**Surrogates and Approximations in Engineering Design**
**Prof. Palaniappan Ramu**
**Department of Engineering Design**
**Indian Institute of Technology, Madras**

**Lecture – 06**
**Introduction to DOE – 1**

As we pointed out you start with the design space, you know the bombs. The first thing that you want to do is, do perform a design of experiment, get few points where you want to run your simulations or experiments ok.

(Refer Slide Time: 00:30).



So, what is it that we try to do in this whole stuff? So, it is observing some pattern and trying to model it or replicate it, getting some simulation point, getting some experimental point ok. So, there are different ways of doing it; that is what the question is, how do you learn and improve. There are multiple ways; first one is to observe significant events.

So, people always say right like you are not designing for failure, you are actually designing against failure, so you should focus on the failure part ok. When will my structures fail and what should I make sure that it will not fail. So, the significant even in this cases will be your failure part. The second one is your expert opinion, previous knowledge legacy.

Now, if you take a company like Ather Energy which released its first electric vehicle a few days ago. They do not have a legacy, they do not have any legacy, they just started 5 years ago, what legacy would they have? Probably the vehicle that they released is there tend to prototype; that is all what, what legacy will they have they will not have any legacy ok. So, what they can do? They can rely is, they can simulate the events ok. So, if you go to any of these car companies today they have these simulation track, track simulation ok, it has all this, they call it a cobbler road ok.

So, it has this gravel based roads and your truck or your car has to go through that and they would have put it, they would have instrumented the entire car starting from your, you know steering column to your you know, your wheel differentials and what all you will look at ok. For instance recently I learned that they even put a sensor on the for the noise ok, they do not want this chattering noise that of course, it is related to the ==m b h== issue, but they do not want that propagate all the way up and they even see by how much your suspensions are going, because they know the landscape, they know where your vehicle is going, they want to understand where, how much your suspension is reacting to that ok. So, those are all beautiful.

So, the point is, I can simulate today, I do not have to wait for 100 years to create all the possible rare events. I can simulate the rare events today and that is exactly the third part ok. So, this is the advantage sometimes using a computer model; that is what they say before bending the bar, you do not need to bend the bar anymore real, you can bend the bar in the computer and as many times as you want, but it is really not that easy. Sometimes when it becomes complex even computer models are expensive, they are less expensive than real experiments, but they are still expensive and we have also moved to a situation, where we call it digital twin ok. You see you about 20 years ago when you present the final element study, people ask what is the validation ok.

Then people will immediately show a cantilever beam with known solution, they will use a machine then they will say this software is validated and then they will start their study, but today you just need to show only a mess conversion study that is all ok, to that extent, in 20 years we have developed confidence in the computer models. please understand, your particular physical model, you might not have confidence, but the software themselves we now have and we have gone to an extent of saying digital twin

this is my physical prototype, here is my computer prototype and this is the digital twin of your real stuff.

So, whatever you do in this its equivalent of doing it in my experiment. So, I can crash this computer online on, on the computer and it is as good as crashing it physically ok. So, to that extent people have developed confidence. So, in that context only most of these things make sense ok. So, I can simulate the events learning by evens and bas knowledge is a slow process of course, and people are still struggling to understand how to import knowledge between, I mean we sometimes even in a systematic institute; like ours when some staff not faculty, but one staff there they usually stay at a place for 3 years or 4 years and after that they go to another department right.

So, still we have this issue of importing their knowledge into the next person, because there is a systematic process right; that is when workflow and all that comes into picture. it is not person dependent anymore, the person is only a data entry operation, everything is recorded there, otherwise they will say that person I remember you send the grades ok, it is somewhere there.

Now, this email can be opened by everyone, whoever has access can open the email and it is on recorder so is the case with workflow. So, the point is this knowledge management is a big deal anywhere ok. So, now, people are trying to see taking brackets from 30 years old, there is a company where you are trying to work, they have like 30 year old bracket, it is the same bracket, but it is evolved, the design has evolved and for 30 year they have used the bracket in that particular area.

Now, what they are trying to do is, doing backward kind of an engineering to understand how this design evolve, because there are about 6 managers who have changed in 30 years and they do not know why your particular Philat was introduced and all changes has resulted in optimization and better performance. There is no doubt, but they are trying to understand that design evolution. Now they are using data based techniques to understand that design evolution ok, so because there is so much of data from their side to understand that ok.

So, the simulation helps. So, structured experimentation with the study of outcomes and what you are going to do with the outcomes will help you here and George E. P. box is a big statistician ok, is a very big statistician, was introduced a lot of design of experiment

a statistical techniques, everything he is like fisher, fisher is fisher information matrix if you know he is like fisher, he is a big guy.

So, what he says is to find out what happens to a system when you interfere with it you have to interfere with it not passively observe it, you cannot wait for 100 years for a rare event to happen if you want to understand why under rare can I mean rare input conditions what happens you have to perform it to understand what happens.

(Refer Slide Time: 06:51)



So, this is this graph is important because this similar to your final element, the truth is like this ok, but this is what your experiment tells you, probably you could have done a little better maybe it was like this, but still you might not be able to capture the truth because the truth is something that you do not know ok.

So, always what we try to do is it is not enough if you get a representation, you also should have some information on how far you are from the actual because its only an approximation its, an approximation is, but how far it is from the approximation that is a good question. So, we experiment, we observed then we conclude this is a cycle of course, what we are trying to do we have already seen Y equal to f of X a well defined experiment will take you very close to your knowledge curve. The whole idea is to get some insight into the process or product, please understand you are not trying to find this and then go and say this is the optima yes you can still do it, but your overall idea is to get some insight into the process or product and then based on this output you can go and

tweak your inputs so that you can change it out to your target values that is in one sense that is also optimization.

(Refer Slide Time: 08:17)



But there is a small problem with screening this is called the design space screening. There is a curse of dimensionality something that we discussed just a while ago about finding gradient how do you find gradients in a numerical sense I do finite difference ok, how many points do you need for fine and difference, I need if I am using central difference I need 3 points, 1 point delta x my positive delta x negative delta x, this is for 1 dimension, for two dimensions how many points you need, I will need 5 points, for n dimensions it will increase in that sense. Similarly if you want to do integration I will do some trapezoidal rule or I will use some Runge Kutta approach something like that when you increase the dimensions you cannot handle this situation.

So, that is called curse of dimensionality when the dimension increases most of your classical techniques will stop or choke, the same is this idea is people kind of wont understand how your response varies in the design domain. So, what they do is simply they try to compute their gradient ok. If your function is going like this ok, it is kind of a uniform slope. So, if you get the gradient information at different points it will show you like this. So, you know that it is a linear function you just need only three points to approximate it done, but the problem that is the idea that is given in this particular slide if you see it says dou f by dou x sorry, dou f by dou x equal to 0 meaning the gradient

information for all x belonging to the designed domain the variable x i can be safely neglected it says x is its not sensitive anymore meaning y is not sensitive to x because when you change x that is what it says when you change x you get 0.

Similarly, if it is a constant that is the example that I gave you, but not equal to 0 then the effect of the variable is linear and it is an additive model you can do an additive model, in case here it is given as a function g of x i ok, so it could be a non-linear function also it is not a linear function anymore for all x belonging where g of x is not a constant because if it is a constant then it is the previous case then f is non-linear in terms of x i this is something that you can tell there is an interesting part g of x i for all where this g of x i is not constant x i x j everything is not a constant then if f is non-linear in x i and involved in interactions with x j that is important ok.

So, this is taking together here it is taking 1 at a time g of x 1 g of x 2 g of x 3 like that here we are taking g of x 1 x 2 x 3. So, there is a difference between the third and the fourth there is interactions are also involved in this particular case. So, but getting this is not easy because you do not have the you have the function f then its somewhere you can do all these things, the point is you do not have this function I told you just forgetting this information you need 3 for a point, please understand this is not for the entire function 3 for a point, you tell me ten points 30 information I need to give you that is not easy.

So, typically you cannot get this entity into a, but then you expect the model to be uniformly accurate across the design space that is also not such a property where you want to get the information is called a space filling property ok, because I do not know how the function varies and I cannot favor a particular. So, now, someone walks into my room let us say and then they say select 1 of the students we will send them to a conference I am just saying some Dean comes in and then he says we chose your class and there is some conference we will response one student select 1 of the students and let us assume that these 3 are not my students otherwise I am naturally biased.

So, I do not know I do not know any of you right. So, similarly these three also I do not know they walked into my room and they said choose 1 of those students I will send them. I cannot choose because I have no information on who you are, how you perform, why you should go to this conference, will you benefit out of these conference I have no

information. So, unless I have those information I cannot say that let us put 40 percent of the points here and the remaining points will use it here I cannot say that. So, you expect the model to be uniformly accurate or uniformly varying across the design space that is what you should do.

So, there are 2 meanings to this particular statement ok, 1 is I have used the word model, but it also useful for its can also be used for function you expect the function to be uniformly accurate then the model that you build also should be uniformly accurate it you cannot say at that point it is not very accurate, but at this point it is that is allowed provided you know this is your zone of interest. If you do not know anything then you cannot say I built a function, which is not accurate there let that place you cannot use the function you can use it at any other place unless you have a bias you should not do that bias means you know then you have the knowledge of that ok.

So, your model your sampling plan also should govern that that is important unless you have bias why will you only choose few points it should be uniformly distributed, uniformly distributed in the sense not uniform distribution not, not uniform sampling ok, in terms of information sense in terms of locating the design of experiment points there should be some kind of a space filling nests. What is the space feeling nests, given the 10 points I want to explore this come goes back to your visitor example ok, given the 20 points I should explore as much as I can within the design space.

(Refer Slide Time: 14:40).



**Types of Experimentation**

- ♦ Trial and Error
  - " Quick and dirty "
  - Get some insight with little time by varying parameters randomly
  - Can't build knowledge with it
  - Often results in "band-aid" fix

- ♦ OFAT – One Factor at a Time
  - Quicker to capture effect of each parameter independently
  - Fails to capture interactions

- ♦ DOE (Design of Experiments)
  - "Statistically designed" experiments representing design space
  - Captures interactions of variables

So, there are different types of experimentation of course, one is the trial and error ok. There are some advantages and disadvantages of that which in a scientific sense this will not work ok, it is a quick band aid fix and all that is fine. There is another one which is called the one factor at a time what this is you change only one factor you freeze the other factor, this is great it gives you some idea about a factor, but you have to do at least two experiments per factor, please understand and it only gives you a linear relationship ok. It could be a non-linear in x 1, but you do 2 points it will only give you a linear relationship, it also does not take into account your correlations because you are fixing x 2 x 3 x 5 x n. So, if you change x 7 what happens to this x one it might not remain the same. So, it is partially only helpful, but if you do not have anything one factor at a time is a great thing to do.

The classical what we call like Taguchi based technique OA and all that has taken a backseat as of now in non manufacturing applications especially with the days with the design of computer design and analysis of computer experiments is statistically designed ok. So, what happens is how do you ensure the space fillingness ok. So, there is some small factor and then they do minimizing the maximum distance between any 2 points. So, there is a criteria that you want to maximize or minimize and you will get a value out of it. The interesting point of the statistical design is you rerun that algorithm you will get another set of points. So, the points that he got will not be the same point that I got and the points that I got will not be the same points that he got.

Which means that value it could be a local optima if you can get an function of that phi or theta that we are talking about it could have multiple optimized and it is also dependent on the points that you start with ok. So, this is an interesting point and this is why it is interesting. So, why is it interesting; every time we were getting different values then how can we reproduce the results. This captures nature that is why people say that you should introduce randomness as early as possible in your design because in reality it is random and see I manufacture a auto I mean automobile and then I say I can only give it to good drivers can I sell a vehicle like that I have no control on who will buy my vehicle and run at anyone who pays 1 lakh rupees can buy my vehicle. I cannot say oh that guy is 18 years old he could be a rash driver I cannot give it to him, this guy looks like a sadhu, so, I will sell it to him. You cannot do all that anyone who gives you 1 lakh rupees you will have to sell your automobile.

So, now the point is with a guy who rides it around about 40 kilometers per hour and there is a guy who rides at about 80 kilometer per hour will suddenly rise it will suddenly drop it will hold it harsh I will have to bring all this variability up front in my design so that I can tune my engine accordingly, otherwise I say that no this automobile was designed for Chennai you cannot go and ride it in Bangalore because the temperatures are different, you cannot say that, it should today they are you know like putting engines inside dusty chambers dust chamber and then they are with the circular fan they are putting dust inside and then it is doing it after like 12 hours of dust exposure they take it out and then they are checking whether the links are working then again they expose it to another 6 hours of dust they are taking it.

So, I cannot say that this will be this will run only in tar roads who knows I will buy it and then I will take it to a village and ride it, but still the vehicle should right and then where you do not have tar roads there is it is exposed to dust so, you it should. So, the point is you should bring all these random conditions up front and this is one of the cases. I am designing it in a random sense. So, I will do to this one design of experiments and I will see if I come to the same conclusion I should come to the same conclusion you cannot say I will have two riders and the efficiency will go down no it cannot it cannot vary too much of course, it will vary a little bit, but it cannot vary too much.
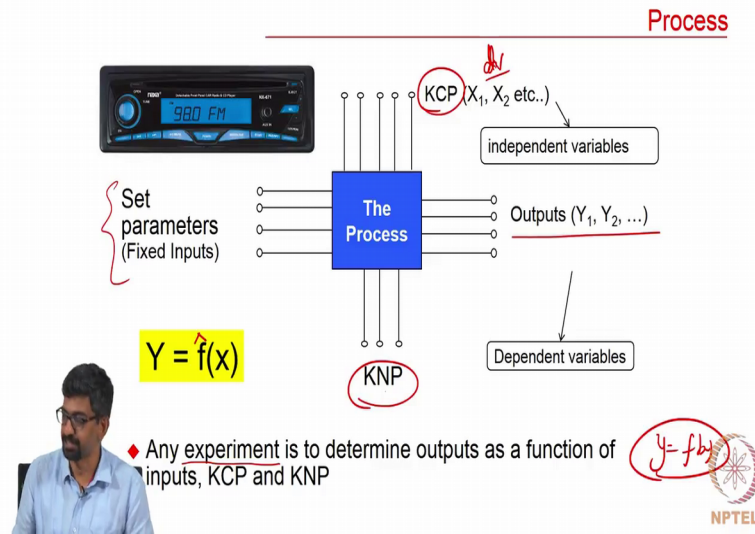
(Refer Slide Time: 19:27).

This is a historical perspective its very interesting someone started with a cure for cure like a particular like a disease kind of a stuff then actually it had this is this is the way that I spoke about fisher, he is a big guy in the scene of experiments ok, but as you can see it actually picked up in the modern era between 71-90 where they wanted to do quality products during and at the end of world war 2 ok, it actually started during world war 2 ok.

Because they wanted to very quickly do it in world war 2 you cannot say write a proposal let us say let us import a vehicle and then let us do it can I do that right it has to be done quickly. So, that is the biggest advantage of this Taguchi stuff from the outcome you can directly you can identify you and then you can control factor analysis and all that is the biggest advantage of that. You can see it actually started in agriculture and then it started going even today. If you want to do soil quality check they do not go and dig it in a land and then bring it ok. There is a scheme which is a very nice design of experiment. You have to dig it in the 4 corners, you will already get in the center of the entire thing that you want to do, you will have to mix it mix all these soil together you will have to sieve multiple times then again there is another experimentation you will have to bring and mix from other another 4 points and then only you they will give you a soil quality ok. You cannot just go dig into one point and then go and give it and that that will become a local effect.

So, you let it do there is a particular way in which you will have to sieve all of them filter meaning filter and then again you lot of mix them then again sieve and then only you will have to go and give it. So, it actually started in different applications as you can see.

So, this is about this is called p diagram. So, if you look at it what they are saying is this is your output Y 1 Y 2, this is your key control parameters X 1 X 2 which we usually call the design variables right. These are the independent variable this guy is our dependent variable something that we do not discuss and slightly out of scope for us is this guy which is called the key noise parameter which I do not have control. I design a wind turbine and the wind turbine is dependent on the wind speeds I know that, but do I have control on the wind speeds I design automobile for my client.

Do I have control on my client? I cannot go right to my do not drive fast if you drive slow and do not go this slow, slow, slow, slow slowly go or the otherwise the suspension will go off you do not have control on that. I write a book do I have control on who reads a book, but that is the problem with students right. They write a paper ok, they do not imagine that it will be read by anyone they think it is the professor's headache to wait it out. So, I need to think is this X 1 or X 2 which plot is this we need to read all that and we need to have our own inferences about that anyway fine.

So, this key noise parameter is something out of scope for us, but this is a general p diagram that we have and then you set your parameter you fix your inputs, this is the process that we are trying to approximate. This example is for a slightly different example on slightly different course on robust design so, but what you are trying to do is any experiment whether computational or a physical is to determine your output as a
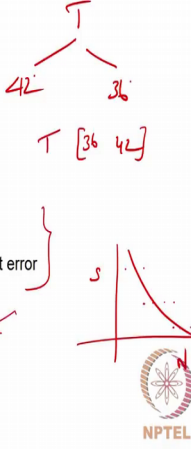
function of input sorry, it is here already X is your input Y is your output. We are going to do f hat which is your approximation.

(Refer Slide Time: 23:16).



These are the key definitions we will not use too much of this, the response is your output variable which is your Y. The factor is what we call the design variables in Taguchi orthogonal array it is called factor and the values that it can take is called level it is here this level and factors are. So, if you take temperature for instance as a factor, it can take 42 degrees, it can take 36 degree. These are called two levels, but in the days what you lose is temperature and then you will give some bonds that is all ok. Then replication this becomes important this is what I told I do 2 different design of experiments I should arrive at the same conclusion otherwise there is a problem ok.
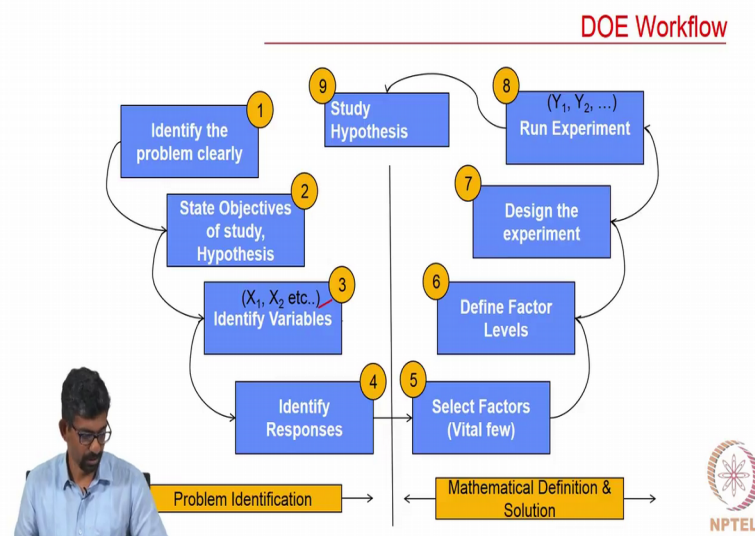
So, I should be able to replicate it interaction becomes important, now people are saying like you know like what is that is what data now you are trying to do ok. For instance a cancer research a problem is because there are so many factors that you do not know ok. So, you they are still trying to identify factors that contribute ok. Now you have a list of factors and you have a list of attributes and then you have a list of responses. The interesting part is you need to find out which interacts with other to give that case. So, there is an interaction between these case and there is an interaction with respect to your output also. So, that is way cancer research is very difficult and the problem is I do not know did anyone listen to this double e, there is this one guy called Vidya Sagar

muttuKamali, he came Indian the controls conference they conducted and then he gave a talk here on, he is an electrical engineer where he is used factor analysis and all that to the core and currently is doing cancer research.

So, it is a very interesting talk what I learned from the talk is why this is very challenging is you have a cancer site let us say, you take three cells from the cancer site and each and each of them have a different characteristic. So, there is no pattern other ways by the time our people would have broken all that thing right because we are so mature now in pattern recognition and data mining, you cannot do any of those stuff for each person why even each person for each site is, it is different the characteristics are different you cannot go with a particular attribute value. So, there is an issue ok.

So, that is when the interaction becomes very important what is the correlation between X 1 and X 2 and how does it affect why and for you to capture this you need lot of data that is one problem in our context like in cancer research we do not have that much data we have a lot of cases, but we do not have a nice data logging system ok. And randomization is required to remove bias why ASME says sorry ASTM standard says that you need to do 3 at least at different design points to do your SN curve is also to avoid bias ok. So, in order to avoid bias we need to do random designs. This is the general DOE workflow.

(Refer Slide Time: 26:27).

So, the deal is you this is like a general design of a I mean new product development equivalent ok identify you will have to identify the problem clearly, you will have to state the objectives what your hypotheses are you will have to identify your variables X 1 X 2 you will have to identify.

So, this becomes important because at this point it is the trivial meaning from that you will have to go to the vital few ok. So, you say what all is responsible then you take it from there to here select the factors few and then you define your factor levels or you say these are the bounds I will use it as this distribution you do your design of experiment run your experiments and then you can do whatever you want study your hypothesis fit a surrogate and all that. I will wrap it up here.