

**Introduction to Internet of Things**  
**Prof. Sudip Misra**  
**Department of Computer Science & Engineering**  
**Indian Institute of Technology, Kharagpur**

**Lecture - 56**  
**Data Handling and Analytics- Part- II**

In this lecture on data analysis and sorry Data Handling and Analytics. In the first part we will focused mostly on data handling, and in the second part we are going to focus mostly on the analytics.

So, having captured the data and storing it in the cloud or in the server or whatever storage mechanisms we have. We now have to use the data, for using it we have to analyze it, we have to analyze the data. So, for this there are different tools, different methodologies that they are; the most common the most primitive once are based on statistical methods; so basic statistical methods can be applied applied on the store data in what order to make more sense out of that data in order to get more insight into that data it is stored.

(Refer Slide Time: 01:20)

**What is Data Analytics**

✓ *"Data analytics (DA) is the process of examining data sets in order to draw conclusions about the information they contain, increasingly with the aid of specialized systems and software. Data analytics technologies and techniques are widely used in commercial industries to enable organizations to make more-informed business decisions and by scientists and researchers to verify or disprove scientific models, theories and hypotheses."*

[An admin's guide to AWS data management]

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

So, we have data analytics. Basically, the data have to be analyzed. So, in the context of IoT people talk a lot about data analytics. So, what is the state analytics? I am going to read one of these definitions. So, data analytics is the process of examining the data sets in order to draw conclusions about the information they contain. That means, what information is contained in this data sets. Increasingly with the aid of specialized systems and software; so

with different specialized software, systems, etcetera to get insight into the data that is existing.

Data analytics technologies and techniques are widely used in commercial industries to enable organizations to make more informed business decisions and by scientist and researchers to verify or disprove the scientific models theories and hypothesis. So, this is basically the premise in which data analytics basically work and the different concerns of data analytics are basically: I mentioned in this particular definition.

(Refer Slide Time: 02:39)

The slide is titled "Types of Data Analysis" in a bold, dark red font. Below the title, there is a bulleted list with checkmarks. The first bullet point is "Two types of analysis". The second bullet point is "Qualitative Analysis", which has a sub-bullet: "Deals with the analysis of data that is categorical in nature". The third bullet point is "Quantitative Analysis", which has a sub-bullet: "Quantitative analysis refers to the process by which numerical data is analyzed". At the bottom of the slide, there is a blue banner with the IIT Kharagpur logo on the left, the NPTEL ONLINE CERTIFICATION COURSES logo in the center, and the text "Introduction to Internet of Things" on the right. A small circular inset image of a man in a white shirt is visible in the bottom right corner of the slide.

- ✓ Two types of analysis
  - ✓ Qualitative Analysis
    - ✓ Deals with the analysis of data that is categorical in nature
  - ✓ Quantitative Analysis
    - ✓ Quantitative analysis refers to the process by which numerical data is analyzed

So, when we talk about analysis in general broadly analysis comes in two forms. So, we can either perform qualitative analysis on the data that has been obtained or we can perform quantitative analysis. So, qualitative analysis basically deals with the analysis of data that are categorical in nature- so qualitative analysis. Whereas, quantitative analysis refers to the process by which numerical methods can be used; numerical data can be analyzed through quantitative analysis.


So, categorical data: qualitative analysis is good enough, for numerical data quantitative analysis quantitative methods are useful.

(Refer Slide Time: 03:31)

**Qualitative Analysis**

- ✓ Data is not described through numerical values
- ✓ Described by some sort of descriptive context such as text
- ✓ Data can be gathered by many methods such as interviews, videos and audio recordings, field notes
- ✓ Data needs to be interpreted
- ✓ The grouping of data into identifiable themes
- ✓ Qualitative analysis can be summarized by three basic principles (Seidel, 1998):
  - ✓ Notice things
  - ✓ Collect things
  - ✓ Think about things

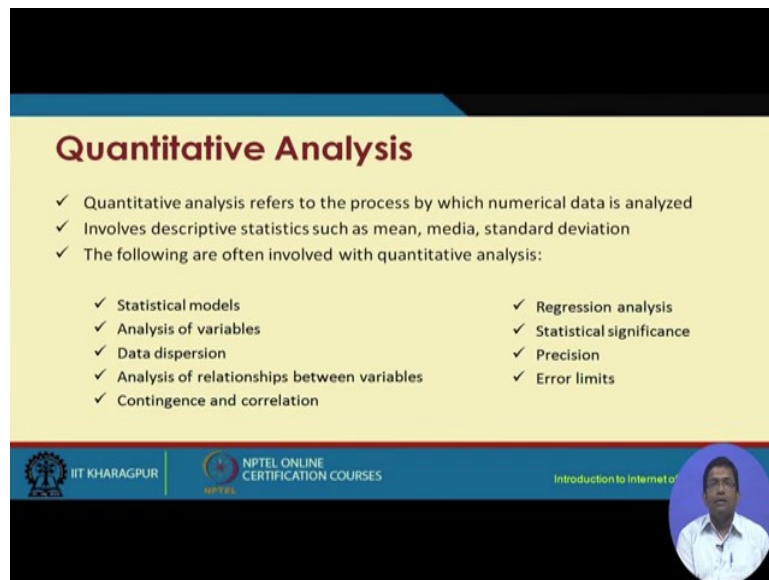
IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things



So, qualitative analysis data is not described through numerical values, but are described by some sort of descriptive contexts such as text. This qualitative data can be gathered by many methods such as from interviews, interviewing different people, from videos, from audio recordings, field notes you know industry manuals and so on.

This data needs to be interpreted; the grouping of the data can be where should be performed into identifiable themes in quantitative qualitative analysis. And the qualitative analysis can be summarized by three basic principles. Notice the things collect the things and think about it. We do not need to get into details of each of these.

(Refer Slide Time: 04:27)



## Quantitative Analysis

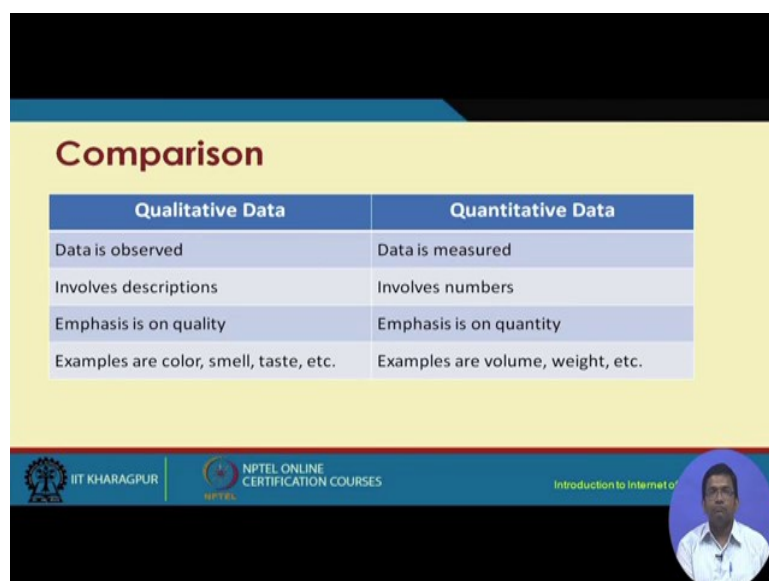
- ✓ Quantitative analysis refers to the process by which numerical data is analyzed
- ✓ Involves descriptive statistics such as mean, media, standard deviation
- ✓ The following are often involved with quantitative analysis:
  - ✓ Statistical models
  - ✓ Analysis of variables
  - ✓ Data dispersion
  - ✓ Analysis of relationships between variables
  - ✓ Contingence and correlation
  - ✓ Regression analysis
  - ✓ Statistical significance
  - ✓ Precision
  - ✓ Error limits

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

NPTEL

Next one is the quantitative analysis. So, quantitative analysis is on the numeric data using different statistical methods such as descriptive statistics, more specifically finding out the mean of the dataset median standard deviation and so on. The following are often involved with the quantitative analysis, statistical models and analysis of variance, then data dispersion analysis of relationship between variables, contingency and correlation, then regression analysis, statistical significance, precision, error limits and so on.

(Refer Slide Time: 05:18)



## Comparison

Qualitative Data	Quantitative Data
Data is observed	Data is measured
Involves descriptions	Involves numbers
Emphasis is on quality	Emphasis is on quantity
Examples are color, smell, taste, etc.	Examples are volume, weight, etc.

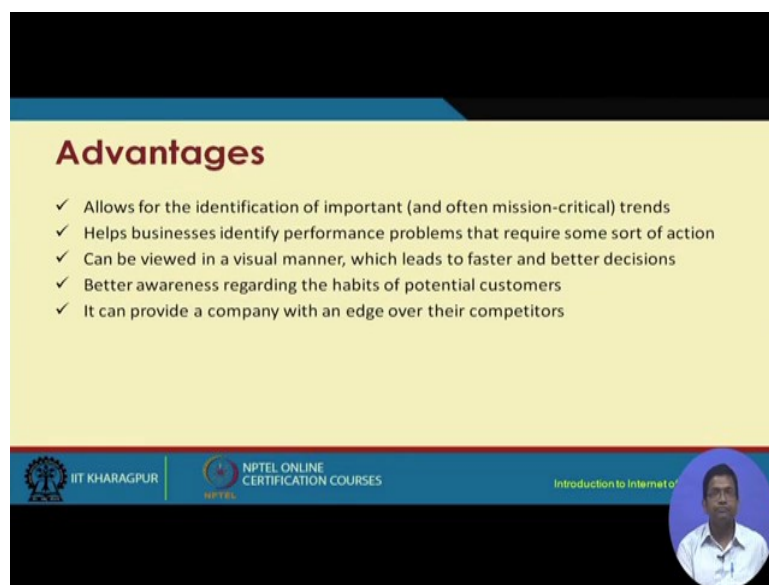
IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

NPTEL

So, these are the different quantitative methods that are used for quantitative analysis of data.

Now comes the comparison between qualitative data and quantitative data. So, qualitative data can be mostly observed, whereas quantitative data can be measured. Qualitative data involves descriptions its more qualitative it involves descriptions, on the other hand quantitative data involves numbers, numeric's and so on. Whereas, in qualitative data the emphasis is on quality, in the quantitative data the emphasis is on quantity. Examples of qualitative data include colour, smell, taste, etcetera which cannot be quantified so easily. On the other hand quantifiable data include volume, weight, etcetera; these are numbers these are figures which can be used to perform different numerics.

(Refer Slide Time: 06:14)



The slide is titled "Advantages" in a bold, dark red font. It lists five advantages, each preceded by a checkmark. The background is a light yellow. At the bottom, there is a blue banner with logos for IIT Kharagpur, NPTEL, and the course title "Introduction to Internet of Things". A small circular inset shows a man in a white shirt.

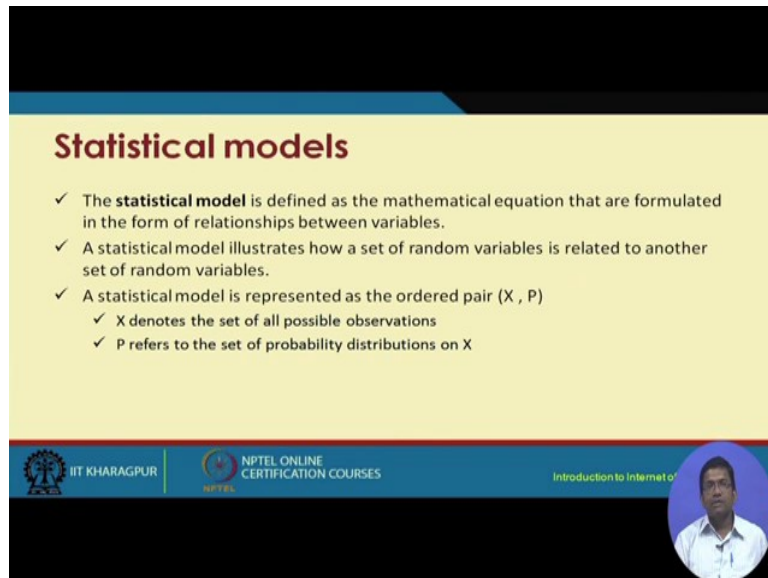
- ✓ Allows for the identification of important (and often mission-critical) trends
- ✓ Helps businesses identify performance problems that require some sort of action
- ✓ Can be viewed in a visual manner, which leads to faster and better decisions
- ✓ Better awareness regarding the habits of potential customers
- ✓ It can provide a company with an edge over their competitors

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

The advantages of data analytics is that: it allows for the identification of important trends, it helps the businesses identify performance problems that we require some sort of action- some prediction can be performed. By analysis of the data fast data something you know, so the businesses they can understand that what has gone wrong quantitatively we can be done or even qualitatively. So, data analytics are useful for that.

So, the analytics can also be performed in a visual manner and that can help in faster and better decision making. Analytics can provide a company with an edge over their competitors.

(Refer Slide Time: 07:00)



### Statistical models

- ✓ The **statistical model** is defined as the mathematical equation that are formulated in the form of relationships between variables.
- ✓ A statistical model illustrates how a set of random variables is related to another set of random variables.
- ✓ A statistical model is represented as the ordered pair  $(X, P)$ 
  - ✓  $X$  denotes the set of all possible observations
  - ✓  $P$  refers to the set of probability distributions on  $X$

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

So, that is the reason actually data analytics has become very popular in the industry, not only in the industrial almost in all spheres of life data analytics has become very popular. And because you can get more insight into what is going on in the processes that are occurring around you.

Statistical different models of statistical; statistical models can be adopted in order to perform quantitative analysis. And a statistical model can is defined as the mathematical equation that is formulated to form the relationship between variables. A statistical model illustrates how a set of random variables is related to another set of random variables. And it is a statistical model is represented as an ordered pair  $X P$ ; where  $X$  denotes the set of all possible observations and  $P$  refers to the set of probability of distributions on  $X$ .

(Refer Slide Time: 08:05)

**Statistical models (Contd.)**

- ✓ Statistical models are broadly categorized as
  - ✓ Complete models
  - ✓ Incomplete models
- ✓ Complete model does have the number of variables equal to the number of equations
- ✓ An incomplete model does not have the same number of variables as the number of equations

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

Statistical models are broadly categorized as complete models and incomplete models. Complete models have the same number of variables as the number of equations. So, the number of equations and the number of variables in the complete models are the same. So, if we have the number of variables equating equating with the number of equations what we have is a complete model. And in an incomplete model the number of variables and the number of equations are not the same- they do not match.

(Refer Slide Time: 08:39)

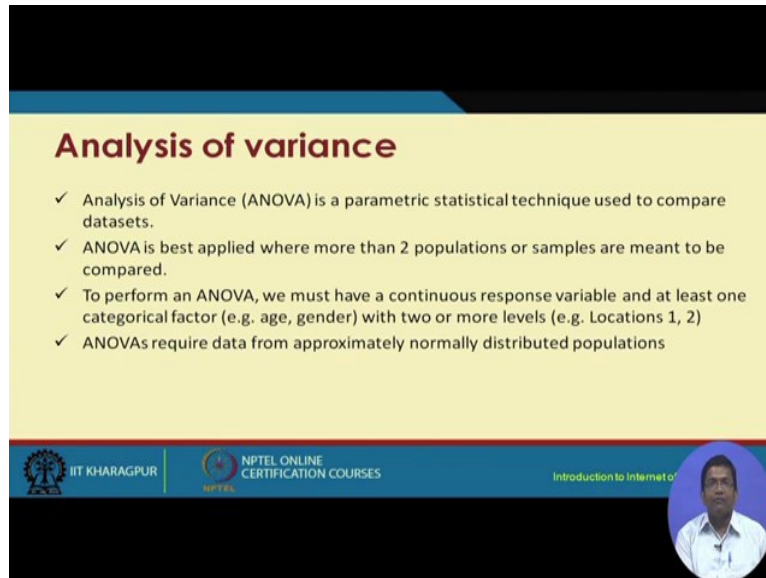
**Statistical models (Contd.)**

- ✓ In order to build a statistical model
  - ✓ Data Gathering
  - ✓ Descriptive Methods
  - ✓ Thinking about Predictors
  - ✓ Building of model
  - ✓ Interpreting the Results

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

So, in order to build a statistical model it is required to gather the data, perform descriptive methods, think about what are the predictors, then build the model and then interpret the results.


(Refer Slide Time: 08:55)



**Analysis of variance**

- ✓ Analysis of Variance (ANOVA) is a parametric statistical technique used to compare datasets.
- ✓ ANOVA is best applied where more than 2 populations or samples are meant to be compared.
- ✓ To perform an ANOVA, we must have a continuous response variable and at least one categorical factor (e.g. age, gender) with two or more levels (e.g. Locations 1, 2)
- ✓ ANOVAs require data from approximately normally distributed populations

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things



Analysis of variance in short known as ANOVA analysis: is a parametric statistical technique that can be used to compare two data sets- two or more data sets they can be compared. So, ANOVA is best applied when more than two populations of samples are meant to be compared. So, we have one dataset, we have another dataset, we want to compare these two populations to see that what is the you know how much is the correlation between these two datasets, what sort of similarity exists between these two database sets.

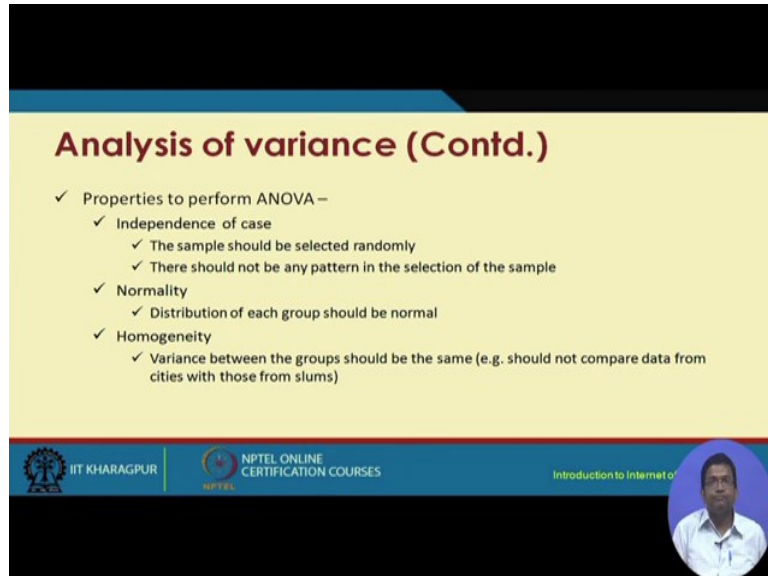
So, to perform ANOVA one has to have a continuous response variable and at least one categorical factor. For example, age gender, etcetera, with at least two or more levels example location 1, location 2, etcetera. So, what it means is basically levels mean that one location: one location Kharagpur another location Kolkata. So, these are two different locations corresponding to two different levels. And categories mean the age is one category. So, with respect to a particular category like age you know. So, at two different locations what is the similarity or what is the dissimilarity; similarly with respect to gender or any other category.

ANOVA requires data from approximately normally distributed population. So, this is a very important assumption or a very important requirement that you know. So, normal distribution



is required for performing; normal distribution of the data set has to be there in order to perform ANOVA analysis.


(Refer Slide Time: 10:42)



**Analysis of variance (Contd.)**

- ✓ Properties to perform ANOVA –
  - ✓ Independence of case
    - ✓ The sample should be selected randomly
    - ✓ There should not be any pattern in the selection of the sample
  - ✓ Normality
    - ✓ Distribution of each group should be normal
  - ✓ Homogeneity
    - ✓ Variance between the groups should be the same (e.g. should not compare data from cities with those from slums)

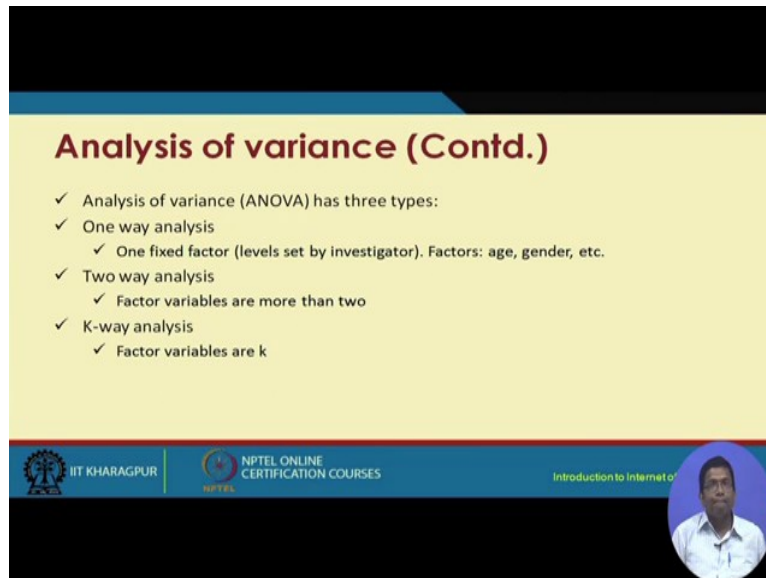
IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things



The properties to perform ANOVA: one is the independence of case, the sample that is selected should be random; random is selected there should not be any bias, there should not be any pattern, in the selected sample. Normality is the second property which constant the distribution of each group should be normal, so normal distribution of the data within the group; and homogeneity which constants variance between the groups and this should be the variance should be the same. So, we should not have a scenario to compare the data from cities with the data from maybe slums areas or maybe the data of Kharagpur compared with the data of Kolkata.

So, because we have a town we have a city so two different datasets you know compare of comparing with each other. So, they have huge variance. And the variance should be as much minimal as possible.

(Refer Slide Time: 11:50)



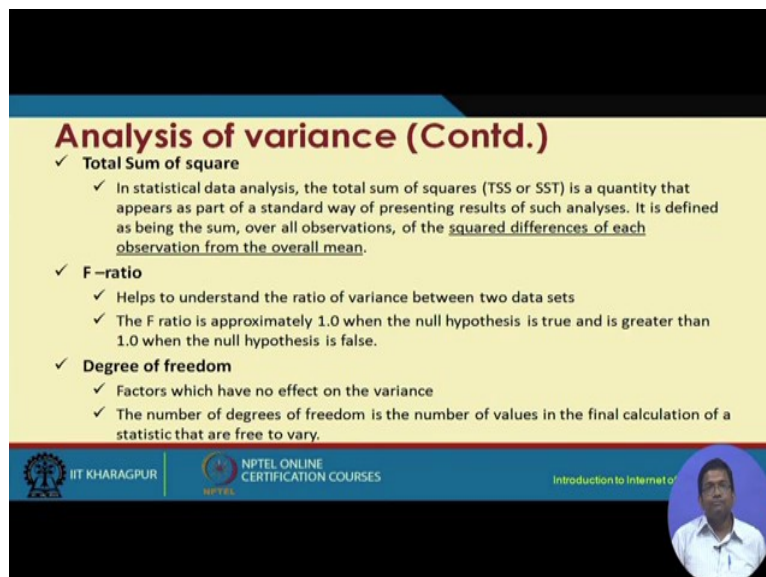
**Analysis of variance (Contd.)**

- ✓ Analysis of variance (ANOVA) has three types:
- ✓ One way analysis
  - ✓ One fixed factor (levels set by investigator). Factors: age, gender, etc.
- ✓ Two way analysis
  - ✓ Factor variables are more than two
- ✓ K-way analysis
  - ✓ Factor variables are k

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

Analysis of variance has three different types: one way analysis which constant one fixed factor. For example, the factors could be age, gender, etcetera. Could be two way analysis where two or more or two factors are going to be involved. So, both maybe both age and gender will be considered in a two way factor two way ANOVA analysis. And it can be k way analysis where k factor variables are involved.

(Refer Slide Time: 12:21)



**Analysis of variance (Contd.)**

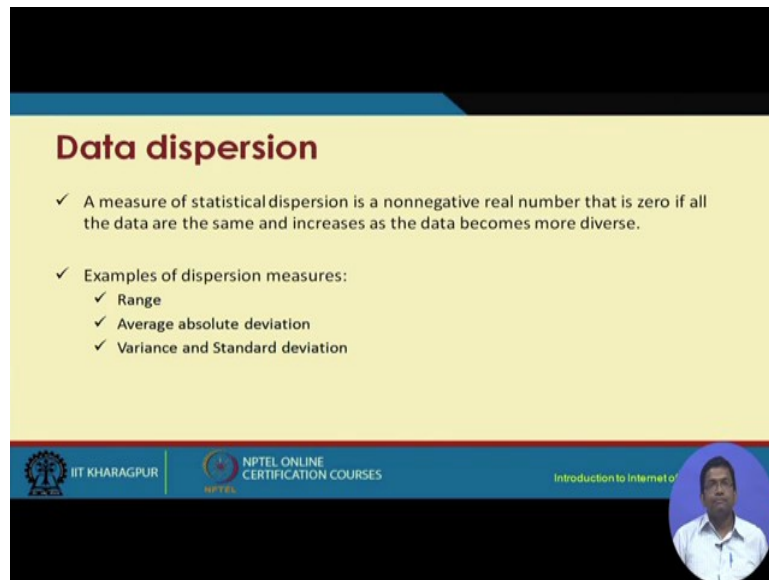
- ✓ **Total Sum of square**
  - ✓ In statistical data analysis, the total sum of squares (TSS or SST) is a quantity that appears as part of a standard way of presenting results of such analyses. It is defined as being the sum, over all observations, of the squared differences of each observation from the overall mean.
- ✓ **F-ratio**
  - ✓ Helps to understand the ratio of variance between two data sets
  - ✓ The F ratio is approximately 1.0 when the null hypothesis is true and is greater than 1.0 when the null hypothesis is false.
- ✓ **Degree of freedom**
  - ✓ Factors which have no effect on the variance
  - ✓ The number of degrees of freedom is the number of values in the final calculation of a statistic that are free to vary.

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

Then there are different ways, different features that are there for performing analysis of variance- total sum of square is 1, f ratio is another and the degree of freedom. So, all of these


things have to be taken into consideration in order to perform and its variance. I am not going to put through these, because this is not essentially a course in statistical methods. And these are all available, but what is important is that ANOVA analysis can be used in order to perform analytics on the data that is obtained from IoT systems.

(Refer Slide Time: 13:01)



**Data dispersion**

- ✓ A measure of statistical dispersion is a nonnegative real number that is zero if all the data are the same and increases as the data becomes more diverse.
- ✓ Examples of dispersion measures:
  - ✓ Range
  - ✓ Average absolute deviation
  - ✓ Variance and Standard deviation

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things | 

The next concept that has to be understood is data dispersion. Data dispersion concerns how much is the dispersion; that means, dispersion is basically a measure of the statistical. So, it is a measure of statistical dispersion a non negative real number that is 0 if all the data are the same and it increases as the data becomes more diverse. Examples of dispersion measures include: range, average, absolute deviation, variance, and standard deviation. So, typically when we talk about dispersion we typically talked about in terms of variance and standard deviation.

(Refer Slide Time: 13:58)

**Data dispersion (Contd.)**

- ✓ **Range**
  - ✓ The range is calculated by simply taking the difference between the maximum and minimum values in the data set.
- ✓ **Average absolute deviation**
  - ✓ The average absolute deviation (or mean absolute deviation) of a data set is the average of the absolute deviations from the mean.
- ✓ **Variance**
  - ✓ Variance is the expectation of the squared deviation of a random variable from its mean
- ✓ **Standard deviation**
  - ✓ Standard deviation (SD) is a measure that is used to quantify the amount of variation or dispersion of a set of data values

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things |

So, how much the deviate from the norm from the norm? So, this is what dispersion constants. So, here as I already mentioned so range and what is meant by it the absolute standard deviation is given. So, the average of the absolute deviation is given, variance and standard deviation here well known methods of dispersion deciding. And these are given over here.

(Refer Slide Time: 14:19)

**Contingence and correlation**

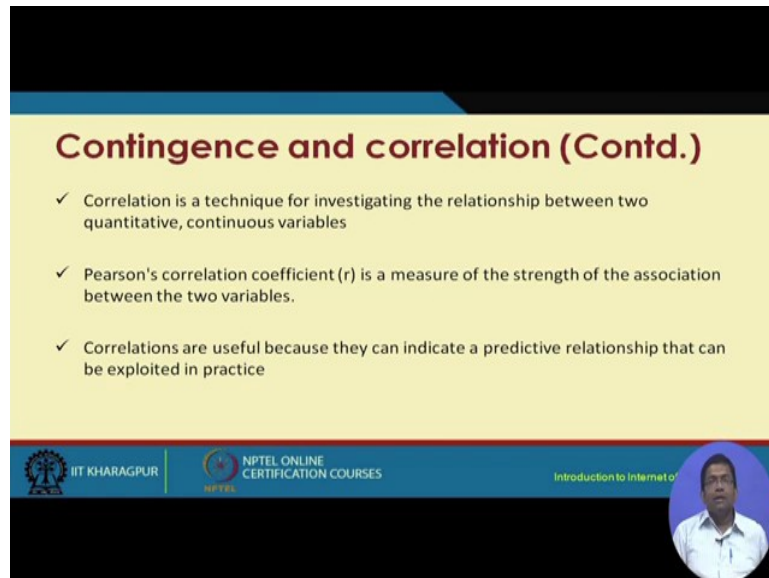
- ✓ In statistics, a contingency table (also known as a cross tabulation or crosstab) is a type of table in a matrix format that displays the (multivariate) frequency distribution of the variables.
- ✓ Provides a basic picture of the interrelation between two variables
- ✓ A crucial problem of multivariate statistics is finding (direct-)dependence structure underlying the variables contained in high-dimensional contingency tables

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things |

Next comes a contingency and correlation. So, in statistics a contingency table is a type of table in a matrix format that displays the multivariate frequency distribution of the variables.

It provides the basic picture of the interrelation between two variables. Correlation is a technique for investigating the relationship between two continuative continuous variables.

(Refer Slide Time: 14:42)



**Contingence and correlation (Contd.)**

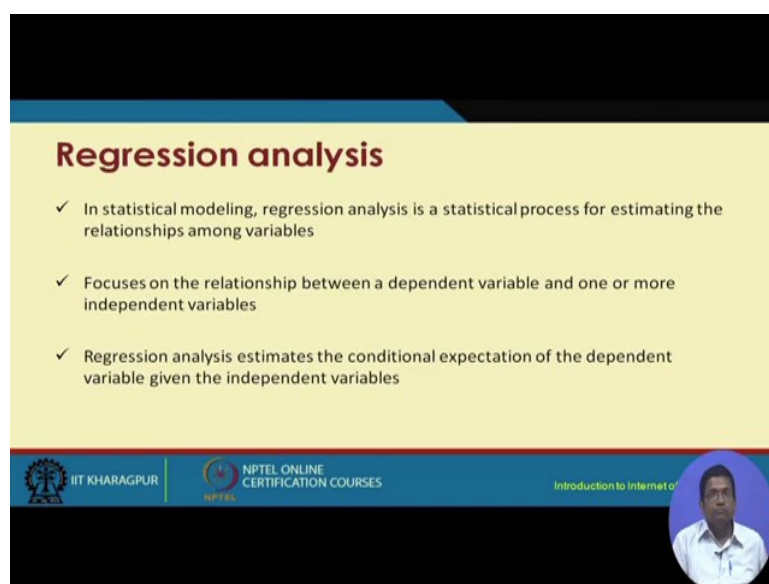
- ✓ Correlation is a technique for investigating the relationship between two quantitative, continuous variables
- ✓ Pearson's correlation coefficient ( $r$ ) is a measure of the strength of the association between the two variables.
- ✓ Correlations are useful because they can indicate a predictive relationship that can be exploited in practice

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

NPTEL

So, they have to be continuous variable this is very important. And how much they are correlated these two variables how much they are correlated and what is the relationship between them. So, a popular measure is the Pearson's correlation coefficient. And it basically measures the strength of association between two variables.

(Refer Slide Time: 15:13)



**Regression analysis**

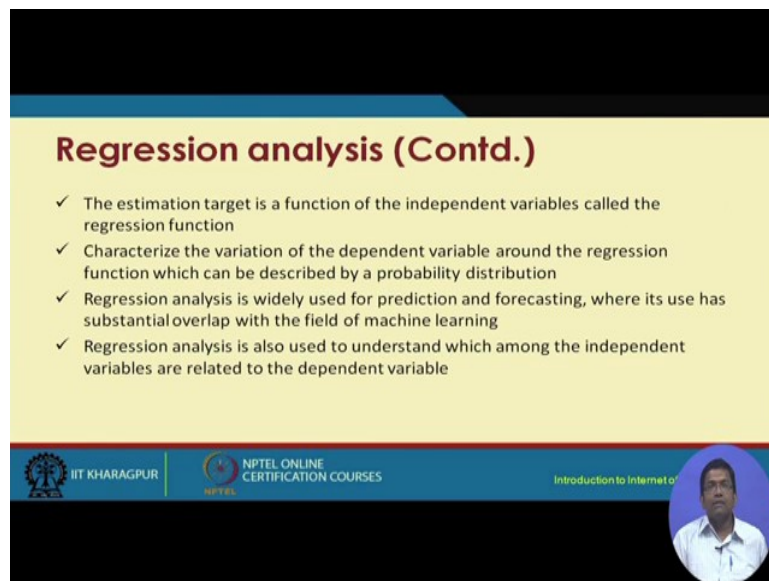
- ✓ In statistical modeling, regression analysis is a statistical process for estimating the relationships among variables
- ✓ Focuses on the relationship between a dependent variable and one or more independent variables
- ✓ Regression analysis estimates the conditional expectation of the dependent variable given the independent variables

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

NPTEL

Then comes regression analysis. So, regression analysis basically tries to estimate the relationship among the different variables. It focuses on the relationship between independent variable and one or more independent variables. So, we have a dependent variable and we have an independent; we have one or more independent variables and how the dependent variable relates to one or more of these variables taken at a time or taken together.


(Refer Slide Time: 15:47)



**Regression analysis (Contd.)**

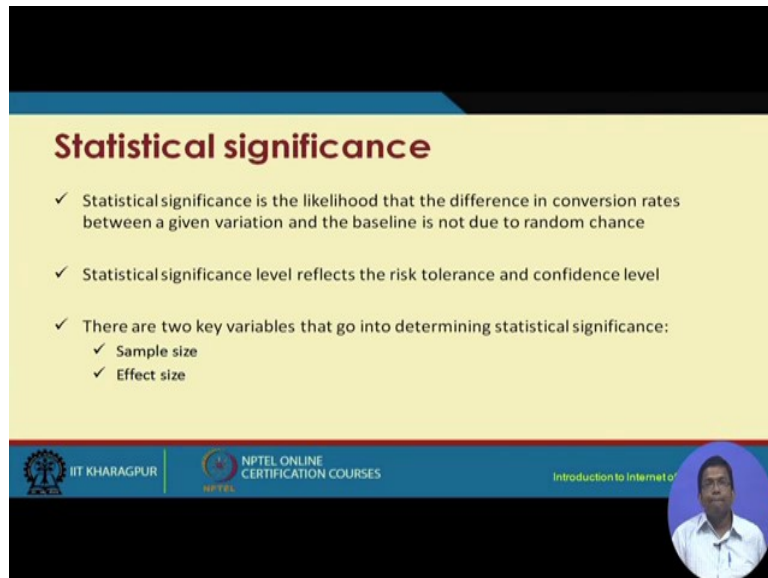
- ✓ The estimation target is a function of the independent variables called the regression function
- ✓ Characterize the variation of the dependent variable around the regression function which can be described by a probability distribution
- ✓ Regression analysis is widely used for prediction and forecasting, where its use has substantial overlap with the field of machine learning
- ✓ Regression analysis is also used to understand which among the independent variables are related to the dependent variable

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things



So, regression analysis estimates the conditional expectation of the dependent variable given the independent variables. The estimation target is a function of the independent variables called the regression function. It characterizes the variation of the independent variable around the regression function which can be described by a probability distribution. So, regression analysis is helpful in different ways: it can be used to understand how the independent variables are related to the dependent variable one at a time or taken together.


(Refer Slide Time: 16:16)



**Statistical significance**

- ✓ Statistical significance is the likelihood that the difference in conversion rates between a given variation and the baseline is not due to random chance
- ✓ Statistical significance level reflects the risk tolerance and confidence level
- ✓ There are two key variables that go into determining statistical significance:
  - ✓ Sample size
  - ✓ Effect size

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

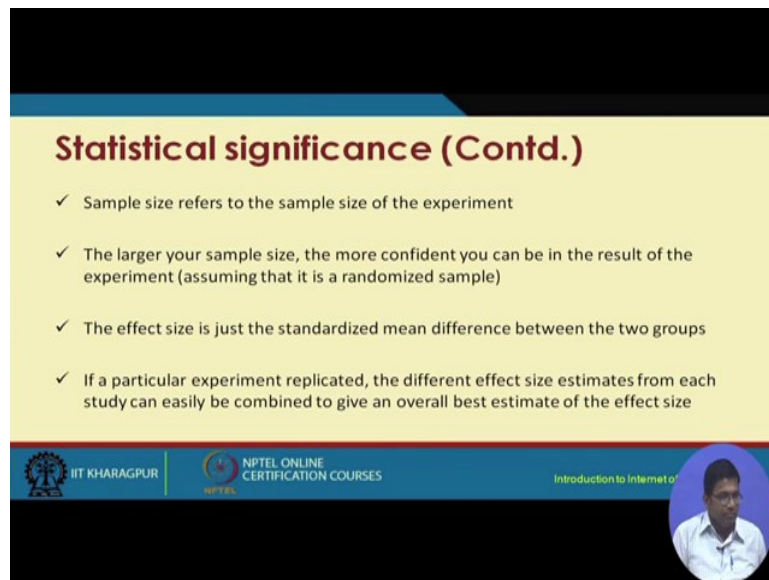


Statistical significance is important. It basically measures the likelihood that the difference in conversion rates between given variation and the baseline is not due to any random chance. So, statistically how much you know the results are significant is something that has to be measured. So, statistical significance basically reflects the risk, tolerance and the confidence level. So, how much is the confidence on the results that are obtained.

So, this is measure through statistical significance. So, typically there are two key variables that are required for determining statistical significance: one is the sample size the other one is the effect size.




(Refer Slide Time: 17:07)



### Statistical significance (Contd.)

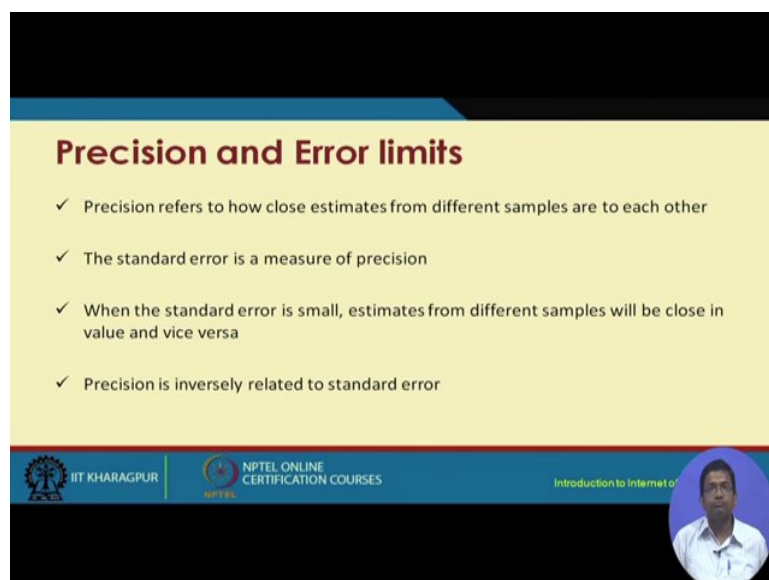
- ✓ Sample size refers to the sample size of the experiment
- ✓ The larger your sample size, the more confident you can be in the result of the experiment (assuming that it is a randomized sample)
- ✓ The effect size is just the standardized mean difference between the two groups
- ✓ If a particular experiment replicated, the different effect size estimates from each study can easily be combined to give an overall best estimate of the effect size

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things



The sample size refers to the sample size of the experiment. The larger the sample size is the more confident one can be on the result of the experiment. And the effect size is just the standardized mean difference between the two groups. So, if a particular experiment is replicated the different effect size estimates from each study can easily be combined to give an overall best estimate of the effect size.


(Refer Slide Time: 17:38)



### Precision and Error limits

- ✓ Precision refers to how close estimates from different samples are to each other
- ✓ The standard error is a measure of precision
- ✓ When the standard error is small, estimates from different samples will be close in value and vice versa
- ✓ Precision is inversely related to standard error

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

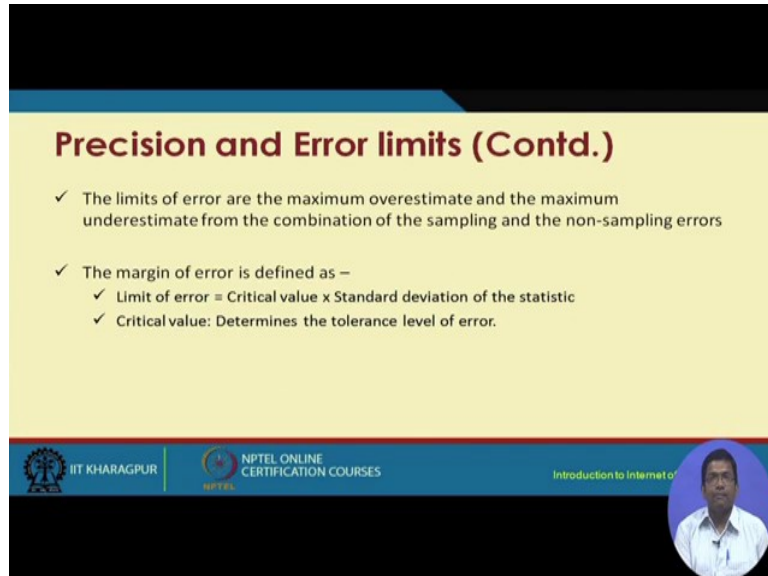


Precision and error limits are important. So, precision basically concerns how close the estimates are from the different samples to each other. The standard error is a measure of the



precision, when the standard error is small the estimates from the different samples will be closed in value and vice versa. So, precision is inversely related to the standard error.

(Refer Slide Time: 18:03)



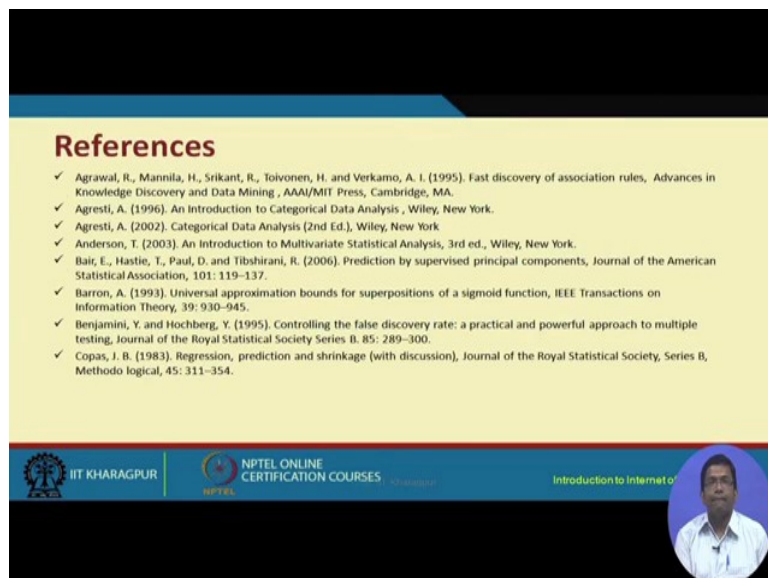
**Precision and Error limits (Contd.)**

- ✓ The limits of error are the maximum overestimate and the maximum underestimate from the combination of the sampling and the non-sampling errors
- ✓ The margin of error is defined as –
  - ✓ Limit of error = Critical value x Standard deviation of the statistic
  - ✓ Critical value: Determines the tolerance level of error.

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

So, this precision and error become hand in hand, the limits of the error and the overestimate and the underestimate are taken into consideration while considering the error limits.

(Refer Slide Time: 18:19)



**References**

- ✓ Agrawal, R., Mannila, H., Srikant, R., Toivonen, H. and Verkamo, A. I. (1995). Fast discovery of association rules, *Advances in Knowledge Discovery and Data Mining*, AAAI/MIT Press, Cambridge, MA.
- ✓ Agresti, A. (1996). *An Introduction to Categorical Data Analysis*, Wiley, New York.
- ✓ Agresti, A. (2002). *Categorical Data Analysis* (2nd Ed.), Wiley, New York
- ✓ Anderson, T. (2003). *An Introduction to Multivariate Statistical Analysis*, 3rd ed., Wiley, New York.
- ✓ Bair, E., Hastie, T., Paul, D. and Tibshirani, R. (2006). Prediction by supervised principal components, *Journal of the American Statistical Association*, 101: 119–137.
- ✓ Barron, A. (1993). Universal approximation bounds for superpositions of a sigmoid function, *IEEE Transactions on Information Theory*, 39: 930–945.
- ✓ Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing, *Journal of the Royal Statistical Society Series B*, 57: 289–300.
- ✓ Copas, J. B. (1983). Regression, prediction and shrinkage (with discussion), *Journal of the Royal Statistical Society, Series B, Methodological*, 45: 311–354.

IIT KHARAGPUR | NPTEL ONLINE CERTIFICATION COURSES | Introduction to Internet of Things

So, as I was mentioning at the outset that there are different statistical tools that additional statistical tools like correlation analysis, regression analysis, analysis of variance can be used

in order to understand to how to get insight on the data that is obtained that is collected. But these are the basic analytic methods.

And what we have not discussed over here and we have intentionally confined ourselves to not discussing things like how text can be how text can be analyzed textual data or how video data can be analyzed and so on. So, that requires video data images can be analyzed and so on; so different other types of data can be analyzed. So, this we have intentionally not discussed, because that requires specialized training in text processing, video processing, image processing, and so on. And we do not want to get into the depth of those types of analytics.

So, these are the differences. So, with this we come to an end. And as I was telling mentioning before data handling, data analytics are very crucial in the context of IoT because lot of data gets generated in the IoT domain. And this data not only have to be analyze, but prior to analyzing they have to be handled. They have to be handled using technologies such as cloud, we have to be handled with technology such as Hadoop and so on.

And once they are handled that means, the data have been stored, they have been cleaned and stored and so on then they have to be analyzed. For analysis we have these different statistical methods, we have different other methods based in machine learning, image processing, video processing, text processing and so on. So, those are advance methods which we do not cover in this particular lecture, in this particular course.

Thank you.