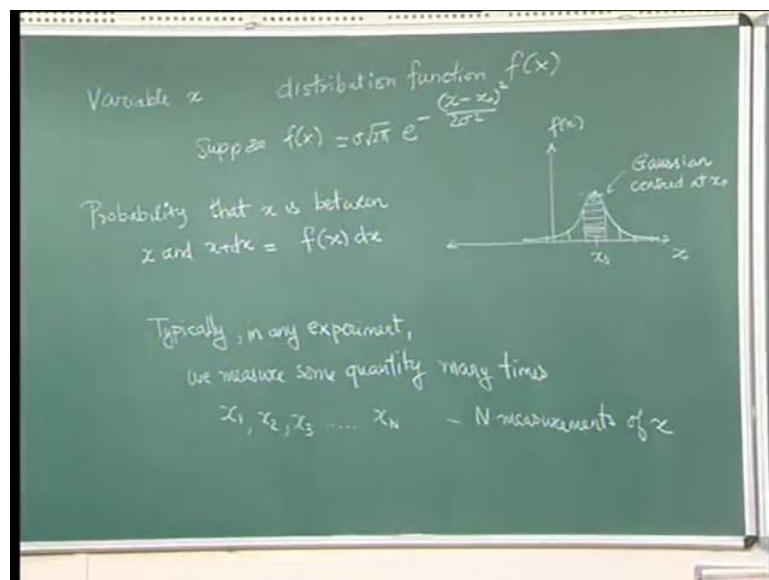**Mathematics for Chemistry**
**Prof. Dr. M. R. Ranganathan.**
**Department of Chemistry**
**Indian Institute of Technology, Kanpur**

**Lecture - 40**

In today's lecture, which will be the last lecture of this course, I am going to talk about a slightly different topic, and this has to do with errors and estimation of parameters. Now, in order to motivate this, let me remind you of an object that you have seen before that is called the distribution function.

(Refer Slide Time: 00:43)



So, suppose you have a variable x and let us say this variable has a distribution function f of x. If the distribution function of this variable is f of x, then you know that you know that this f of x will be a function for example, suppose f of x is equal to e to the minus x minus x 0 square by 2 sigma square into sigma root 2 pi. So, suppose this is a Gaussian distribution and suppose, so then what you know that f of x has this form where x versus f of x looks like this and if this is x 0 then it looks like a Gaussian function centered at x 0.

So, it looks like a Gaussian centered at x 0. Standard deviation is sigma. So, then what do we mean by distribution function? What we mean is the following, so probability that x is between x and x plus dx is equal to f of x dx. So, the probability of finding x in an interval is f of x times dx. So, then f of x is the probability density or the probability

distribution function. So, now what do we mean by this distribution function? What the way to think about this is the following, that suppose you have this suppose you have this variable x and you had a way to measure it.

So, you measure it once you will get some value of x, you measure it again you will get another value you keep repeating this measurement many times. So, if you repeat the measurement many times then most of the times you will get values of x in this range, in the range whether distribution function is high or more. More likely to get values of x in this range and or less likely to get values where the distribution, where the value of the distribution function is very low.
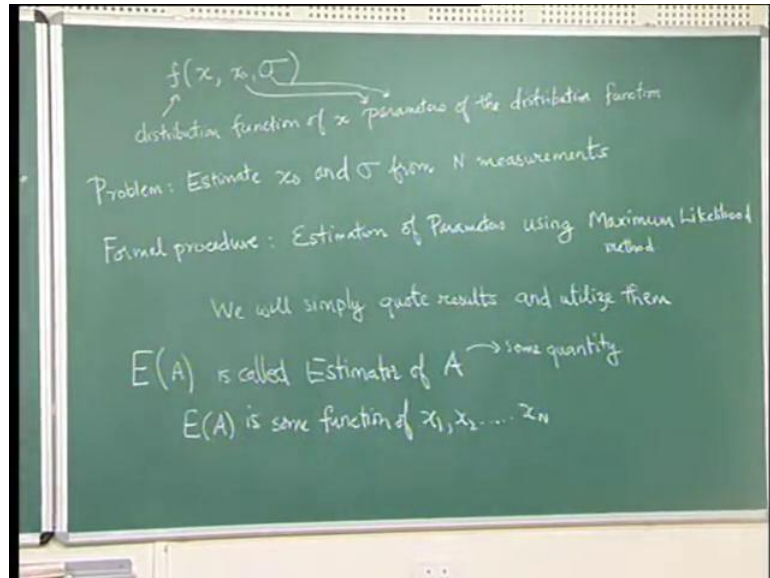
So, if you break this space into intervals like this, then most of the times you will be in this interval. So, most of the times you will be in this interval, most of the time you measure x you will get somewhere around here, then lesser number times you will get, in this lesser number of times you will get in this and even fewer number of times in these. Now, this is what will happen if you make a very large number of measurements. So, if you make an, how large should the number of measurement be, the number of measurements should actually be infinity.

If you make infinitely make many measurements then the number of times you get x in any interval will be proportional to f of x times the size of the interval or actually integral f of x dx so that in interval. Now, typically, so in any measurement, in any experiment we measure some quantity many times. So, now I want you to imagine that you are measuring this quantity x, so you do an experiment and you measure this quantity x. You measure it once, first time you measure it you get a value x1, second time you measure it you get a value x2, third time you measure it you get a value x3 and so on and you do it N times you get these N values.

So, N measurements of x, so this is what you do in any experiment you might be a, if you doing a titration experiment in the lab you try to repeat it many times. So, same way any experiment that you do you try to repeat many times till you get significant number of measurements. Now, the question is what we do with this typically? What we do is we take an average and we take standard deviation and so on. So, now the question is why do we do those things and what is the correct procedure for doing that? So, what we are trying to do is the following, we are trying to, we you have x has a distribution

function but you do not know that the mean is, you do not know what the standard deviation is.

(Refer Slide Time: 07:02)



So, you have, x has a distribution function f of x. Now, x has two parameters it has x 0 and it has sigma. So, this is distribution function of x and these are unknown parameters or these are parameters of the distribution function. So, now your problem is basically you want to estimate these parameters, you want to estimate these parameters from these N measurements. So, the problem, estimate x 0 and sigma from N measurements. So, this in some sense the goal of your experiment is to get the values of x 0 and sigma, so you want to get the values of x 0 and sigma.

Now, we have sort of turned around the problem instead of saying that we want to know the value of x, what we say is we want to know the mean value of x and you want to know what is the standard deviation of x. So, if x, if you say that if x is distributed, if you say that the distribution function of x has this form where it has, where it spiked at x 0 and has a standard deviation sigma, then you want to estimate x 0 and sigma based on these measurements. So, this is the, this in some sense is the crocks of what we are doing in all our measurements, what we are doing in any measurement is to estimate the parameters.
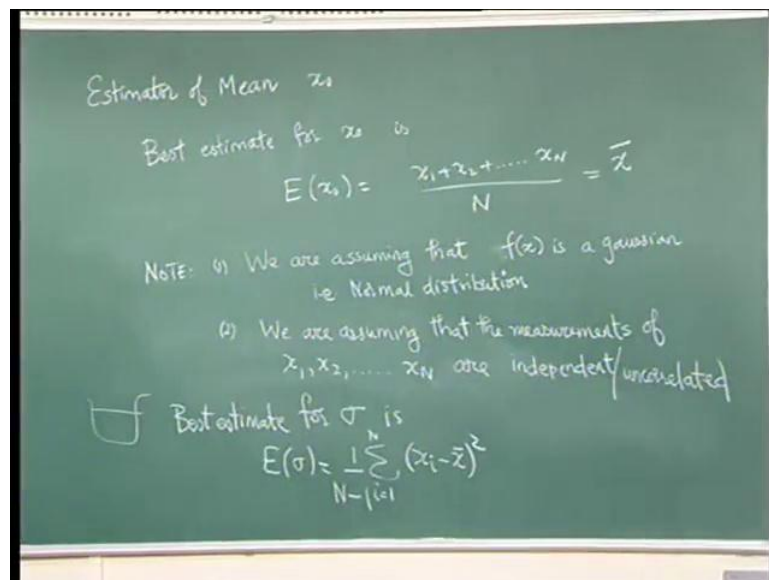
The parameters that we are estimating are related to the distribution function. Now, you might be, you might say that you want to know what is the mean of the distribution

function, you are only interested in, sometimes you are only interested in the mean of distribution function. You want to know what is the most likely value of that function which is the also the same as the mean in this case. So, how do we go about estimating x, x 0 and sigma from these measurements? Now, the formal procedure is here, is an area of work called estimation of parameters, it is using maximum likelihood methods.

Now, we will not do the formal procedure. So, this is not something that we will be doing, we will simply use the results, quote results and use them. So, now let us see what are the results that we are interested in? So, the first thing is to, is we have an idea that let us say E of A is called estimator of A. And E of A is some function of all these variables. So, E of A is some function of x 1, x 2 up to x N.

So, suppose you have a quantity A. So, A is some quantity, some quantity when I say some quantity, it might mean he, it might be average, it might be standard deviation etc. And E of A is some function of E of A is called the estimator of A. The estimator of A is some function of these variables, of these measured quantities. So, let us go ahead and do some estimates. First one, first estimate will be doing as a mean.

(Refer Slide Time: 12:39)



So, the first estimate will be doing is to give the estimator of x 0. So, we want to estimate the value of x 0 from these measurements, from all these measurements. The estimator, the if you follow the procedure using the maximum likelihood method which I have not done, but using a formal procedure you can show that the best estimator for x

0 is E of x 0 is equal to x 1 plus x 2 x N by N. So, the most likely value of x 0 is just the mean of these measurements. So, and this is something that makes that is very obvious that suppose I give you a bunch of measurements.

Then if I ask you what is the most likely value of this variable you will always say it is the mean of those variables, of those measurements. So, this is the most likely estimator of these, most likely estimator of the mean. So, note I want you to note couple of things, first one, we are assuming that f of x is a Gaussian, f of x is a Gaussian. In the sense, we are assuming that in other words normal distribution, this is called the normal distribution. So, we are assuming that the variables are distributed according to a normal distribution.

Second, and this is the more, this is the more important thing to keep in mind whenever we, you are making measurements is that we are assuming that the measurements of x 1, x 2 up to x N are and I will use two words here first is independent and I will use another word uncorrelated. So, they are they basically mean the same thing and both these are very commonly used terminology. So, you are assuming that all these measurements are completely independent that have been the first time you measure x that does not affect what you get the next time. The next time you measure it they have no affect what you get the third time.

So, each of these, each of these measurements are completely independent of each other. Now I want to say that you know this is, you have to be very careful to ensure that all these measurements are independent, you have to do your experiment very carefully. So, this property of these variables is a direct test of how well your experiment is done. And if you do it the right way then you will truly get independent or uncorrelated measurements. Sometimes you might be force to do your experiment in ways where you do not get exactly independent or uncorrelated measurements.
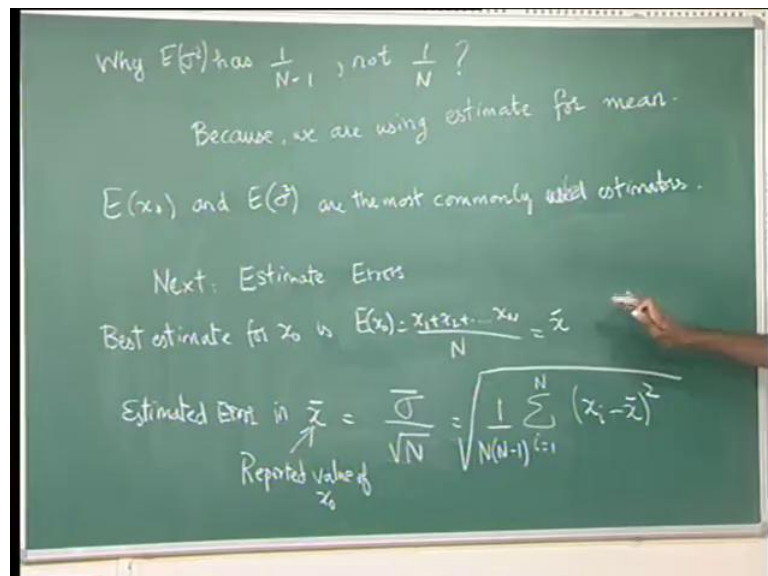
Then in that case you cannot use this, you cannot use this estimator for average. Just to give an example of measurement where in which your various values are not independent. Just imagine that you are doing a titration and let us assume that you have this beaker and in this beaker you measure the volume. So, the first time you measure a volume then you throw this beaker away and you again use it. So, if you throw the sample in the beaker and then you use a beaker again then the next time you measure it

you are likely to have some small amount of that substance, some small amount of the substance then you titrate a volume will turn out to be less.

So, if the flask in which you are carrying out, the carrying out the asset based titration if that is being reused every time then there is slightly to be some amount of the asset or base left in that and that is going to give systematic errors. So, then you, your measurements are not completely independent of each other. There are many such examples you can think of where your measurements might not be completely independent of each other but in such cases you do not use this averaging procedure.

So, you try as far as possible to make sure that your measurements are completely independent of each other and then you can use such an averaging procedure. Next, you want the estimate, the best estimate for sigma is estimator of sigma is equal to sum over i equal to 1 to N x i. I will call this x bar x i minus x bar square. And I will divide this by n minus 1. So, this is something that you have seen before. So, the best estimator for sigma is x i minus x bar square divided by N minus 1. And you should pay attention to that fact that it is N minus 1 and not N. Now, I will just try to give you an idea why it should be N minus 1 and not N.
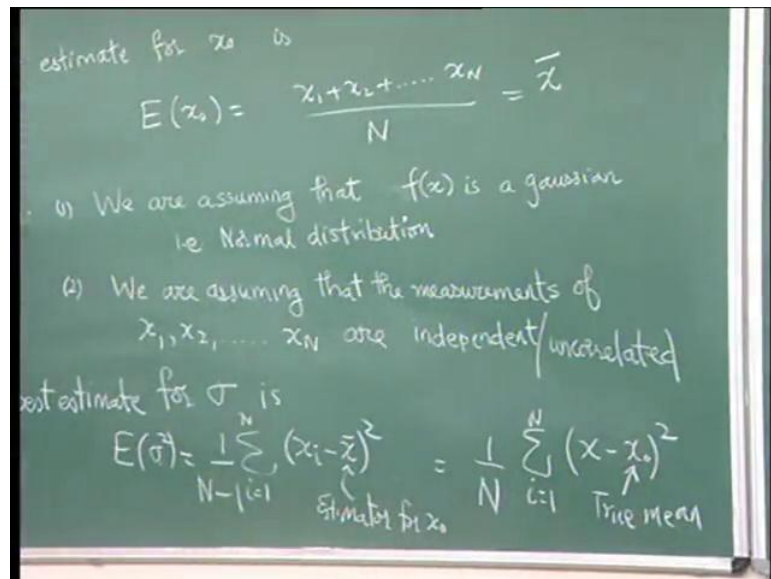
(Refer Slide Time: 20:09)



So, why sigma square, this is the best estimator for sigma square, the variance. Why sigma square has 1 over N minus 1 not 1 over N? Because it looks like each of this is a square deviation, so the average square deviation should seem like 1 over N times this.

But in fact the best estimator for sigma square is not the average square deviation but it is the sum of square deviation divided by N minus 1. So, now actually you can show that this is the best estimator of sigma using the maximum likelihood method which I have not discussed but there is a formal way to derive this.

Now, I am not going to do the derivation but I want to point out one important thing here which will help you understand why you should have N minus 1. That is the following. Now, suppose now in this estimator for sigma square you have an x bar and x bar is actually an estimator x bar itself is an estimator. So, this is an estimator for x 0.

(Refer Slide Time: 21:48)



So, bar is an estimator for x 0. So, now if instead we had used sum over i equal to 1 to N x minus x 0 square. So if you had used sum over i equal to 1 to N x minus x 0 square this is the true this I will call the true mean. So, if there was some way you knew the actual value of the mean, if there was some way you actually knew what the mean was not and the true mean will any general not be equal to x bar. So, if you actually knew the real mean then you would use 1 over N.

But since you estimated the mean based on these measurements, since you already estimated the mean based on these measurements then it turns out that all these N values of these deviations are not completely independent they are only N minus 1, of them that are independent. Because all these, because the sum of all these deviations should

be 0 because that is a definition of xi bar. So, therefore, whenever you use N estimator for x 0 whenever you estimate the average you should use 1 over N minus 1.

If you know, the true mean by some other method if there is some method that helps you actually calculate the true mean without using these measurements then you use 1 over N. So, for example, if you say that you know that the height of some person is 6 feet 1 inch and let us say you know that it is exactly 6 feet 1 inch. Now, you make a measurement, you make 5 or 6 measurements and each time you measure you will never get 6 feet 1 you might be or you will get something spread around 6 feet 1 and you take an average of those measurements.

Then it would not be exactly equal to the 6 feet 1 it might be slightly different depending on the quality of your measurement. But you know that the true mean is 6 feet 1. So, then you can do your standard deviation using this estimate or using this estimate. So, if you know the true mean you use this estimate if you do not know the true mean then you use estimate N minus 1. And this is a very important idea and so, I will just say why not 1 of by N? This is because we are using estimate for mean. So, since you are using an estimate for this mean then this should have 1 over N minus 1.

So, now that you know we have that two most important estimates and these are the most common estimates. So, E of x 0 and E of sigma square are the most common estimates most commonly used estimators. So, most cases you assume that your data has satisfies or normal distribution and you assume that your, each of your measurements is completely independent of each other and you follow this procedure and get estimates for the mean and estimates for sigma square. So, how do you use this?

The way you use this is very simple if somebody ask you what is your best, what is your measured value then you say the best guess for the measure value is this average. Then if somebody ask you what is the variance then you say the best estimate for the variance is this. So, that is how these quantities are used. Now, I mention something about independent and uncorrelated, it is again, it is really a test of how well the experiment is done to make sure that your measurements are independent and uncorrelated. There is there are other cases when you actually have correlated data that you want to measure.

You want to measure the actual correlations, but those are not things that I will be covering. So, next what we want to do is to estimate the errors in this measurements. So,

next estimate errors. So in these two in this x bar and sigma we estimated the parameters of the distribution function. Now, we want to estimate what are the possible errors in these estimates. So, estimate errors in these estimates itself. And this is the next topic that is error estimation. So, I will just give the result for the estimation of the error in these measurements.

So, the statement is so, best estimate for x 0 is e of x 0 is equal to x1 plus x2 x N by N. And the estimated error in this x bar, x bar. So, then suppose somebody ask you the question, what is what you what your estimation of erroneous x bar can be how erroneous x bar can and this is given by sigma divided by square root of N. Is given by sigma divided by square root of N where sigma square I will call this sigma bar square and this is sigma by N this is equal to sum over i equal to 1 to N xi minus x bar square divided by N, N minus 1 the whole thing under root.

So, the estimate for the error, the error estimate in x bar, so this is the error estimate in the reported mean value. So, this is the, this will be the typically the reported value of x 0. So, this is the error estimate in x bar in the reported value of x 0, this is sigma bar divided by square root of N. Now, again this is a very, this is a very tricky issue that if you had noticed is that why do you have the square root of N in the denominator. So, why isn't just sigma a measure of the estimate in x bar of the error in x bar.

The answer is why the reason why it is not just sigma is that we, this is the estimated error in the reported value of the average and suppose I make a single measurement, suppose I made only one measurement. Then the error estimate for that one measurement will be related to sigma. See error estimate for any one of these measurements is related to sigma.
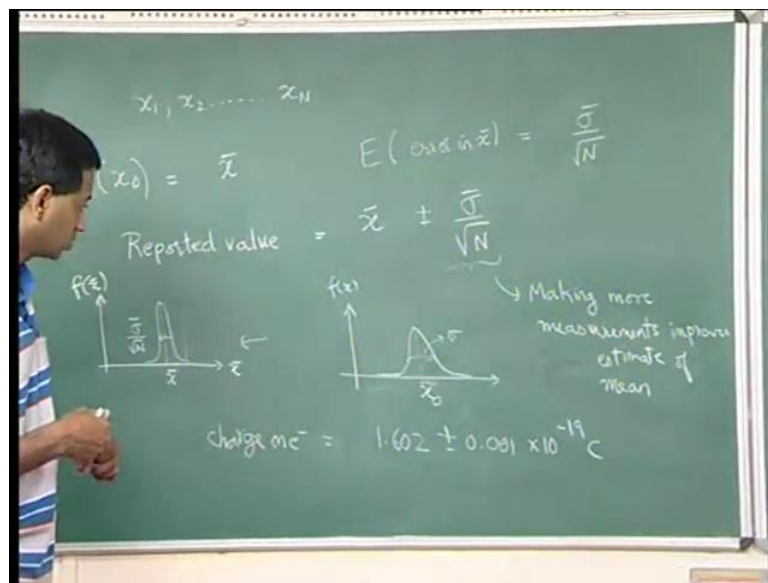
But the error estimate for the average is much smaller than the error estimate for a single measurement. So, what I mean is, suppose you have these N measurements and somebody ask you estimate x 0 and you just choose an estimate which is just one of those measure. If you just say x 1 is an estimate for x 0 that is also an estimate for x 0.

But it is not a very good estimate because the error in x1, the estimated error in that will be of order sigma which is much greater than this estimated error in this. So, this idea of how to get your data and calculate the best estimate for average and calculate what is

the error estimate in that is a very important role in fact, any measurement you take you should always be careful to do this analysis.

So, now next what I want to show is how you take this analysis, how you take various measurements and how you report them. Let us say, if you want to report them in a lab report or in a scientific journal or in a book then how should you, what should your report. So, if you have a bunch of measurements and then how you should you report those the result of those measurements.

(Refer Slide Time: 33:01)



So, suppose you have a bunch of measurements x 1 x 2 up to x N. You want to report your data, then what you say is that your best, you will say that x 0. So, the best estimate for x 0 is x bar and then what you want to say is, so this is the best estimate for x 0 and the best estimate for the error in x bar is equal to error estimate. So, that is sigma bar by square root of N. So, then what you do is you say that your reported value is equal to x bar. So, you say that my value is x bar and then you say plus minus sigma bar by square root of N. So that is how you report the errors.
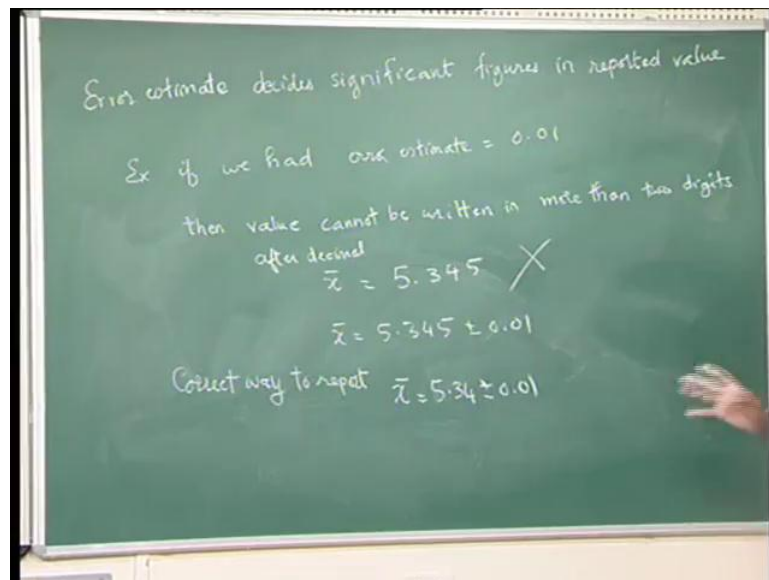
Now, what this means is that the implicit in this is the assumption that your x your f of x is distributed as a Gaussian with mean x bar, and this with mean x 0 and standard and the standard deviation. So, this is of order sigma. So, now what this means is that when you look in this range from plus minus sigma bar by square root of N. Then you can say that a large number of the values will be within this range. So, if you make, so making

more measurements improves error estimate. So, in other sense it decreases the error estimate that is it improves the estimate of mean.

So, as you go as you make more and more so, what we are saying here is that if you what your values of f x bar you are reporting values of x bar and this is your estimate, is based on the average. Now, let us say you make ten measurements you end up here you make eleven measurements you might end up somewhere else you make twelve measurements you might end up somewhere else. This width will be, this width of order sigma bar by square root of N. So, as you make more measurements your x bar will also change and your, but your width will become narrower.

This plus minus, so you can only be sure of x bar up to this interval. So, you might have seen some reported experiments, so reported experiments you might have seen for example, charge on electron is equal to, so 1.602 plus minus 0.001 into 10 power minus 19. So, you might have seen things like this, so you give the, you give the estimate and you see how much the error is?
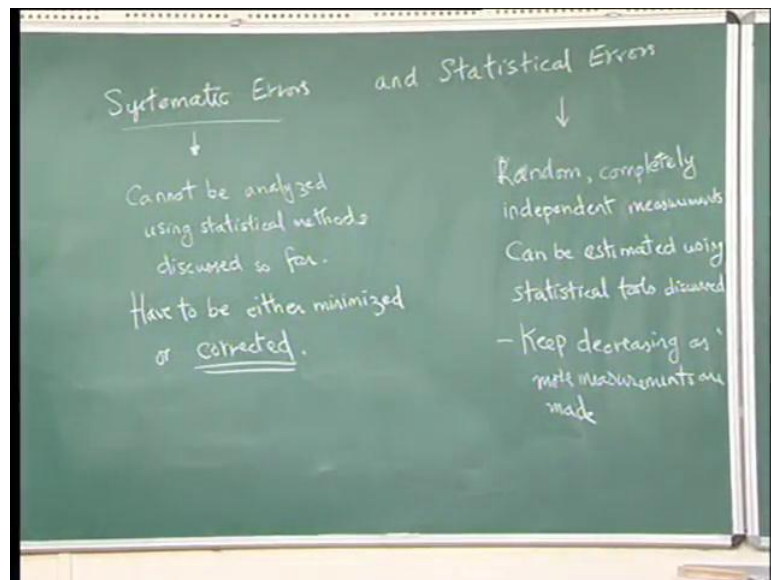
(Refer Slide Time: 38:01)



Now, this error estimate decides significant figures in reported value. So, the error estimate is what decides how many significant figures you will use in the reported values? So, if your for example, if we had error estimate is equal to 0.01 then value cannot be written in more than two digits after decimal. So, you cannot have something like 5.345 so, you cannot have x bar equal to this not possible, so this is not possible.

Because your error estimate is plus minus 0.01, so if you say x bar equal to 5.345 plus minus 0.01, then what it means is you do not know whether it is 5.33 or 5.35? Then it does not matter what you say is, so this 5 is completely immaterial.

So, the correct way to report is report x bar is equal to 5.34 plus minus .01. 5.34 or 5.35 plus minus 01. So, you can never be sure whether it is 5.35 or 5.5, it can be anywhere between 5.33 and 5.35. So, most of the time you measure it you will get between 5.35 and 5.33. So, this is something that you should keep in mind whenever you write your laboratory reports you pay special attention to using the correct number of significant figures in those measurements and also always try to estimate what your error is.

Even if you made a measurement only three times you go ahead and you calculate the error estimate. So, calculate the sigma bar square root of N using the formula of your discussed. So, now the last thing I want to talk about is what are the sources of error in your measurements and in particular. I will discuss the two kinds of errors and it is again extremely important that you keep a track of these in any measurement you make.

(Refer Slide Time: 41:23)



So, the two commonly encountered errors in any measurement are what are called the systematic errors and the second kind what are called the statistical errors. So, the systematic errors are truly errors they are not completely independent they are errors because you are making some genuine mistake which you cannot help making. So, for example, if you are measuring the density of water, now you cannot control let us say

your apparatus does not allow you to control the temperature outside. So, you are measuring the density of water and you cannot control the temperature outside. So, when the temperature outside changes your measurement will give you will get a different measurement. So, then you are generating errors in your experiment because of changes in temperature and the temperature outside. Let us say it changes in some systematic fashion so, then you will get some systematic errors.

So, these are typically cannot be analyzed using statistical methods discussed so far. So, these errors are not you cannot so, whatever systematic errors you get you cannot analyze them using the statistical methods like mean and variance and so on. So, these are these have to be minimized, these have to be either minimized or corrected. So, whenever you know that you are going to get some of these errors then you correct your results for these errors. So, for example, if you let us take the example you have measuring the density of water in a lake now that will change as the temperature outside changes it might change.

So, you correct for the fact that the temperature outside is changing and you corrected in the estimate, and you can use various ways to correct it. But this is completely different from these statistical errors. So, the statistical errors these are random, they are completely independent of so, one measurement is completely independent of another. So, completely independent measurements, if this is can be estimated using statistical tools discussed. So, you can estimate the statistical errors using these tools that we have discussed like the mean and the variance. But all your systematic errors are once that you have to correct for.
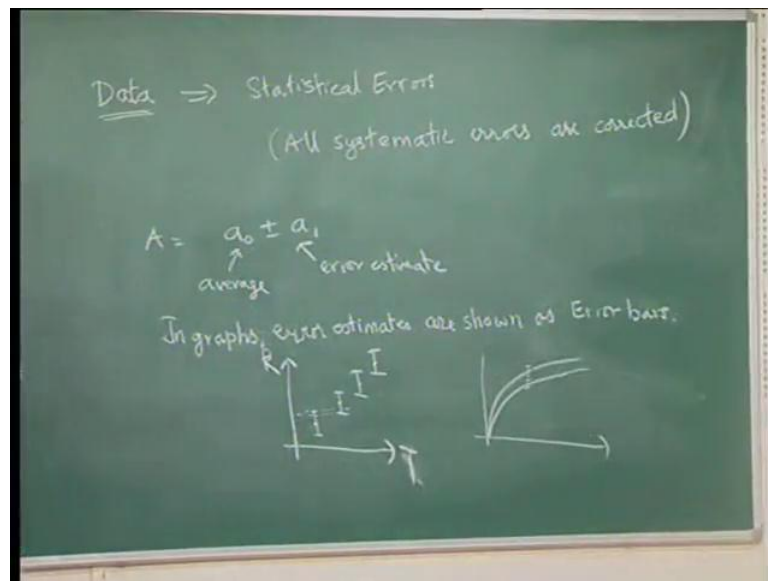
This and its very important to distinguish between these kinds of errors and so whenever you make any measurements you want to either minimize the systematic errors or you want to correct for them. Similarly, when you do statistical errors you want to estimate them and if you make more measurements then the statistical errors keep going down. So, these keep decreasing as more measurements are made. This is typically true, just not, it is not always true.

Typically they keep decreasing as you make more measurements. Whereas, these the typically do not depend on whether you make some or more measurements. So, will keep this in mind that we are always try to minimize the systematic errors and or you

correct for them and whereas, statistical errors are those that are estimated and they can be reduced by making more and more measurements.

So, whenever you do any experiment, you always identify what are the sources of error in your experiment. So, it is a standard exercise to go ahead and identify what are the various sources of error. Classify those errors into systematic and statistical errors and then you make this correction.

(Refer Slide Time: 47:52)



The final thing I want to mention is what is the phrase that is used is data. Data implies statistical errors, so suppose somebody says I have data on this measurement then the normal implication of data is statistical errors not systematic errors. Systematic errors are usually corrected data, so statistical errors and all the systematic errors are corrected. So, once you have corrected your measurements of all systematic errors, what you have left is only statistical errors, and then you can call them data, and then they only have statistical errors.

Though we used most commonly used is the normal distribution, so we fitted our data we estimated parameters for a normal distribution, you can also estimate parameters for other also other distribution. I would not be doing that but all these can be done using the maximum likelihood method. So, to summaries what we have we, I want to mention that the maximum likelihood method is a way to estimate the errors. Even in what we

have seen is we have not actually gone into the details of maximum likelihood method. But we have taken the results and seen how to use them in any practical applications.

So, whenever you look at any reported values you will always see the value of A is equal to some number is a0 plus minus a1 and this is our average and this is the error estimate. And more another very common thing that you will observe is that in graphs error estimates are shown as error bars. So, for example, you might have a graph of values and what is done is typically… So, suppose you measure the rate constant as, or your rate constant as a function of temperature then you will get various numbers. And so you will get a graph with various numbers.

Now, each of these numbers corresponds to one data, so this is a result of many measurements. So, each of these numbers is a result of many measurements and so this is the mean of all those measurements at this temperature and then there is usually error estimate. So, the error estimate is shown like this. So, that means this measurement is this is the mean value and so it is plus minus this so plus this to minus this. So, then these error estimates are often shown in measurements. These error estimates are very important at times so, when you want to say whether the difference between two values is significant.

So, notice that the error bars, when the error bars overlap then you say that there is no difference between these two values, these two values are almost the same. Whenever the error bars overlap so, the, this goes all the way down to here and this goes up to here so, then in this region there errors overlaps. And when you see that there errors estimate overlap then you, what we say is that you cannot tell whether this is here of here you cannot tell where it is. So, you cannot tell that this is different from this you cannot tell that these two measurements are different.

So, this is another important part of analyzing various experimental data. So, if you get two graphs if you plot your results and you get one graph like this and another graph like this then the first thing you should do is you should look at the error estimate. If this error estimates, if the error estimates are smaller than the, if the error estimate do not overlap, then these two are different. So, error estimate do not overlap, then these two values are different, is if the error estimate overlap then you say that you cannot tell that those two values are different.

So, then you can do two things the error estimate overlap you make more measurements so, when you make more measurements these error estimates will go down. So, when you make more measurements the mean value might go up or down but the error estimate we usually goes down and then you can tell when these values are significantly different.