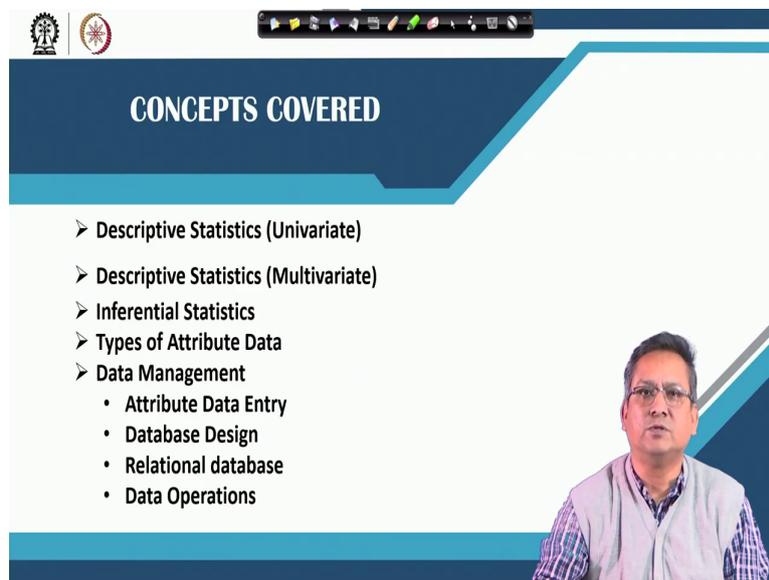


Geo Spatial Analysis in Urban Planning
Prof. Saikat Kumar Paul
Department of Architecture and Regional Planning
Indian Institute of Technology, Kharagpur

Lecture – 11
Attribute Data Management and Data Exploration

Welcome back dear students. We are into the module 3 now and we shall talk about Attribute Data Management and Data Exploration in this particular lecture.

(Refer Slide Time: 00:27)



The slide features a dark blue header with the text 'CONCEPTS COVERED' in white. Below the header is a list of topics, each preceded by a right-pointing arrowhead. The topics are: Descriptive Statistics (Univariate), Descriptive Statistics (Multivariate), Inferential Statistics, Types of Attribute Data, and Data Management. Under 'Data Management', there are four sub-topics: Attribute Data Entry, Database Design, Relational database, and Data Operations. In the bottom right corner of the slide, there is a small video inset showing a man with glasses and a light-colored shirt, presumably the professor, speaking.

- Descriptive Statistics (Univariate)
- Descriptive Statistics (Multivariate)
- Inferential Statistics
- Types of Attribute Data
- Data Management
 - Attribute Data Entry
 - Database Design
 - Relational database
 - Data Operations

So, the concepts that we are going to cover today is about Descriptive Statistics and we are going to talk about the Univariate Statistics, we are also going to talk about touch upon Multivariate Statistics, we would see what is Inferential Statistics, we shall also look into the types of Data Attributes and how the data is managed in GIS.

So, how we do the Attribute data entry, how we do the Database design, what is a Relational database and the I mean different types of data operations regarding joining of the tables. So, we shall look into these following concepts.

(Refer Slide Time: 01:09)

Descriptive Statistics (Univariate)

- Cannot be used for analyzing nominal or categorical variables
- Summarize observations – with reference to the distribution of a variable –
- histogram of classes - grouped values – showing frequency of occurrence
- Number of classes - determined as a function of number of observations and range of values
- Mean - measures of central tendency in a distribution

$$\mu = \frac{1}{n} \sum_{i=1}^n z_i$$

- Median - middle value when all values are ordered from smallest to largest
- Mode - most frequently occurring value
- Mean is sensitive to outliers, median or mode reduces impact of outliers

NPTEL Online Certification Program
IIT Kharagpur

So, talking about Descriptive or Univariate statistics, it is used I mean for the ordinal variables and the ratio variables and it cannot be used for nominal or categorical variables. We shall see what are these nominal and categorical variables. Some of you might be knowing about these who are into I mean data management and I mean database design and all these things.

So, it basically summarises the observation with descriptive statistics, it gives a summary of the observations and it is with reference to the distribution of a variable. So, whenever we have variables and we have distributed data or discrete data, it is lumped or grouped into

different classes or grouped values and we have frequencies of occurrence for each of these groups which can be shown as histogram of the different classes.

Now, I mean the number of classes that we would categorise our data is determined as a function of the number of observation and the range of the values. Now, talking about the measures of central tendencies, we have three measures which you already know about. The first one is known as Mean which is I mean the average of all the sets of observation.

So, if we have n number of observations, we sum up the I mean observations all the “ n ” observations and divide it by the number of observation to give us the mean which is denoted as μ . Now the next measure of central tendency is the Median which is the median middle value when all values are ordered if we I mean stack up the values from the smallest to the largest.

So, you can find out the middle value and in case if you have even number of data sets, in that case what we do is we take the two middle values and I mean take mean of those. Now, you have different types of algorithms in computer I mean software analysis. So, we can I mean categories the data or I mean save the data in a I mean smallest to largest I mean ordering. So, those different algorithms are available where in if we have a very large data set, it can be ordered and then we can find out the median.

Now, mode is the most frequently occurring value in a data set. So, whatever data value is the most frequently occurring data value I mean that we take it as a mode. Now, if we have outliers in the data, if we have some data values which are extreme data values I mean which is beyond the I mean central any of the central tendencies and which has very high values, so in those cases what happens this outliers will add up significantly to the mean. So, in a way the mean gets I mean affected by this kind of outliers. So, mean median or mode value it reduces the impact of the outliers.

(Refer Slide Time: 04:52)

Descriptive Statistics (Univariate)

- **Standard deviation** – degree of dispersion around the mean
$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (z_i - \mu)^2}$$
- **Variance** - expectation of squared deviation of a random variable from its mean - how far a set of (random) numbers are spread out from their average value
- **Coefficient of skewness** - measures lack of symmetry in data distribution – measured using Pearson's moment coefficient of skewness
$$\text{skewness} = \frac{n}{(n-1)(n-2)} \sum_{i=1}^n \left(\frac{z_i - \bar{z}}{s} \right)^3$$

NPTEL Online Certification Course
IIT Kharagpur

So, I mean we talk about the descriptive statistics. So, we have three measures basically we talk about the standard deviation where in we have the degree of the dispersion around the mean and it is represented as sigma which is the difference of the square I mean of your observed values to the mean are divided by the total number of observations and it is taken as a square root.

We also have a term which is known as Variance which is I mean extensively used in GIS or otherwise in descriptive statistics. So, I mean it is the expectation of square deviation when the variable is random with respect to its mean, now I mean it gives us a spread out of those random numbers with respect to the average value with respect to the mean values. Talking about the coefficient of skewness when we are having a distribution a, when we are having a histogram we can see that it could be in the form of a inverse bell curve.

So, this curve could be symmetrical in nature or it could be asymmetrical in nature. So, the coefficient of skewness, it measures the lack of symmetry or the presence of symmetry in the data distribution and it is measured using Pearson's moment coefficient of skewness which is given by this particular equation shown here.

(Refer Slide Time: 06:32)

Descriptive Statistics (Multivariate)

- Cannot be used for analyzing nominal or categorical variables
- Explores relationship of two or more variables to each other - Scatter plot
- Include more than one independent variable allowing simultaneous assessment of interrelationships between several variables
- Correlation and Regression - exploration of nature of relationship between two or more variables and strength of the relationship between them
- Line of best fit (equation of the line – regression equation) - fit a line through the points on a scatter plot
 - line is as close as possible to all points
 - represents the trend in the data

NPTEL Online Certification Courses
IIT Kharagpur

Now, going to Multivariate Statistics, generally I mean we are dealing with data which is not univariate, but it is multivariate where in we have one dependent variable and we may have multiple independent variable.

So, I mean your univariate statistics cannot be used to analyse such situation. So, I mean again when we like we were talking about univariate data in case of multivariate analysis as well I mean it is not possible for us to analyse the nominal or categorical variables, so we can as I

was telling that we can explore if there is a relationship between two or even more than two variables.

So if you have two variables, then it would come as a scatter plot you can have abscissa and ordinate x-axis and y-axis and the data plots would data values would be scattered in as paper pot, I mean diagram as I mean as dots in the x and the y axis. So, I mean if you have a three dimensional data set in which you may have one independent, dependent variable and two independent variable, in that case it would look like a cloud of points and similarly, if you have more number of variables it becomes a figure becomes complicated.

So, I mean we can include more than one variable and but the advantage of this multivariate data descriptive statistics is that we can simultaneously assess the interrelationship between multiple number of variables. Now, to do this assessment what we do is, there are two methods. We are talk about the correlation and we talk about regression.

So, in correlation and regression we nature I mean we explore the nature of relationships between the different variables and we try to assess a what is the strength of relationships between these variables. So, when we have these variables, what we can do is we can fit a line to these data points and this line is the line of the best fit, that is the line which passes through this scatter plot is worked out. So, it is as close as possible to all the points and it represents I mean this particular line or the best fit line it represents the trend in the data.

So, suppose we have this data points and we I mean try to create this data points scatter plots in this particular x, y and z values two variables which is y and z. So I mean if we try to do a best fit plot, then what happens is you can find out the distances and these distances would be minimised, these distances from this points would be minimised, so that the this is the best fit line.

(Refer Slide Time: 10:47)

Descriptive Statistics (Multivariate)

- Cannot be used for analyzing nominal or categorical variables
- Explores relationship of two or more variables to each other - Scatter plot
- Include more than one independent variable allowing simultaneous assessment of interrelationships between several variables
- Correlation and Regression - exploration of nature of relationship between two or more variables and strength of the relationship between them
- Line of best fit (equation of the line – regression equation) - fit a line through the points on a scatter plot
 - line is as close as possible to all points
 - represents the trend in the data

$$\hat{z}_i = \beta_0 + \beta_1 y_i$$

β_0 and β_1 are Intercept and Slope coefficients

The slide features a scatter plot with a blue regression line and pink data points. The axes are labeled 'y' and 'z'. The background has a light blue grid with various icons representing data and statistics.

Now, if we see the equation of the best fit line, it is in the form of y equals mx plus c where in we have the slope and we have the I mean intersect. So, this in this particular equation that is z_i equals to I mean β_0 plus $\beta_1 y_i$ in this β_0 represents the slope and the β_1 represents the intercept coefficients. So, we see this particular equation in this your z_i is equal to β_0 plus $\beta_1 y_i$.

This is the equation of this particular line which is the line of best fit in which it is as close as possible to the to all the points. So, we shall see how we minimise this distance from all these points. So, in this case this is of the form of y equals to mx plus c . So, we see the β_0 and β_1 are the intercept and the slope coefficients.

(Refer Slide Time: 12:06)

Descriptive Statistics (Multivariate)

➤ Line of best fit (equation of the line – regression equation)

➤ **Slope**
$$\beta_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(z_i - \bar{z})}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

➤ **Intercept**
$$\beta_0 = \frac{\sum_{i=1}^n z_i - \beta_1 \sum_{i=1}^n y_i}{n} = \bar{z} - \beta_1 \bar{y}$$

$$\hat{z}_i = \beta_0 + \beta_1 y_i$$

β_0 and β_1 are Intercept and Slope coefficients

The slide features a scatter plot with a blue regression line and a video inset of a man speaking. The background includes icons of a gear, a tree, and a microscope.

Now, we had talked about the I mean the best fit line. So, this is the regression equation. This equation that we had talked about is the regression equation. Now talking about the slope of this particular line that is the beta 1 it is I mean calculated using this particular equation I mean it gives you the ratio of the rise along the x to the ratio of rise along the y axis. So, in this case it is given by this particular equation. So, similarly we can work out the calculate the intercept by this equation. Again it is I mean subtracting your z mean minus beta 1 into y mean.

(Refer Slide Time: 12:55)

Descriptive Statistics (Multivariate)

➤ **Correlation Coefficient (r)** - measure of the nature and strength of the relationship between variables - degree to which points scatter around the regression line

$$r = \frac{\sum_{i=1}^n (y_i - \bar{y})(z_i - \bar{z})}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2} \sqrt{\sum_{i=1}^n (z_i - \bar{z})^2}}$$

r ranges from -1 to +1
+ve values of r indicate positive correlation
-ve values of r indicate negative correlation

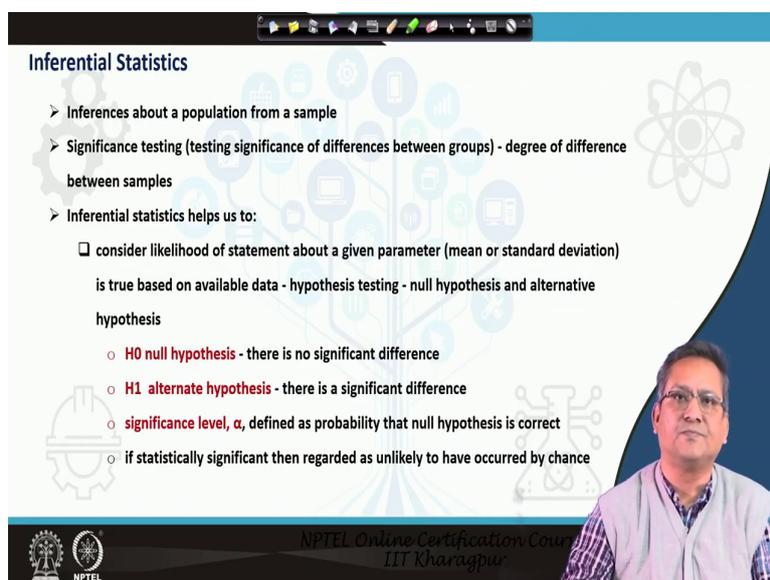
➤ **Coefficient of Determination (r²)** - indicates the goodness of fit - squared correlation coefficient - indicates that % of variation in data that can be explained by line of best fit

NPTEL Online Certification Course
IIT Kharagpur

Talking about the correlation coefficient which gives us the strength and the nature of the relationship between the variables, so it gives us the degree to which the scatter points are around the regression line. So, it is a ratio again and it is worked out using this particular equation and we can see that the r that is the correlation coefficient, it ranges from minus 1 to positive 1 value.

Now, when we have positive values of r, it indicates a positive correlation and when we have negative r values as negative, it indicates we have negative correlation. Now, the coefficient of determination that is also known as r square or in some cases you will see there are values of adjusted r square, it would give the indication of the goodness of the fit. Now, it is the square of the correlation coefficient that we had seen earlier and it I mean gives us the percentage variation in the data that is explained by the line of the best fit.

(Refer Slide Time: 14:15)



Inferential Statistics

- Inferences about a population from a sample
- Significance testing (testing significance of differences between groups) - degree of difference between samples
- Inferential statistics helps us to:
 - ☐ consider likelihood of statement about a given parameter (mean or standard deviation) is true based on available data - hypothesis testing - null hypothesis and alternative hypothesis
 - **H0 null hypothesis** - there is no significant difference
 - **H1 alternate hypothesis** - there is a significant difference
 - **significance level, α** , defined as probability that null hypothesis is correct
 - if statistically significant then regarded as unlikely to have occurred by chance

NPTEL Online Certification Course
IIT Kharagpur

Now, talking about inferential statistics when we have a population huge population, we cannot do a sampling of the entire population. So what we can do is, we can take a sample of the population and that sample has to be significant. So what we do in inferential statistics, we tried to see or do the testing of significance. So, this testing of significance is between the different groups that is the we try to measure the degree of difference between the different samples in the population. Now, this inferential statistics it helps us to consider the likelihood of statement about a given parameter.

So, when we are talking about a hypothesis before doing a research, so we bank upon this inferential statistics to see whether the hypothesis is correct or it is incorrect. So, we talk about two cases in terms of hypothesis 1 is the null hypothesis and one is the Alternate hypothesis. So when we do this inferential statistics, we can conclude that the if there is significant difference between the two groups I mean between the statistical parameters of the

two groups that is the mean or the standard deviation if the two groups were significantly different for the population.

So it gives us the null hypothesis and the other one is your alternate hypothesis that is when the significant there is a significant difference in the population, it gives the alternate hypothesis and we also have the significance level which is denoted as alpha and it is defined as the probability that null hypothesis is correct and if it is statistically significant I mean it is unlikely that the observation or the sample would have occurred by chance.

(Refer Slide Time: 16:36)

The slide is titled "Inferential Statistics" and contains the following content:

- Standard Error** – confidence interval - small standard error gives greater confidence that sample mean is close to population mean
- Formula: $SE_{\mu} = \frac{s}{\sqrt{n}}$
- Definitions: μ is the mean average, s is the sample standard deviation, n is number of observations in the sample
- t-statistic** - assesses differences between two different sample means
- Formula: $t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$
- Analysis of Variance (ANOVA)** - compare variation within data columns to variation between data columns

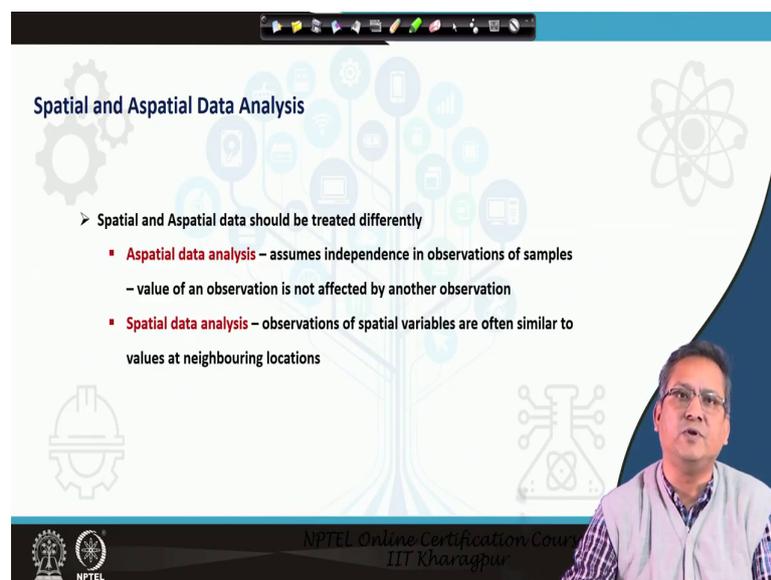
The slide also features a presenter's video feed in the bottom right corner and the NPTEL Online Certification Course logo at the bottom.

So, this is how we use the I mean or we can use your univariate statistics, multivariate statistics or inferential statistics. So, we have several other matrix in inferential statistics. First is the standard error where in we try to measure the confidence interval which is the I mean which gives us the sample mean and the population mean difference. So, if this error is small

that means there is greater confidence that the two means are closer to each other. Now, the standard error is the ratio of standard deviation to the square root of the number of the observations in the sample.

We also have a measure which is known as t-statistic which assesses the difference between the two sample means. It is calculated using this particular equation which is $t = \frac{\bar{x} - \bar{y}}{s \sqrt{\frac{1}{n} + \frac{1}{m}}}$ into root over n minus 2 divided by 1 minus r square. Now, we also analyse we can also analyse the variation within the data columns and between the data columns. So, it can be done using this particular method which is known as ANOVA which is analysis of variance.

(Refer Slide Time: 18:07)



Spatial and Aspatial Data Analysis

- Spatial and Aspatial data should be treated differently
 - **Aspatial data analysis** – assumes independence in observations of samples
 - value of an observation is not affected by another observation
 - **Spatial data analysis** – observations of spatial variables are often similar to values at neighbouring locations

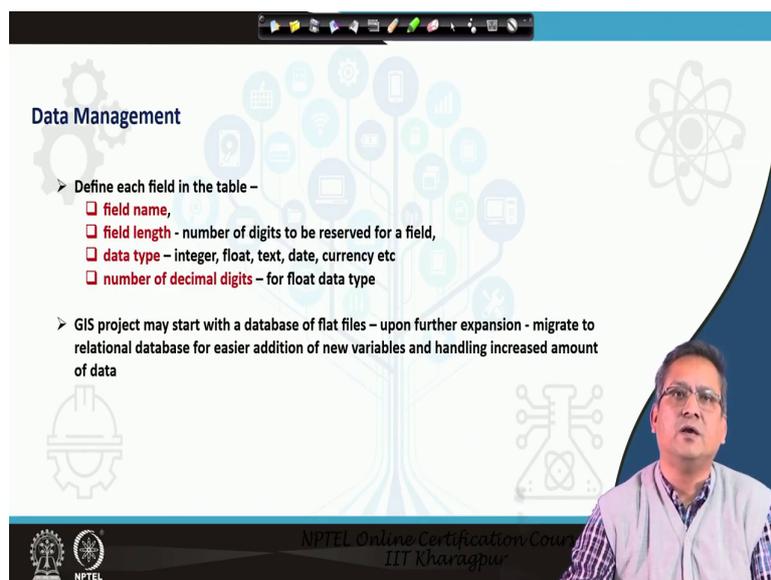
NPTEL Online Certification Course
IIT Kharagpur

So, apart from this we can also I mean find out how your spatial data and aspatial data are different from each other. So, if we run the statistical inferences on spatial data, we can see that there is a difference between the spatial data sets and the aspatial data sets does. It should

be treated differently, in a way that your the premise when we are dealing with aspatial data, any sort of statistical analysis that we had talked about whether it is univariate, multivariate statistics in that case the assumption is that there is independence in the observation of the samples that is the value of an observation is not affected by the other observations in the given data set.

But when we are talking about spatial data analysis in real world, we would see that observations in the spatial variables in the spatial domain, they are often similar to the neighbouring values. So in a way when we are to analyse the spatial data, we should have a different approach, a distinctly different approach than the approach that we use for a spatial data analysis.

(Refer Slide Time: 19:39)



The slide is titled "Data Management" and features a background with various icons representing data and technology. The content is as follows:

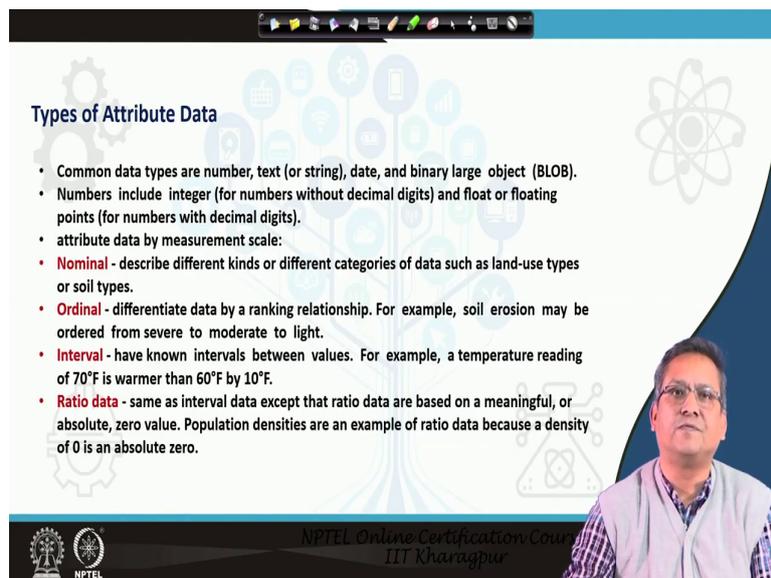
- > Define each field in the table –
 - ❑ **field name**,
 - ❑ **field length** - number of digits to be reserved for a field,
 - ❑ **data type** – integer, float, text, date, currency etc
 - ❑ **number of decimal digits** – for float data type
- > GIS project may start with a database of flat files – upon further expansion - migrate to relational database for easier addition of new variables and handling increased amount of data

The slide also includes the NPTEL logo in the bottom left corner and the text "NPTEL Online Certification Course IIT Kharagpur" in the bottom right corner. A video inset in the bottom right shows a man speaking.

Now, talking about Data Management we can define the each field in the data table where in we give the field name, we can specify the field length which is the number of digits to be reserved for a field. We can also give the data type whether that data is integer, it is float or it is a text information, whether it is a date or a currency information and we can also give the number of decimal digits in case the data is a float data type.

So, I mean we when we are doing this analysis, initially we can start with flat files. So I mean when we have more amount of data, then we can get into a relational database. If we have a huge data set, we can create a relational database and we can it becomes easier to handle in such a case when we are using the relational database management for handling huge amount of data.

(Refer Slide Time: 20:51)



The slide is titled "Types of Attribute Data" and features a list of data types. The background is light blue with various icons representing data and technology. A video feed of a presenter is visible in the bottom right corner. The slide content is as follows:

Types of Attribute Data

- Common data types are number, text (or string), date, and binary large object (BLOB).
- Numbers include integer (for numbers without decimal digits) and float or floating points (for numbers with decimal digits).
- attribute data by measurement scale:
 - **Nominal** - describe different kinds or different categories of data such as land-use types or soil types.
 - **Ordinal** - differentiate data by a ranking relationship. For example, soil erosion may be ordered from severe to moderate to light.
 - **Interval** - have known intervals between values. For example, a temperature reading of 70°F is warmer than 60°F by 10°F.
 - **Ratio data** - same as interval data except that ratio data are based on a meaningful, or absolute, zero value. Population densities are an example of ratio data because a density of 0 is an absolute zero.

NPTEL Online Certification Course
IIT Kharagpur

So talking about the different types of attribute data, we have the number text date binary or large binary large object which is also known as BLOB data. When we are talking about numbers, we would be dealing with either I mean real numbers, integer numbers. So, we can have integer values or we can have float values. So, when we are coding the data specially for the raster type of data, we have to be very careful about what kind of data set we are creating because say suppose if you have a float operation and your data type you are defining it as integer, it would discard all the float values.

So, we have to be very careful about the I mean the kind of data that we are handling. So, specifically in case of raster. So, similarly we can have float data. So, the data resolution also would come into picture when we are dealing with raster data set. It could be a 2 bit data, it could be a 4 bit or 8 bit data, 12 bit data depending on the I mean the size of the I mean difference of the highest values and the lowest values.

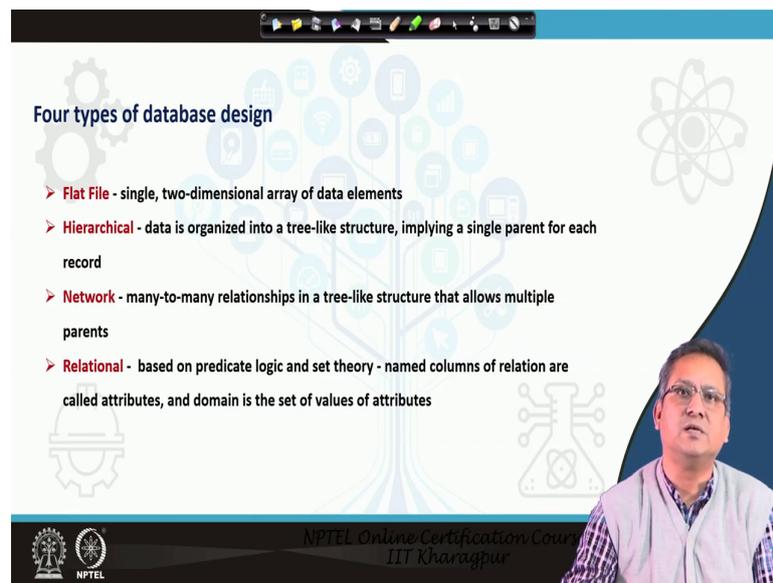
Now, attribute data is measured by different scales. So, we have different types of attribute data. The first one that we come across is the nominal data which describe the different kinds or categories of data. So, examples of this kind of data could be your land use data or data are such a soil data, so where in you give the categories of land use or the categories of soil. The next data that we deal with is the ordinal data. So, it I mean differentiates the data by a ranking relationship.

So, we can say suppose we have a data set, so we can quantify the intensity like in this particular case we are talking about soil erosion. So, whether the erosion is severe or moderate, we can also talk about proximity matrix like it is very near or far or I mean we can have different types of measures and rank the relationships. Now, the interval data are the intervals between different values have known intervals, like suppose we can categorise the temperature data into different groups and we can have group where in we can have a cold data temperature which are normal or comfortable and temperatures which are warmer.

So, we can categorise the data and these are known as interval data. The last one which we often used for data modelling in GIS is the ratio data and this is the most powerful tool which

we have. So, we can use both integer type of your numerical data or float type data when we work with ratio data sets. So, your I mean it is the ratio basically. So, I mean for an example we can talk about the population density in different words which is an example of ratio data.

(Refer Slide Time: 24:33)



Four types of database design

- > **Flat File** - single, two-dimensional array of data elements
- > **Hierarchical** - data is organized into a tree-like structure, implying a single parent for each record
- > **Network** - many-to-many relationships in a tree-like structure that allows multiple parents
- > **Relational** - based on predicate logic and set theory - named columns of relation are called attributes, and domain is the set of values of attributes

NPTEL Online Certification Course
IIT Kharagpur

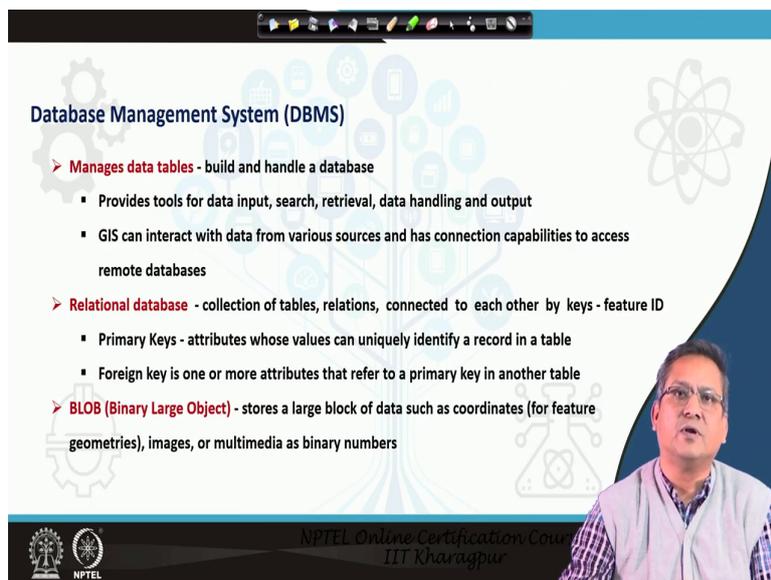
Now, there are four types of database design whenever we take GIS tasks, we can create different types of database depending on the size of the database. So, if the database is small we can create a flat file nomenclature where in we have a single file and the data is arranged in a two dimensional array of data elements. The next one is the hierarchical data in which the data can be organised in a tree like structure and it implies that there is a single parent for each record.

The next data type is the network data type. So, it is a modification, further modification of the hierarchical data set and it embodies many to many relationships in a tree like structure in a

hierarchical structure. So, this data structure or this type of database design it allows for multiple parents. So, if you have leaf nodes you can have multiple parents to those leaf nodes.

The last type for the database design is the relational database which is the most powerful data I mean database design when you have very big data sets. So, we generally used or relational database management system to model the entire data to do the designing of the entire data sets. So, I mean it is based on predicate logic and it is based on set theory. So, we have name columns of relation which are called attributes and domain and this domain is the sets of values of the attributes.

(Refer Slide Time: 26:27)



The slide is titled "Database Management System (DBMS)" and contains the following content:

- **Manages data tables** - build and handle a database
 - Provides tools for data input, search, retrieval, data handling and output
 - GIS can interact with data from various sources and has connection capabilities to access remote databases
- **Relational database** - collection of tables, relations, connected to each other by keys - feature ID
 - Primary Keys - attributes whose values can uniquely identify a record in a table
 - Foreign key is one or more attributes that refer to a primary key in another table
- **BLOB (Binary Large Object)** - stores a large block of data such as coordinates (for feature geometries), images, or multimedia as binary numbers

The slide also features a video inset of a man speaking in the bottom right corner, the NPTEL logo in the bottom left, and the text "NPTEL Online Certification Course IIT Kharagpur" at the bottom center.

So, we can have database management systems which work in the back end when it handles your GIS data sets the attribute data. So, it builds and handles the GIS database. So, these DBMS tools they provide I mean solutions for data input for search and retrieval for data

handling and for generating output from your queries, your GIS can interact from I mean data from multifarious sources.

So, I mean we can connect to remote data bases and this GIS has the capability to access such data bases I mean multiple data bases. So to I mean connect this multiple databases, we would need to have a unique field which is known as keys. So, we have a relational database where in we have collection of tables and they would be connected to each other by a feature id which is known as keys and the relations are built into it.

There are different types of keys. First is we talk about the primary keys I mean these are the attributes whose values are unique and it can be identified as a record in a table. We also have a foreign key which is one or more attribute that refer to the primary key in another table, another reference table in which some other data sets are available.

We had also talked about the BLOB which is the binary large object which stores of huge block of data I mean for example it could be the coordinates of the I mean the points or the lines I mean for I mean generally we store the feature geometries as block files. So, it could be also images or it could be multimedia data and these I mean files generally store this coordinates as binary numbers.

(Refer Slide Time: 28:55)

Database Management System (DBMS)

Types of Relationships

➤ **One-to-one relationship** - one and only one record in a table is related to one and only one record in another table

NPTEL Online Certification Course
IIT Kharagpur

Now, the types of relationships that we have in the database management system could be of the four types. The first one is the one-to-one relationship. So, each record in the table is related to one and only record, one and only record in another table. So, if we have two tables it is a one-to-one relationships I mean you do not see multiple connections in this type of relationships.

(Refer Slide Time: 29:27)

The slide is titled "Database Management System (DBMS)" and "Types of Relationships". It lists two types of relationships:

- > **One-to-one relationship** - one and only one record in a table is related to one and only one record in another table
- > **One-to-many relationship** - one record in a table may be related to many records in another table

The diagram below illustrates a one-to-many relationship. It shows two vertical columns of rectangular boxes representing records in two different tables. The left column has 4 boxes, and the right column has 6 boxes. Lines connect the top box of the left column to the top two boxes of the right column, and the second box of the left column to the next two boxes of the right column. This visualizes one record from the first table being linked to multiple records in the second table.

NPTEL Online Certification Course
IIT Kharagpur

The next one that is the one-to-many relationship. In this one record of the table is related to many records in another table. So, you can see for each of these particular records they are related to multiple records in other table.

(Refer Slide Time: 29:44)

Database Management System (DBMS)

Types of Relationships

- **One-to-one relationship** - one and only one record in a table is related to one and only one record in another table
- **One-to-many relationship** - one record in a table may be related to many records in another table
- **Many-to-one relationship** - many records in a table may be related to one record in another table

The diagram illustrates a many-to-one relationship. It shows two tables represented as horizontal bars. The left table has five bars, and the right table has three bars. Lines connect each of the five bars on the left to one of the three bars on the right, demonstrating that multiple records in one table are related to a single record in another table.

NPTEL Online Certification Course
IIT Kharagpur

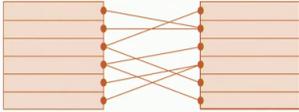
Now, there could be many-to-one relationship. In this many record in the table I mean your attribute table may be related to one record in another table. So, in this case you have your attribute table, GIS table and you may have another database where in you see there are multiple connections from your input database to your to another database another data table.

(Refer Slide Time: 30:13)

Database Management System (DBMS)

Types of Relationships

- **One-to-one relationship** - one and only one record in a table is related to one and only one record in another table
- **One-to-many relationship** - one record in a table may be related to many records in another table
- **Many-to-one relationship** - many records in a table may be related to one record in another table
- **Many-to-many relationship** - many records in a table may be related to many records in another table



NPTEL Online Certification Course
IIT Kharagpur

The next one is many-to-many relationships in which many records in a table would be related to many records in another table. So, this is how these two tables are. They would have their identifiers or the keys that we have talked about. So, I mean they would be related to each other in a many-to-many relationship.

(Refer Slide Time: 30:34)

Database Management System (DBMS)

Joining Attribute Data

➤ Non-spatial table to a feature attribute table - feature attribute table (has primary key) is the origin and the other table is the destination (has foreign key)

Primary key										Foreign key									
key	state	population	%	male	female	difference	urban	rural	area	key	state	urban	rural	area	density				
1	Uttar Pradesh	19912141	26.2	1041708	9511181	1048579	5021	153448	4427923	24919	826								
2	Madhya Pradesh	11217403	9.28	5824506	5411277	411179	502	634548	592731	10731	361								
3	Bihar	10459442	8.6	5427817	4881295	465882	513	537520	1172609	9413	110								
4	West Bengal	9127613	7.54	4682527	4487088	242539	504	422197	293400	8572	103								
5	Madhya Pradesh	7262809	6	3761236	3504593	257093	501	523279	1097666	39942	256								
6	Tamil Nadu	7214730	5.56	3613735	3600955	12820	506	1718120	1494729	130558	535								
7	Rajasthan	6544837	5.46	3335917	3209140	253537	528	1534238	1760178	142339	201								
8	Karnataka	6102527	5.06	3096917	3010940	88817	523	2726220	2328129	19376	112								
9	Gujarat	6045982	4.99	3145206	2894642	2452828	519	1467917	2571261	190204	336								
10	Andhra Pradesh	49181799	4.08	2473828	2464871	8937	504	1477639	1463010	146205	301								
11	Odisha	4174428	3.47	2121118	2076302	45054	529	1401113	499164	133797	262								
12	Telangana	3519378	2.51	1756478	1748900	21478	540	1121513	1380665	114840	137								
13	Kerala	3440261	2.76	1622412	1771649	-151237	524	1740508	5351271	38910	879								
14	Haryana	2708134	2.22	1391613	1307769	82944	542	2250969	292762	7974	414								
15	Assam	3105526	2.58	1531843	1536113	-47320	504	2679626	438754	76418	137								
16	Punjab	2774118	2.29	1403465	1310873	92512	495	1711890	1016749	50162	500								
17	Chhattisgarh	2542428	2.12	1281819	1271283	10536	491	1760610	981638	12516	109								
18	Haryana	2055482	2.09	1044734	1104738	-60006	479	1815493	882158	44212	571								
19	Delhi (UT)	1678741	1.39	887726	790915	98871	484	144227	1240730	1484	1127								
20	Jammu and Kashmir	1242103	1.06	646462	595640	50822	489	114020	144206	22234	31								
21	Uttarakhand	1008232	0.83	511777	496519	15254	461	702083	306189	53483	189								
22	Himachal Pradesh	684602	0.57	348173	336279	9944	472	617305	68754	55673	123								
23	Tripura	607187	0.5	324716	279561	45155	465	272021	89026	10464	100								
24	Meghalaya	266889	0.22	141812	147057	-16775	469	218971	59038	22429	112								
25	Mizoram	285794	0.24	143887	141767	2120	465	289924	822112	22127	120								
26	Nagaland	178202	0.16	104449	110253	-7074	451	148861	51741	16579	122								
27	Goa	145845	0.12	739140	73485	19735	471	155434	86080	3702	334								
28	Arunachal Pradesh	138727	0.11	72112	68915	4897	438	108165	31346	8743	17								
29	Andhra Pradesh (UT)	134763	0.1	62511	62462	49	437	36441	81012	478	290								
30	Mizoram	109706	0.09	55319	54387	1342	476	52621	56197	2281	52								
31	Chandigarh (UT)	105450	0.09	58963	47677	10586	431	2926	102462	114	922								
32	Sikkim	61071	0.05	32470	28707	3563	464	42041	11274	3766	60								
33	Andaman and Nicobar	88581	0.03	20871	17730	25161	476	14441	13033	6249	46								
34	Dadra and Nagar Haveli	44329	0.03	19120	49549	-43811	274	18324	23829	491	490								
35	Daman and Diu (UT)	24247	0.02	12681	12646	35	414	6811	15246	113	2100								
36	Lakshadweep (UT)	64473	0.01	33123	31350	1773	543	1421	3038	12	2013								

Origin Destination

NPTEL Online Certification Courses
IIT Kharagpur

Now, talking about database management system we can join the attribute data I mean we can do a non-spatial join of the attribute table. So if we have multiple attribute table, in this case you can see that we have the population tables for different states of India in which we have population, we have the male and the female population, we have the difference between the male and the female and the sex ratio in the first table which has a primary key and in the next table, we have the total population which is urban and the total population which is rural we have the area as well as the density and it has a foreign key.

So, in this in the earlier slide we had talked about the primary key and the foreign key. So, you can see how they are located and they would be related or in any query or search operation or processing further processing these two keys would be related for in a non-spatial join. So, we can see how we can do a non-spatial join.

(Refer Slide Time: 31:55)

Database Management System (DBMS)

Merging Attribute Data

- **Join** - brings together two tables by using keys or a field common to both tables and appends all columns from one table to the other based on Key
- **Relate** - temporarily connects two tables by using keys or a field common to both tables
- **Spatial Join** - uses a spatial relationship to join two sets of spatial features and their attribute data

Layer 1
● Population (Sum)
● Urban (Sum)

Spatial Join - Proportion of Urban Population in India

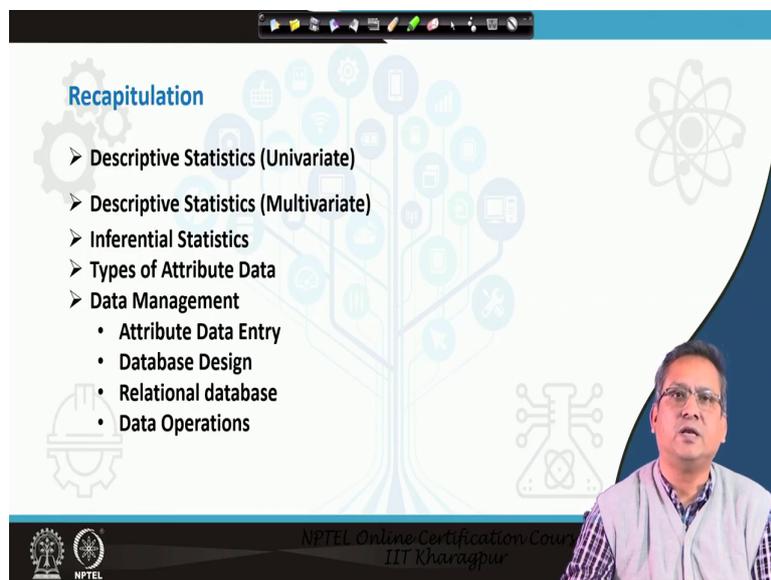
NPTEL Online Certification Courses
IIT Kharagpur

So, for merging the attribute data there are few options which are available in GIS softwares. So, most of these packages have got these operations. So, first one is the join operation where in the two tables using the keys, the common keys is are join and the columns are appended from one table to the another table. So, for the in the input table the I mean it would become an extended table, where in the all the fields would be appended.

The next one, next way to merge the attribute data is the relate operator, where in it temporally connects two tables by using the common keys or the fields that we have seen are common tool, both the tables. Now, the 3rd one is the spatial join which uses a spatial relationship to join the two data sets of the spatial features as well as their attribute data. So, we can see an example of spatial data, a spatial join.

So, in our earlier slide we had seen the data pertaining to the population statistics for different states of India. So, you can see here a spatial join operation has been done in the blue pie chart, you can see that it is the population and the orange one gives you the fraction of the urban population in the different states. So, you can identify the states where in which has the highest amount of urban population those two tables were joined together using a spatial join and we can see the output here.

(Refer Slide Time: 33:54)



The image shows a presentation slide titled "Recapitulation" with a list of topics covered. The slide is part of an NPTEL Online Certification Course from IIT Kharagpur. A presenter is visible in the bottom right corner of the slide frame.

- Descriptive Statistics (Univariate)
- Descriptive Statistics (Multivariate)
- Inferential Statistics
- Types of Attribute Data
- Data Management
 - Attribute Data Entry
 - Database Design
 - Relational database
 - Data Operations

NPTEL Online Certification Course
IIT Kharagpur

So, recapitulation of what we have covered today, we have talked about Univariate Descriptive Statistics, where in we had talked about the central tendencies of mean, median and mode for a univariate data which has only one attribute column. We had then talked about multivariate data analysis, we had talked about inferential statistics, we had talked about the different types of attribute data and then we have finally talked about data management.

We had talked about how we can do the data entry and what are the different types of data. In that we had talked about the database design, we had talked about the relational database, the data operations. So, thank you for your patient hearing till we meet again in the next lecture.

Thanks so much.