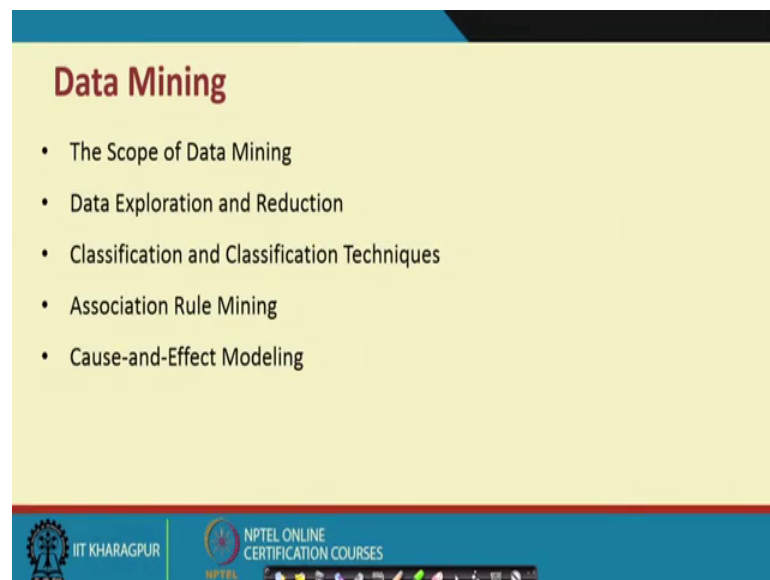**Business Analytics for Management Decision**
**Prof. Rudra P Pradhan**
**Vinod Gupta School of Management**
**Indian Institute of Technology, Kharagpur**

**Lecture – 39**
**Predictive Analytics: Data Mining**

Hello everybody this is Rudra Pradhan here. Welcome to BMD lecture series, and today we will continue with the predictive analytics, and that too coverage on data mining. In fact, this particular technique or this particular tool is very vast, and very useful for business problems; where we have lots of complexity and the kind of dynamics.

And with the help of this particular technique again we like to understand the particular data structure, and then create a kind of environment through which we can do better predictions and better forecasting's. In fact, in the data mining we have a couple of items through which you can actually analyze the kind of problems, and then come out with a kind of solution through which the particular decision will be very effective.
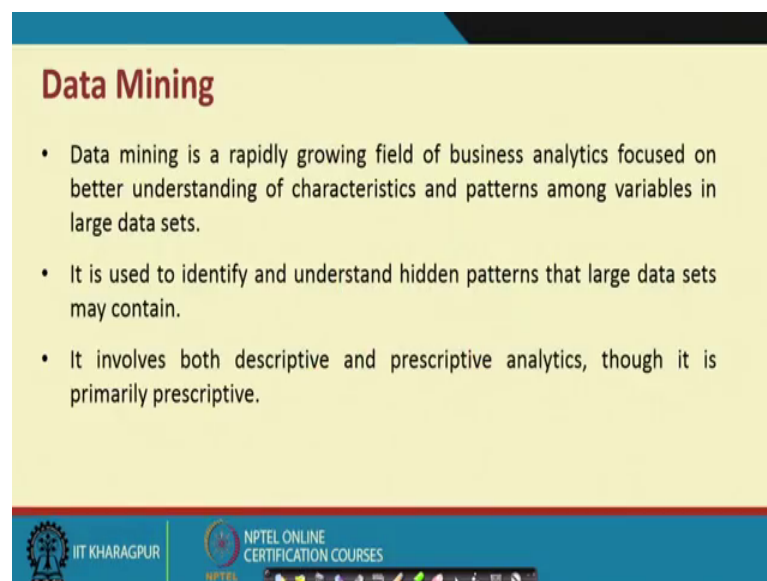
(Refer Slide Time: 01:06)



In the data mining so, the typical discussion will be the scope of data mining, then we will go for data explorations, and data reductions and again so, classification and the classification techniques. So, it is with respect to association rule mining and cause and effect modeling. In fact, we have already discussed couple of predictive analytics starting

with a simple regressions that association, classification, cause and effect modeling in the context of you know regression modeling.

So, that means, altogether we have already discussed certain techniques in the predictive analytics, which are actually part of this particular data mining process. But here one thing we like to give high stress on the kind of cluster analysis, again this is a kind of classification.

So in fact, we have already discussed so many classification techniques like random forest and support vector machines, and here we have a kind of beautiful classification structure through cluster analysis, and for this kind of environment; means, business environment this particular classification technique can give better kind of inference through which you can understand the problem more accurately. And then we will create a kind of structure through which you do the prediction and the kind of forecasting's. So, typically we start with what is exactly the data mining.

(Refer Slide Time: 02:46)



And then we will connect with some of the techniques through which you can actually understand the kind of data structure. And then analyze the problem as per the particular management requirement.

So, what I have mentioned already. So, it is a kind of beautiful technique through which we business analytics can be focused through which to understand the particular business

problems as per the big data set, and the kind of complex number. If you know, means multi multiple number of variables through which you can analyze the particular problem. It is used to identify and understand the kind of hidden insights through in the data and the kind of the kind of variables which you have identified.

And then we develop a kind of set up and the structure through which you do the kind of prediction, it is a actually more or less same like neural network structure random forest and support vector machines or any kind of regression kind of techniques; that means, actually the fact is that we have a large data, and that too with a more number of variables.

And in one instance we have a kind of structure then using the data we like to just test and verify and find out the particular moulding set up through which you can better prediction and better forecasting for the future. And in another environment sometimes we do not find a any kind of insights in built insights initially. And through which you can just verify and then develop a kind of set up through which you do the prediction.

But in the second case typically the structure is it to develop a kind of system, and set up through which we have to again trained then verify then test and then finally, do the kind of prediction and forecasting as per the particular requirement. And cluster analysis and that too data mine in general data mining process is like that to understand the particular data set, and the understand the particular problem with respect to more number of variables. And the kind of dynamics, and then develop a kind of system and structure through which you can actually go for better predictions better forecasting's as per the particular management requirement or problem requirement.

So, it also involves both descriptive and prescriptive analytics; that means, the things which you are discussing right now is a kind of predictive analytics structure. But what I mentioned earlier that all analytics are very closely connected. And in fact, we are still we are discussing predictive analytics. But we need actually descriptive structure and predictive structure; that means, the previous or the previous kind of set up and the kind of future set ups through which actually we can analyze the problem more accurately you know as per the particular requirement, and the kind of best management needs.

(Refer Slide Time: 06:01)



So, likewise so, we have actually kind of system. And so, far as scope is concerned that to in a in the data mining. So, what I have already mentioned. So, we can actually go for data exploration; that means, more visualize and more insights to in the data and sometimes we will go for data reduction process.
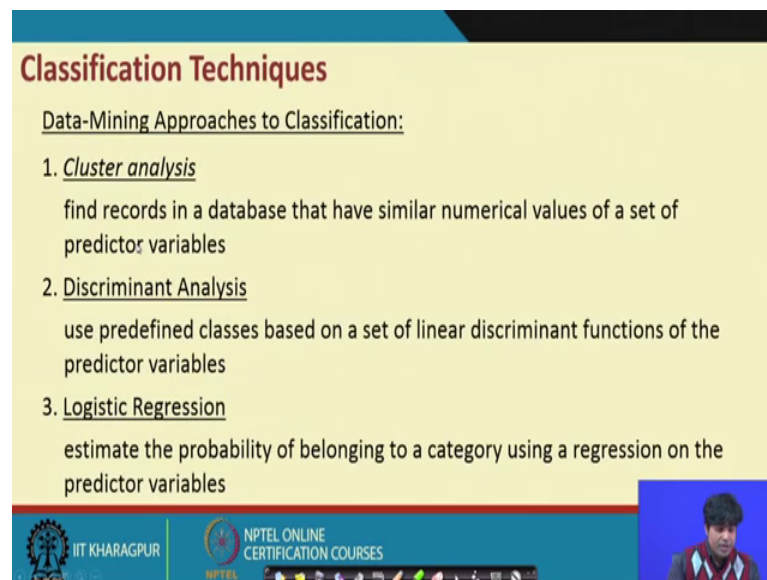
In fact, there is a kind of technique called as factor analysis through which the large set of variables can be transferred into a kind of you know, pools through which you can actually have a better clusters then better structures to understand the problems and to go for the kind of predictions. But in this particular lectures we typically go for cluster analysis to with a particular set up we develop a kind of cluster to understand more accurately, the data more accurately the problem, and then the kind of prediction requirement.

So, we have data exploration data reductions data classifications, then association cause and effect modeling. More or less we have already you know we have already covered all these things little bit high and low depending upon a kind of items which we have discussed earlier.

And in this case just we connect since we are discussing actually classification clusters in this particular you know, unit starting with neural network then random forest support vector machines. So, in the similar clusters we are just connecting a technique called as cluster analysis. Then we like to know how the problems can be analyzed more

accurately more you know, in a interesting way by developing a system. And developing a kind of classification through which you can understand the issue more accurately and then go for the kind of management requirement.

(Refer Slide Time: 07:54)



And ha in fact, in the in the kind of classification techniques, particularly the things are very connected with cluster analysis, discriminant analysis and logistic kind of regressions. means the problems which we have already discussed here exploration classification; that means, see here we what we have actually here. So, this is actually one kind of package the kind of you know. So, data exploration and data reductions, then the classification association cause you know. So, these are all actually lots of you know, tool tools are there lots of tools are here also, lots of tools are here also, lots of tools are here also.

So, technically these are all the base through which actually to visualize the data understand the data get insights from the data, then understand the problems. And then develop a system and structures or develop the kind of algorithms through which actually you can predict the environment as per the particular requirement. And then the decision by default will be very effective. So, comparing these kind of structures, if you go little bit more. So, then you will get to know the that in the kind of classification structure, and the counter parts. So what we have actually you know connecting here; the cluster analysis discriminant analysis logistic regressions.

So, that means, the these are the items so, where actually a means these items can be actually connected and with their. In fact, in a kind of similar kind of techniques pools, or techniques basket, where we have a kind of problems, where more number of variables are there having more number of data.

And within the more number of variables, we have some kind of quantitative variables some kind of qualitative variables. And then how we can actually recognize this particular problems depending upon a particular response variable whether it is qualitative in nature, whether it is in a kind of quantitative in nature then, and develop this system through which you will do the predictions as per the particular problem requirement.

And in this context, cluster analysis discriminant analysis and logistic analysis, logistic regressions are in a kind of more highlighted pictures right. And in fact, in the previous kind of lectures we have already connected or we have already discussed the kind of discriminant analysis, logistic regressions. In fact, in the that is called as a qualitative response regression modeling.

In fact, we have discussed binary stress models that is the linear probability model, which is actually more or less similar to the kind of discriminant analysis, then we have we have discussed logistic regressions followed by logistic functions then similar to we have discussed actually the kind of normal distribution followed by normal distribution that is profit model. And again in the similar scenario we are actually now highlighting the kind of you know structure called as cluster analysis. This is again the kind of classification structure. And through which you will analyze the understand the problem and analyze the problem as per the particular requirement.

Here the means these are the items or these are the techniques predictive analytics tools or techniques, where some of the variables you know, will be kind of qualitative and may be binary in structure or categorical in structure. Then depending upon the kind of in data set and the kind of variables. So, you first check the kind of structure then you can pick up a particular technique through which you can understand the data visualize the data properly.

Then may be go for some kind of reduction mechanism. And then a built a kind of system through which you can do the better kind of prediction better kind of set ups

through which you can actually analyze the problem, and come with effective management decision.

(Refer Slide Time: 12:02)



And in fact, since our highlights of this particular lectures though we have already discussed discriminant analysis logistic structure that is qualitative response modeling.

And in this case, we typically highlight or give high stress on cluster analysis. when you talk about the cluster analysis; that means, with a kind of large variables, and large data set we like to prepare a kind of cluster where we can bring the kind of kind of link or structure, where it is a kind of homogenous kind of set ups through which you can do the better prediction, and the kind of better forecasting. And in fact, in the in the kind of you know classification structure that too the particular kind of structure called as cluster analysis.

So, here we have 2 different mechanisms through which you can understand this particular you know predictive analytics tools. And then connect a problems, and check how it is very effective for the management requirement or business requirement. So, there are 2 2 different structures through which we will do the clustering. And one is called as a hierarchical clustering, and then k means clustering. And in the hierarchical clustering's, we have a 2 different methods altogether.

So, agglomerative methods and divisive methods through which you can do the kind of structuring, and then we will understand the system and then find out a kind of set up through which you will do the predictions. So, what I will do? So, I will give you little bit exposure about the hierarchical clustering, and k means clustering.

And then we pick up a kind of problem, then we connect with the both the techniques, and then we see how this this particular technique can be very effective to understand the data understand the kind of problems. And then come with a kind of structures where we can do the effective predictions, and effective kind of forecasting as per the particular management requirement, and the kind of business requirement.

(Refer Slide Time: 14:15)



So, let us come with first the kind of clustering with you know, the kind of structure which we called as hierarchical clustering. And in the hierarchical clustering, particularly what we have mentioned earlier so, we have a 2-different set ups altogether, agglomerative and divisive. And in the case of first methods; so, there are 2 items through which you can analyze and to understand the system that is called as dendogram.

And then Euclidean distance and it is a kind of you know, graphical visualization through which you can understand the particular structure, and then do the set ups for the prediction and forecasting. So, this is what the Euclidean distance and in fact it is all actually what we can say that you know.

(Refer Slide Time: 15:00)



It is a kind of mathematics through which we can understand the particular structure. And this is what the Euclidean distance, and this is what the particular structure. And I will let the dendogram structure through which again you can analyze or you can understand the particular structure. So, this is what simply through mathematics, we can understand the Euclidean distance. And then we use this particular structure through which you can do the better clustering to understand the problems, and the develop the set ups for future prediction and future forecasting's.

So, accordingly so, what I will take you to the particular structure to understand how it is actually effective as per the particular requirement in the clustering. you know, in the cluster analysis particularly this particular method so, we have a single linkage clustering complete linkage clustering average linkage clustering.

So, like we hierarchical clustering. So, we have a different kind of you know models through which you can understand the data, or understand the insights from the data understand the problems then develop the structure through which the manage the problem can be analyzed more accurately as per the management requirement or the kind of business requirement.

So, in any case it is a kind of different means, it is a having a flexibility systems. Within a particular technique, we have a lots of flexibility to analyze the data or get the insights

or explore the insights best insights as per the particular, prediction requirement and the kind of management requirement.

(Refer Slide Time: 16:49)



So, accordingly s, let us start with a simple examples, and let us say this is a kind of college and university data. And we like to classify the college and university with respect to several attributes in this particular examples we have here's schools, and then different types schools and colleges. And then we have a different attributes like radiance curve or acceptance rate expenditures students. Or this particular problem we have we have already discussed with a different technique context and problem context. But same problem we are now putting here in the clustering analysis.

Now, here idea is that that with respect to all these attributes, we have to define different clusters, then means what is the management requirement that once you develop different clusters as per the particular requirement, then the policy structure the strategy structure will be accordingly different. And as a result, the or in the entire management system will be very effective and efficient if you develop different clusters. And then plan you know plan accordingly for a particular clusters instead of you know, thinking about to or planning something in a kind of aggregate level.

So, in that context cluster analysis and the data mining package will be very effective for the management requirement, and for the kind of management decisions. So, likewise we

have actually. So, let us see here; so, this is what the problem, and then what will you go so, if you analyze.

(Refer Slide Time: 18:28)



So, this is what the dendograms, the dendograms structure will be like it is like decision tree structures like what we have already discussed in the case of you know, random forest, and through which how we will connect properly with these attributes, then we will classify the particular structure. And develop the kind of clustering through which you can analyze the problem.

So, this is a simple look about the dendograms through which clustering can be done with respect to college and university, and that too with respect to different attributes. So, I will take you to a kind of software's, then I will highlight with a problems, and then I will show you how the dendogram can be coming and classifying the things or the clustering the things as per the particular business requirement or the problem requirement.

(Refer Slide Time: 19:21)



So, then accordingly so, this is another kind of visualization; that means, we have a flexibility system now, if you choose certain dynamics in the kind of structure, then you will get different kind of result then finally, we have a plenty of different structure. Or models or set ups, then we pick up a particular set up which is very effective for the particular problem, and the kind of management requirement.

(Refer Slide Time: 19:45)

(Refer Slide Time: 19:46)



So, these are all data set, and through which actually we analyze the problem. And then finally, we have in this case in the particular analysis we have 4 different clusters 1, 2, 3, 4. Depending upon the kind of college and the university, and then and the kind of attributes.

So, here in this 4 clusters. So, the first clusters is having 23 colleges, and second cluster is having 22 colleges, third cluster is having 3 colleges, and 4th cluster is having one college one college. So, that means, technically now, so far as management requirement is concerned management decision is concerned. So then effective kind of strategy or effective kind of policy can be implemented differently instead of putting the general kind of strategy or general kind of policy.

So, in that context cluster analysis is always at harsh labels to understand the system. And predict the kind of system as per the requirement.

(Refer Slide Time: 20:45)



So now with this particular data so, we have actually this classification, these are in the cluster 3. So, these are university is coming. And this is what the attributes through which you do the kind of classifications.

(Refer Slide Time: 20:59)



So, we will we will we will again analyze with with a kind of software's, and this another kind of examples, this is a credit approval case.

So, credit approved and not approved like in the case of discriminant analysis logistic regression. So, this is also similar kind of structure. So, here so, whether credit approved

or not approved, and that too that too exclusively depends upon several factors and several attributes. So, in this kind of example we have several factors here. So, these are all factors credit approved yes or no. Then we have several factors and through which we have take the decision so; that means, we have actually data set.

And then we have to connect a model then after that we can analyze that how a this kind of effective structures, and that too how is the kind of decision about the credit approval with respect to these indicators. So now, now with the help of these particular structures we can actually get better insights into understand the particular technique.

(Refer Slide Time: 22:04)



So, this is how attributes, that too in the independent clusters.

(Refer Slide Time: 22:08)



Now, having a kind of classification structures, let us say this is bar and through which. So, we have a kind of means if you specify a particular attributes, let us say here in this case it is a 640, then on the basis of that you know, you can give a kind of you know, structure whether to accept and reject. So now, having the kind of bar fixation.

So, these are the kind of means data set which are actually approved and the another group of data vary is not; that means, it is a kind of classifier you find out a kind of kind of structure, through which entire system will be classified into 2 different structure, and then you will analyze as per the particular requirement.

So, that means, it is a kind of interesting structure, through which actually you can develop this, another way of classification.

(Refer Slide Time: 23:00)



That means, if you change the dynamics at this particular, systems then you will find the particular kind of classifier through which you can understand the data then predict the environment as per the particular business requirement.

(Refer Slide Time: 23:12)

(Refer Slide Time: 23:14)



So, likewise you know, this is what you understand the particular systems and, and then you analyze the problem as per the particular requirement. In the classification structure like we have actually discussed in the case of neural network. In fact, in the random forest and support vector machines; so, every times you can since we have a large data then you know, we first pick up data develop the model, then train the models and then validate and test them finally, fix the kind of structure through which you can go for the effective prediction and effective, forecasting and then effective management decision.

So, it is also same thing here. So, we have here also the training structure validation structure and testing structures. So, this will give you kind of efficient system through which we can actually we will confident that the particular structure which we developed particular everything which we have developed is very effective of the kind of problem, and the kind of management. So, as a result means these are all more or less, same and through which actually you can work and then you analyze.

(Refer Slide Time: 24:25)



## Classification

Example (contd.): Partitioning Data Sets in *XLMiner*

▸ Partitioning choices when choosing random

1. Automatic  60% training, 40% validation

2. Specify %  50% training, 30% validation, 20% test (training and validation % can be modified)

3. Equal # records  33.33% training, validation, test

▸ *XLMiner* has size and relative size limitations on the data sets, which can affect the amount and % of data assigned to the data sets.

So, there is no hard and fast rule what should be the kind of training set, and what should be the validation stress, like we have already highlights this is in the case of neural network random forest and support vector machine.

So, anyway this is what actually; you can actually fix up the particular structure, how many data points you can put for training how much data point we can put for testing. But ultimately, we need a efficient structure through which you can actually look for the efficient prediction efficient forecasting and the kind of management requirement.

(Refer Slide Time: 25:00)



## Classification

▸ Example:  Classifying New Data for Credit Decisions Using Credit Scores and Years of Credit History

▸ Use the Classification rule discussed earlier:

Reject if 0.095(credit score) + (years of credit history) ≤ 74.66

| | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 55 | New Data to Classify | | | | | |
| 56 | | | | | | |
| 57 | Homeowner | Credit Score | Years of Credit History | Revolving Balance | Revolving Utilization | Decision |
| 58 | 1 | 700 | 8 | $21,000 | 15% | |
| 59 | 0 | 520 | 1 | $4,000 | 90% | |
| 60 | 1 | 650 | 10 | $8,500.00 | 25% | |
| 61 | 0 | 602 | 7 | $16,300.00 | 70% | |
| 62 | 0 | 549 | 2 | $2,500.00 | 90% | |
| 63 | 1 | 742 | 15 | $16,700.00 | 18% | |

And in fact, some classification with the particular data. So now, the decision can means, the kind of model can be different if you again change the particular attribute.

(Refer Slide Time: 25:07)



Then fix the kind of range, and again you will find 2 different classification group; where it is accepted and the kind of rejected.

(Refer Slide Time: 25:20)



So, this is again similar kind of structuring. And another kind of structure called as k means structuring, and here's it is also more or less same. So, how you have to develop a system, through which you can actually understand the particular insights develop a

particular clustering, through which you can analyze the problem more effectively, and develop the kind of structure through which you can actually take a kind of best management decisions. So, in order since we have already discussed 2 kinds of you know, clustering techniques.

So, simple structure then the k mean structures. And what I will do? So, I will take you to a kind of problem, and then I will highlight how clustering can be done, and then what kind of insights or we can observe from the data set and in the kind of problem, the you will find the kind of effective system from the data set, and the kind of variables structure. And then we will think about the kind of management decision, and the kind of management requirement.

(Refer Slide Time: 26:29)



So, so, what will you do? In the in this kind of to understand the clustering analysis. So, I will take a here the particular problem. So, here variables used in this particular discussion is this state total population, and net domestic migration civilian move, net international migrations periods of birth periods of death then residential populations 60 under 65, and above 65. So, these are the things in the problem basket.

And we have a big data set and then how you have to actually pick up the particular variables, and how you have to clustering. Then that is that needs actually planning and that needs efficient kind of understanding. And of course, it is with respect to management objectives or business objectives. So, that do the kind of classification do

the kind of clustering. Without clue or knowing the particular hint you may not go for proper classification or kind of proper clustering.

So, obviously having or large data set and having you know, information about all kind of variables to this particular problem to analyze the problem. So, then you think about to what kind of technique and what kind of structure, or what kind of algorithm you have to develop. So that the particular system can be very effective for the kind of prediction, and the kind of management requirement; so, in order to understand, so what I will do with respect to this particular.

(Refer Slide Time: 28:02)



So, this is the original data set and so, then we will go for the kind of analysis.

Say, with respect to both the clustering techniques. So, in order to understand. So, I will take you to the kind of excel sheet. So, the same data here.

So, you will find. So, this is a large data set. And having actually you know plenty of information; so whatever I have already highlighted here, this particular structure. So, the same data pool here's, I am actually taking it to the excel sheet. And then, you go to the excel start this same actually software's you can use. In fact, this can be also done through r software. But clustering technique is actually very useful tool, and most of the software's is having this particular module.

And once you understand, then whether you use r software or you know, you use excel start or use SPSS or starta, there is no such kind of you know, problem. Because it is the indication more or less sames, and the outcome through any software. You will derive is also more or less same just you understand the kind of structure, then on understand the problem then the software will help you get the result. So, that we can analyze the problem more effectively, and then you think about the kind of management decision as per the particular business requirement or the kind of problem requirement.

So now understanding the particular problem. So, what will you do here. So, in the kind of clustering analysis, because we have a plenty of techniques same problems can be used by different techniques. And then analyze as per the particular requirement, but since we are in this particular, class we are our discussion is more or less on clustering

analysis. So, we like to analyze this problem through clustering analysis, and what we have already discussed that, in the in this particular class we have discussed 2 different clustering technique one is called as agglomerative hierarchical clustering. And then k means clustering. So, let us start with the first k means clustering.

(Refer Slide Time: 30:11)



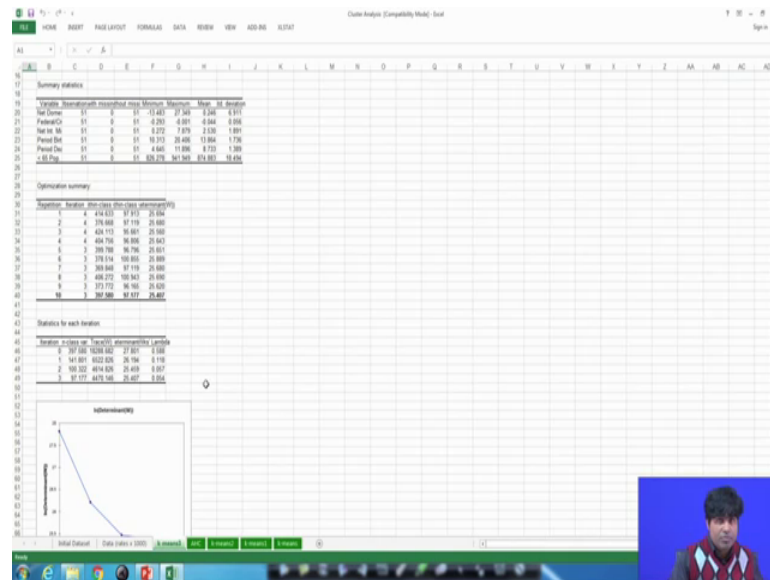And for that we can click here, and then the input box will be coming like this.

So now it is actually clean. So, what will you do? So, you just indicate how you can do the clustering with respect to all these attributes. So, what will you do? So, you just indicate the kind of variable sheet to these are the variables, which you have given the kind of indications. And then in the row variables also we can clean and give the kind of indications in the row variables. We will give the kind of state indications.

So, that the clustering will be done on the basis of state, and these are the attributes through which you would do the kind of clustering like, which you have already highlighted with you know, with 2 different examples that is with respect to college and university and the kind of credit approval yes and no situation.

Yes, here also same things and after giving all the indications. So, this is how the kind of management objective, once you understand the problem fixing the kind of objective having the data and having the variable information. Just you have to connect, how? as per particular requirement so; that means, you need a complete planning before you

come to the software's. And then because software will only be giving you the output through which you can analyze the problem and think about the prediction and think about the kind of insights. But ultimately the planning and the structure which you need to actually operate. So now, after putting all these things we just put and then continue.

(Refer Slide Time: 31:46)



So, this is how the k means clustering's result is coming. So, as usual. So, every techniques you can get a summary statistics, but as descriptive statistic. That is how descriptive analytics is descriptive analytics is the first priority to understand before you go to the predictive analytics and prescriptive analytics. And in fact, in this case we have lots of these are all basic statistics, and this is how the kind of optimization summary and in the statistic of each iteration. These are all k means clustering different class and different iterations, and then again, we will have actually some kind of statistics about the a clustering's.

(Refer Slide Time: 32:27)



So, this is how the indication through which the iteration can be stopped, because once you start like that. You know, clustering can be done. So now on the graphically; if you and a particular structure is coming like this, so then that will give you the signal that your clustering is done, and that too that is actually effective kind of structures. Once you find all these things, then this is how the kind of summary sheet about the validations, and within the class between the class and between the class 75 percent 76 percent, then this is 24 percent and again. So, we have a lots of classifications with respect to different attributes. So, these are all various classification.

In fact, from this result, you can just guess that manually it is very difficult to done. This kind of you know, clustering and this kind of classification until and unless you connect with a particular software, and understand the particular technique.

(Refer Slide Time: 33:30)



And this is how distance between the class enters. And finally, and these are all different clustering with respect to different attributes.

(Refer Slide Time: 33:40)



And by the way you will find plenty of clustering like this, with respect o particular you know, problem and the kind of the kind of cities. And with different attributes like this, this is the class, and then this how the centroid. And so, of course, you know; that means, technically so, this is how the effective clustering can be you know, having to understand

the problem, and to understand the structure through which analyze the problems right ok.

So, what will you do is so, these are the results through which actually you are getting through k means, clustering so; that means, different iterations and different training structure is developed. Then finally, you pick up a particular structure through which you can have actually have a better clustering as per the particular requirement, and then go for the kind of prediction. And the kind of management decision again the same problems can be also analyzed through the other you know, different techniques.

And that too the kind of structure which we have discussed called as dendogram. And the kind of Euclidean distance let us say, go for the particular data set again. So, this is the original data set against you can go to the excel start data analysis, and then this will be under analysing data, and we have already solved this problem through k means clustering.

Now, we can analyze this problem, other mechanism that is agglomerative hierarchical clustering's so; that means, we have discussed k means clustering. And hierarchical clustering k means, clustering actually it is a very different kind of clustering through different iterations. And then final kind of structure through which you can get the better structure or better prediction and against the same structure here. So, we are giving the kind of you know, indication about the variable same the kind of you know, attributes which we have already highlighted here. And then you just put and against continue, you will get similar kind of results like previous case, we have here also some descriptive statistics.

So, in the descriptive statistics. So, these are all various kind of structure, summary sheet, and then the more interesting part of this particular you know, hierarchical clustering is the kind of label bar chart.

(Refer Slide Time: 35:58)



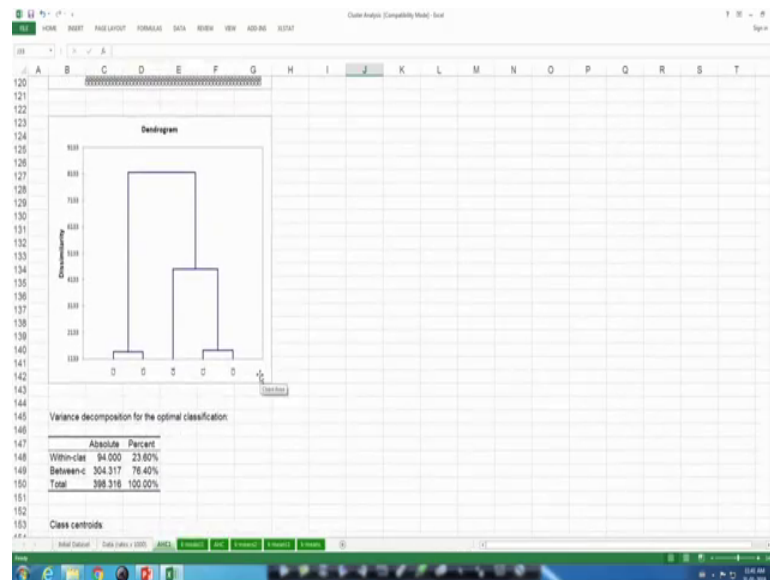And then the kind of dendogram.

(Refer Slide Time: 36:01)



See here this is a actually the kind of you know, structure through which you can actually clustering the kind of you know, data get understand the insights. And build a kind of classification through which you can do the prediction so; that means, it is a kind of the particular technique is kind of effective tools through which you know, you can better understand the data; visualize the data, and then have a kind of cluster. So, that the management problem can be analyzed more effectively as per the business requirement,

and the kind of problem requirement, and do the predictions and do the kind of effective management decision.

So, that is why means it is a very easy to get ok.

(Refer Slide Time: 36:43)



So, this is a dendogram another look and these are all like, this is also coming under k means clustering's.

(Refer Slide Time: 36:46)

So, this is how the validation about the means, various variance de composition; that is called as robustness of you know, this particular results. And then finally, there are different you know, class and then these are all with respect to different attributes. And it is more or less same like as the k means clustering. But in the hierarchical clustering, one of the most important thing is the kind of the graphical structure to understand the particular, link and the kind of clustering.

That is the dendogram and then accordingly we can understand, this is what the distance between that is the Euclidean distance, which is actually since we have actually 5 different class. So, we have a actually the kind of matrix here 5 into 5 and; obviously, 1 to 1. So, the distance will be 0 then 2 to 2; obviously, the distance will be 0 again 3 to 3 distance.

So, as a result all the diagonal elements will be having actually 0, figures and others depends upon how the distance is coming actually to a particular, point. So, these are all central objects and that too with different attributes again.

(Refer Slide Time: 37:59)



So, distance between the central objects again. So, the particular diagonal elements will be 0 as per the Euclidean distance measures, and these are all various clusters class wise. And that too with descriptive statistics, these are all class wise descriptive statistic. And then we can get better to understand now you see here in the third class, we have actually sufficient number of kind of data points. And then the observations, and the as a results

you can actually this is more, summary statistic so; that means, it is a more interesting is the kind of this is the kind of classification through which actually, clustering is done and the kind of analysis is actually done and accordingly.

So, this is what the kind of you know, structure.

(Refer Slide Time: 38:49)



And what we have already discussed this is how the graphical understanding about the you know, validations to justify. That you know, this is what the effective clustering through which you can do the kind of you know, understand the particular you know, problem and do the kind of prediction as per the particular you know, business requirement and problem.

(Refer Slide Time: 39:12)



So, these are the results, which we have already seen from the you know, spread sheet and. And I am just bringing to highlight it this these are the results through which you know, the problem can be analyzed on or can be understand, you can understand better. And through this particular clustering, and that too through k means clustering and through the kind of the kind of hierarchical clustering.

(Refer Slide Time: 39:34)



And the corresponding to the kind of class classification, means clustering techniques like k means, clustering and hierarchical clustering.

So, there are also some of the classification techniques under the data mining; which can also be very effective, through which you can understand the problems get the insights. And do the kind of predictions and forecasting. So, what I have mentioned earlier the kind of discriminant analysis, it is like linear probability model or binary choice model.

So, where l is the kind of classifiers that discriminants the between 2 different group subjects 2 different attributes. So, like this is the kind of you know, case and then similarly so, we can have here ok. So, this is how the kind of structuring, and then you can understand the particular structure, and with the basis of this particular discriminant analysis.

(Refer Slide Time: 40:45)



And the same the kind of a credit decision; which you have discussed in the case of clustering analysis so; that means, it will give you 2 different structure through which you can understand the particular requirement. And as a result, so, we have 2 different models for yes situation, and for no situation; that means, it will it will be classified into 2 different groups with respect to the particular you know, indicators.

So, this is another kind of you know, effective tools in the data mining.

(Refer Slide Time: 41:16)



Through which you can understand the problem get the insights go for the kind of effective predictions, and the kind of you know, management decision qualitative response structure through which we can we have actually 3 different models to do the kind of you know, means to understand the kind of structure connect, the kind of model through, which you can do the better visualization.

(Refer Slide Time: 41:22)



Ah like the kind of clustering, then the kind of discriminant analysis. So, we have the kind of logistic structure, which we have which we have already discussed in a separate

lecture. And that too in the it is a so, either you can say choose, the binary choice models or logistic models. So, or the kind of probability models ruled by you know, normal distribution function, then you come with a kind of system through which you can create a kind of set up and do the better prediction.

(Refer Slide Time: 42:05)



And better kind of management decisions. So, this is a simple simpler structure of logistic regression, which we have already discussed and as again. So, these are all attributes, and then the particular dependent variable structure will be classifier through which you can actually have a different kind of you know, understanding as per the particular you know, problem; and then get the insights as per the need.

(Refer Slide Time: 42:29)



## Association Rule Mining

Association Rule Mining:

- Seeks to uncover associations in large data sets
- Association rules identify attributes that occur together frequently in a given data set.
- Market basket analysis, for example, is used determine groups of items consumers tend to purchase together.
- Association rules provide information in the form of if-then (antecedent-consequent) statements.
- The rules are probabilistic in nature.

IIT KHARAGPUR  •  NPTEL ONLINE CERTIFICATION COURSES

So, again so, in the data mining we have actually association structure through which you can understand the particular systems of course, we have already gone through you know, like covariance and you know, correlations, and these are the tools also can be used in the data mining to understand the particular you know, problems get the insights. And then create a structure through which you can do effective prediction, and the effective kind of you know, forecasting and that too as per the particular you know, problems need and the kind of management needs.

(Refer Slide Time: 43:05)



## Association Rule Mining

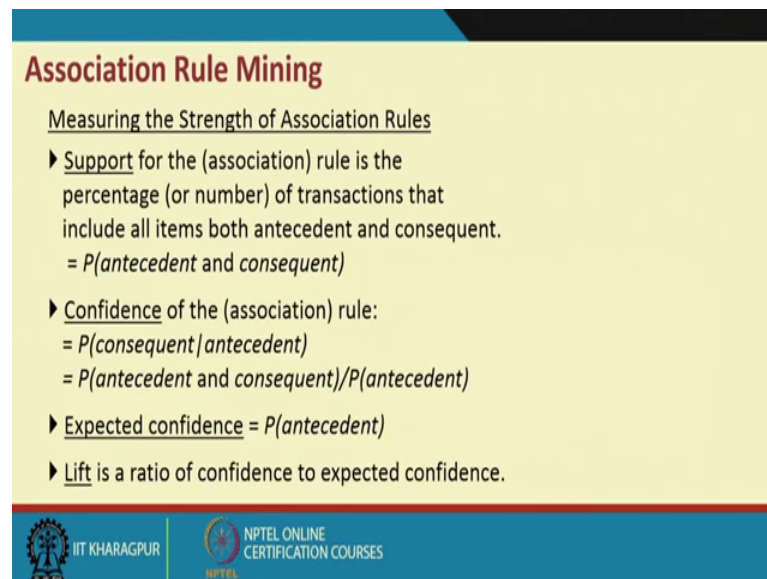Example: Custom Computer Configuration

(PC Purchase Data)

- Suppose we want to know which PC components are often ordered together.

| | A | B | C | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | PC Purchase Data | | | | | | | | | | | |
| 2 | | | | | | | | | | | | |
| 3 | | Processor | | | Screen Size | | | Memory | | | Hard Drive | |
| 4 | | | | | | | | | | | | |
| 5 | Intel Core i3 | Intel Core i5 | Intel Core i7 | 10 inch screen | 12 inch screen | 15 inch screen | 2 GB | 4 GB | 8 GB | 320 GB | 500 GB | 750 GB |
| 6 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
| 7 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| 8 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| 9 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 |
| 10 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| 11 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 |
| 12 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| 13 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
| 14 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |

IIT KHARAGPUR  •  NPTEL ONLINE CERTIFICATION COURSES

So; that means, we have association rules and this is a simple example through which you can actually develop a kind of you know, structures association structures; that means, we have done through clustering, we have done through classifications, again we can do analyze this problem through association technique. Then you develop a system through which you can do effective predictions and effective forecasting's.
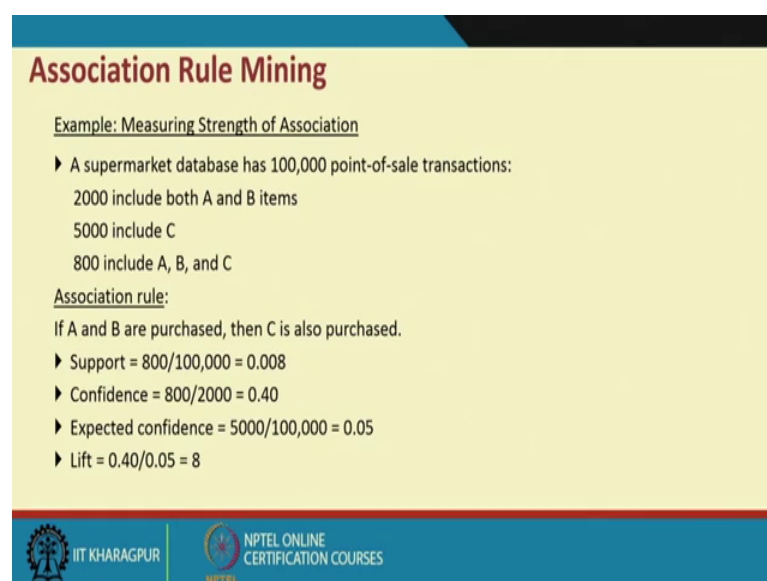
(Refer Slide Time: 43:25)



So, this is a simple example to understand against to problems.

(Refer Slide Time: 43:30)

(Refer Slide Time: 43:31)



And these are all association structures.

(Refer Slide Time: 43:32)

(Refer Slide Time: 43:33)



And again so, there is a kind of cause and effect modeling which a in fact, the logistic structures, and binary choice models itself is a cause and effect modeling. But here also there are different steps of causality models are there through which you actually you can understand the particular you know, problem and again observe the better insights. And then think about the kind of prediction structures, and the kind of management requirement.

(Refer Slide Time: 44:07)

So, that means, technically what we have mentioned we have actually plenty of you know, structure through which actually understand the problem get insights.

(Refer Slide Time: 44:09)



### Cause-and-Effect Modeling

Example (contd.): Using Correlation for Cause-and-Effect Modeling

| A | B Customer satisfaction | C Employee satisfaction | D Job satisfaction | E Satisfaction with supervisor | F Training and skill improvement |
|---|---|---|---|---|---|
| 2 Customer satisfaction | 1 | | | | |
| 3 Employee satisfaction | 0.49335 | 1 | | | |
| 4 Job satisfaction | 0.15169 | 0.84044 | 1 | | |
| 5 Satisfaction with supervisor | 0.49598 | 0.88132 | 0.60680 | 1 | |
| 6 Training and skill improvement | 0.53231 | 0.82866 | 0.71062 | 0.76970 | 1 |

Correlation analysis does not prove cause-and-effect but we can logically infer that a cause-and-effect relationship exists.
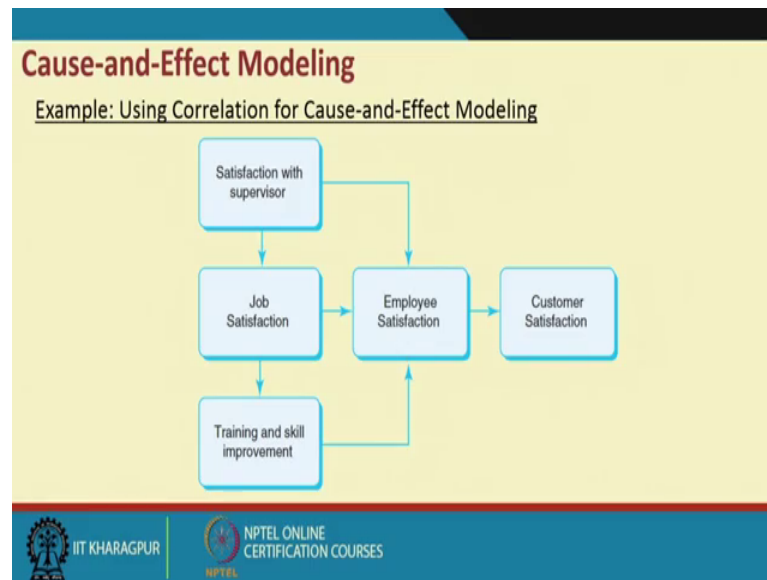
IIT KHARAGPUR    NPTEL ONLINE CERTIFICATION COURSES

Then you know, think about the final set up through which you can do the effective structure, and do the kind of prediction and forecasting as per the particular you know, business requirement and the kind of management requirement.

So, that means, actually data mining is a kind of systems, we have actually n number of you know, tools are there analytics tools are there starting with you know, simple classification clustering. Then the kind of association cause and effect models and depending upon a particular problem structures the kind of data structures then the kind of understanding, the kind of business objectives or management objectives.
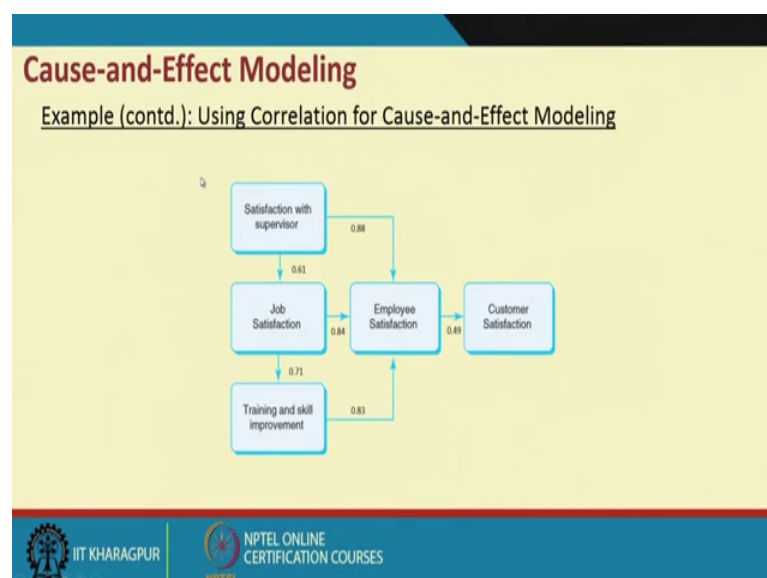
So, you pick up a particular tools and then analyze the particular problems, then get the better insights and then develop the system exactly as per the particular you know, business requirement and management for instance in the cause and effect model.

(Refer Slide Time: 45:16)



And this is how the simple example and where the customer satisfaction is the kind of you know, major effect. And then this is derived through job satisfaction training and skill improvement satisfaction with supervisor employee satisfaction. So, that means, in this particular case the most important thing is how to develop the particular you know, structure and then data will help you to verify the particular you know, structure and then we can actually validate and then use this particular model for the prediction. And you know, forecasting or as per the particular you know, management requirement again with the data.

(Refer Slide Time: 45:54)

So, we can have here the kind of you know, outcome and these are all correlation output through which you know, because we are discussing the cause and effect and association through which you can do the kind of better prediction. Now with the data, what we have observed here the correlation coefficients are very high as a result. So, the particular model's you know, for this particular problem is very effective for the kind of you know, management requirement.

So, likewise we can actually think about the particular problems, check the data structure you know, pick up a particular technique, and then develop the system, which can be very effective as per the particular you know, problem requirement or you know, management requirement. With this we will stop here.

Thank you very much; have a nice day.