocks and 'a' treatments per block then we need to carry out the total of 'ab' experiments. Now, we are talking about blocks and we are no longer referring to the term repeats. So, we have to talk and analyze the effects of blocks in our experiments. In the fertilizer allotment experiment we have 'a' fertilizers and 'b' blocks.

So that we need ab bags of fertilizers or specifically 'b' bags of each of the 'a' different brands of fertilizer. So, no big problem if you are having 3 plots or rather 5 plots and 3 different brands of fertilizers you should have a total of 15 bags of fertilizers and each brand you should have 5 bags.
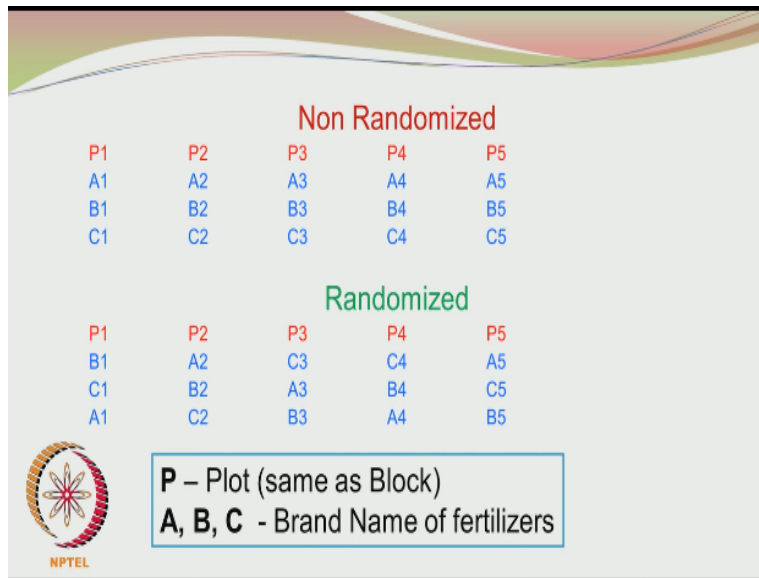
**(Refer Slide Time: 43:56)**



Now, each of the 'a' brands of fertilizer is alloted randomly in each of the 'b' plots. What does this means? Looking too many 'a' and 'b' can become a bit confusing what I am simply trying to say is within each plot of land you allocate the different brand of fertilizers randomly. So, I will just illustrate it in a moment.

**(Refer Slide Time: 44:33)**

**Non Randomized**

| P1 | P2 | P3 | P4 | P5 |
|----|----|----|----|----|
| A1 | A2 | A3 | A4 | A5 |
| B1 | B2 | B3 | B4 | B5 |
| C1 | C2 | C3 | C4 | C5 |

**Randomized**

| P1 | P2 | P3 | P4 | P5 |
|----|----|----|----|----|
| B1 | A2 | C3 | C4 | A5 |
| C1 | B2 | A3 | B4 | C5 |
| A1 | C2 | B3 | A4 | B5 |

**P** – Plot (same as Block)
**A, B, C** - Brand Name of fertilizers

So, let us look at randomized design and the non-randomized design. So you are having 5 plots of lands. Let us call them as P1, P2, P3, P4, P5 and of course we are not putting the blocks the blocks are P1, P2 so on to P5. They are not randomized, please note that. Within each plot you can randomize the allocation of the fertilizer. Let us say that this plot P is divided into 3 portions and this portion is given A1 then B1 then C1 in the next portion again.

It is again A2, B2, C2 so you can see that the fertilizer A is assigned to one end of the plot. Fertilizer B is always assigned to the middle portion of the plot. Fertilizer C is always assigned to the other extreme of the plot. Well, you can let your imagination run wild and think that there is a river in this region and so because of the proximity of the river the rich alluvial soil, what not? May be A1 fertilizer will get a boost because of these favorable conditions.
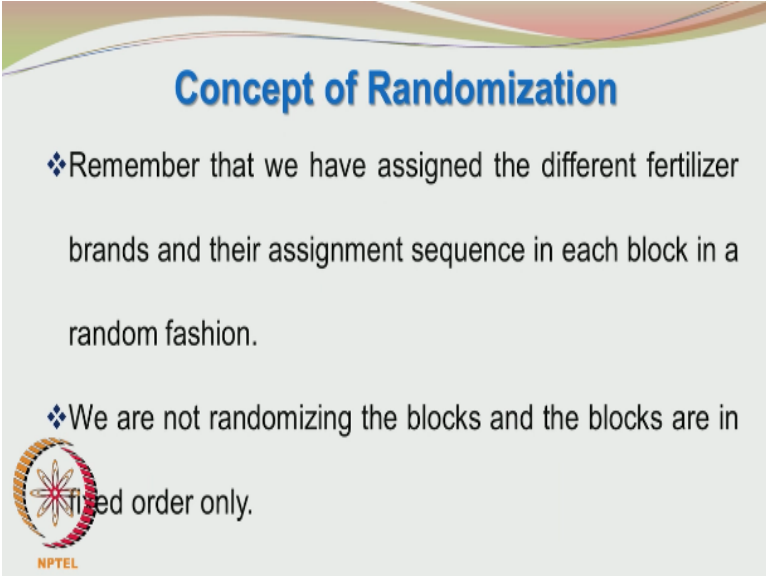
And seem to perform very well. Whereas fertilizer C on the other hand which is furthest from the river may be closer to barren and rocky regions and so it may appear as if fertilizer C is performing poorly. So, the proximity to the river or the distance from the river is a nuisance factor. So to mitigate that effect you are assigning these fertilizer bags randomly in each plot. Well that is not the only affect, well we are trying to get rid of that.

By randomizing the distribution or allocation of the fertilizers. The next source of variability is the variability caused by the different plots. So you have plots one P1, P2, P3, P4, P5 may be plot

1 you are having more man power plot 2 you are having less man power. So, these plots themselves are different so in order to account for the effect of the variability in plots you are calling them as block.

So, P1 is block 1, P2 is block 2 so on to P5 is block 5. So, P is standing for plot and is as representing a block. A, B and C are the different brand names of fertilizers. You can see that in the randomized design you do not have any systematic organization of the fertilizers you have B1, C1, A1, A2, B2, C2, C3, A3, B3 and so on. So, it is pretty random.

**(Refer Slide Time: 48:15)**



## Concept of Randomization

❖Remember that we have assigned the different fertilizer brands and their assignment sequence in each block in a random fashion.

❖We are not randomizing the blocks and the blocks are in fixed order only.

NPTEL

So, please note that we are not randomizing the blocks and the blocks are in fixed order. Sometime it may not be possible to randomize the blocks. How can you move one plot of land to another portion and so on any way let us not get too carried away by all this. We have to realize some limitations also.

**(Refer Slide Time: 48:37)**

**Concept of Randomization**

There may be differences between the blocks and differences within each of the block where different treatments are being applied.

So there may be differences between the blocks and differences within each of the block where different treatments are being applied.

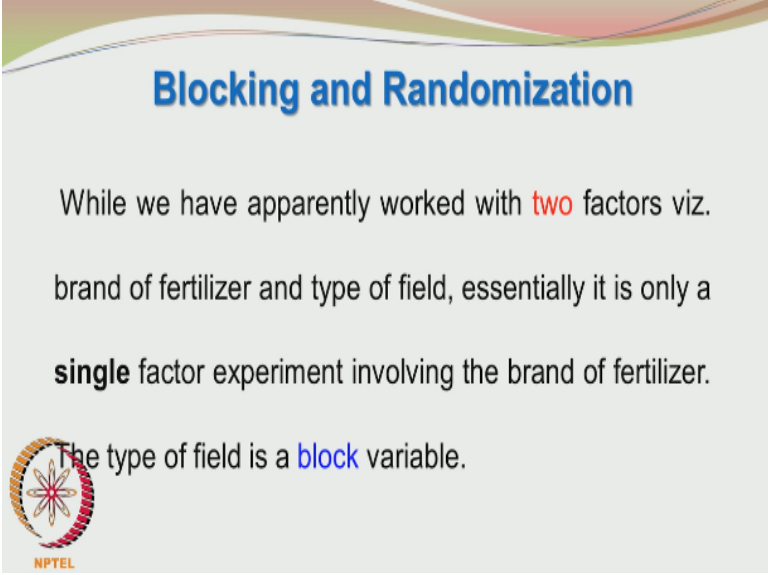**(Refer Slide Time: 49:01)**



**Concept of Randomization**

However, differences within the blocks may only be attributed to differences between the three brands of fertilizer as the experiments within the block were carried out in a random fashion.

However, differences within the blocks may only be attributed to differences between the 3 brands of fertilizer as the experiments within the block are being carried out in a random fashion. So, across the blocks there may be the effects of blocks also coming into the calculations which we will be soon resolving but within each block the output is only governed by the type of fertilize.

Any nuisance factor which we are not able to account for is dampened out or filtered out by the concept of randomization.
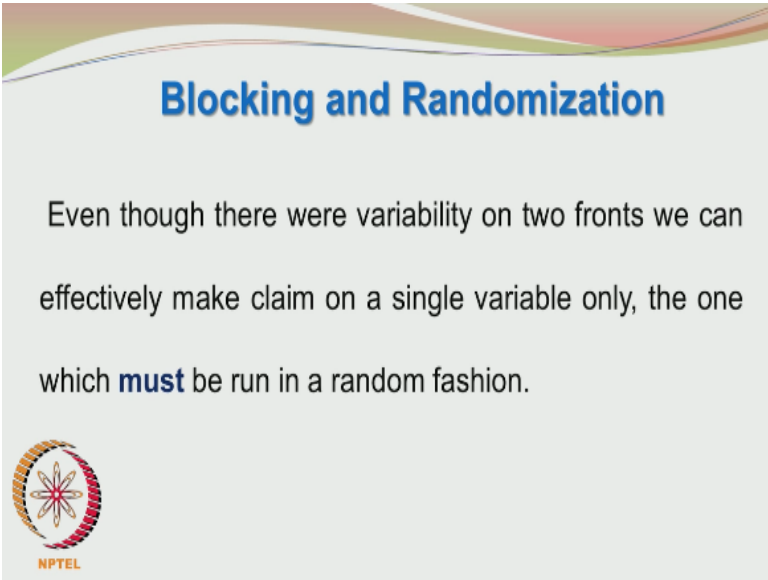
**(Refer Slide Time: 49:54)**



So, it looks as if we are working with 2 factors. The first factor is the brand of fertilizer and the second factor is the type of plot or type of field. However, please remember that we are only working with the single factor we are analyzing a single factor experiment, a fixed effect experiment at that and the additional variability cost by the plots are analyzed in terms of blocks. So, we are having a block variable and then we are also having a treatment or a factor.

**(Refer Slide Time: 50:41)**

So, there are variability on 2 fronts we are effectively making a claim on a single variable only. We are not comparing different plots of lands. We are only comparing different brands of fertilizers. So, our conclusion and hypothesis is based on a single factor that is the type of fertilizer. It so happens that the means that we are implementing on is variable and we have to account for that variability.

So you cannot consider it as an additional factor. So, we have covered a lot of ground so we will take a break now and it is important that we refresh the concepts that we have seen so far. We will next move on to mathematical analysis of these concepts. Thank you.