

# Day - 10: Understanding Data and Descriptive Statistics from Theory to Python Implementation

23rd September 2024

## Introduction to Data Types:

### What is Data?

A collection of Facts, number of qualitative Characteristics that represent information about an event, object or system...

Data is the core building block that allows for the extraction of insights, decision-making, and predictions in Data Science

### Why is Data important in Data Science?

- Foundation of Data Science
- Making Data-Driven Decisions
- Transforming Raw Data into Insights

### Examples:

- Risk Management
- Fraud Detection
- Inventory Management
- Customer Segmentation
- Renewable Energy Forecasting
- Autonomous Vehicles
- Performance Analysis

### Types of Data and their Uses:

<u>Aspect</u>	<u>Data Type</u>	<u>Description</u>	<u>Examples</u>	Representation	Advantage	Disadvantage
<u>Source</u>	Primary Data	Data collected directly by the researchers for specific purpose	Surveys, Experiments, Interviews and sensors	Line Charts, Bar Charts, Histogram	1. Relevance 2. Control over Data Quality	1. Time-Consuming 2. Costly
<u>Source</u>	Secondary Data	Pre-existing data collected for different purpose, repurposed for analysis	Government Reports, Academic Datasets	Heat map, Box Plots, Violin Plot, Area Charts	1. Easily Accessible 2. Cost-Effective	1. Relevance and Fit 2. Quality and Accuracy Concerns 3. Need for Cleaning
<u>Frequency</u>	Frequency Data	Data Showing how often an event or observation occurs...	Customer name, Qualitative Responses	Histogram, Bar Charts, Frequency Distribution	1. Insightful for analysis 2. Facilitates Decision Making	1. Limited Context 2. Potential for Misinterpretation

<b><u>Aspect</u></b>	<b><u>Data Type</u></b>	<b><u>Description</u></b>	<b><u>Examples</u></b>	<b><u>Representation</u></b>	<b><u>Advantage</u></b>	<b><u>Disadvantage</u></b>
<b><u>Frequency</u></b>	Non-Frequency Data	Data that does not represent counts or occurrences...	Qualitative Nature, Categorization	Contour Plots, Scatter Plots, Violin Plots	1. Rich in Detail 2. Flexibility in Analysis	1. Subjectivity 2. Challenges in Quantification
<b><u>Nature</u></b>	Quantitative Data	Numerical information that can be measured and expressed as numbers.	Measurable Values, Types - Discrete and Continuous	Histogram, Scatter Plot, Box plot	1. Statistical Analysis 2. Objective measurement	1. Limited Context 2. Potential Misinterpretation
<b><u>Nature</u></b>	Qualitative Data	Consists of Non-Numerical information that represents descriptions, characteristics, or qualities...	Descriptive Nature, Categorization	Bar Chart, Pie Chart, Stacked Bar Chart	1. In-depth insights 2. Flexibility in Approach	1. Subjectivity 2. Challenges in Quantification
<b><u>Measurement</u></b>	Discrete Data	Consists of numerical information that can only take specific, distinct values, often representing counts of item values.	Specific Values, Countable Events	Bar Charts, Dot Plots	1. Simplicity in Analysis 2. Clear Categorization	1. Limited Range of Values 2. Potential loss of Information
<b><u>Measurement</u></b>	Continuous Data	Refers to numerical information that can take value within a given range.	Range of Values, Measurement Precision	Histograms, Frequency Curves	1. Detailed Analysis 2. Rich Insights	1. Complexity in Interpretation 2. Potential Measurements Errors
<b><u>Structure</u></b>	Structured Data	Refers to information thats organized in a defined format like in a matrix (rows and columns)	Sales database, Excel Spreadsheet	Line Chart, Heat Map, Bar Chart	1. Ease of Access and Analysis 2. High Data Integrity	1. Limited Flexibility 2. Inability to Capture Complex Data
<b><u>Structure</u></b>	Unstructured Data	Refers to Information that does not follow a specific format or organization	Social Media Posts, Customer Feedback Emails	Word Cloud, Network Graph, Histogram	1. Captures Complex Information Rich Source of Information	1. Complex Analysis 2. Storage and Processing Challenges
<b><u>Category</u></b>	Numerical Data	Consists of values represented as numbers. an further be classified in to discrete and continuous types.	Sales Figure, Temperature	Line Charts, Histograms	1. Robust Analysis 2. Detailed Insight	1. Sensitivity to Outliers 2. Complexity
<b><u>Category</u></b>	Categorical Data	Represents groups or	Gender, Product	Bar Charts, Pie Charts	1. Simplicity 2. Clear	1. Limited Statistical

Aspect	Data Type	Description	Examples	Representation	Advantage	Disadvantage
		categories of information, often non-numeric in nature. Describes qualitative attributes of characteristics. divided into ordinal and nominal data.	Category		Comparison	Analysis 2. Potential Loss of Information

### Common issues of the Data:

Data Type	Common Issues	How to Overcome
<b><u>Primary Data</u></b>	<ol style="list-style-type: none"> <li>1. Time consuming and costly to collect</li> <li>2. requires well-defined collection to avoid bias</li> <li>3. May lack historical data for comparison</li> </ol>	<ol style="list-style-type: none"> <li>1. Use Automated Tools for data Collection</li> <li>2. Implement standardization methods to reduce bias</li> <li>3. Complement Primary Data with secondary data for historical Events</li> </ol>
<b><u>Secondary Data</u></b>	<ol style="list-style-type: none"> <li>1. May not be specific to researcher's needs</li> <li>2. Could be outdated or irrelevant</li> <li>3. Often Requires Cleaning and reformatting before use.</li> <li>4. Potential Issues with data accuracy or trustworthiness</li> </ol>	<ol style="list-style-type: none"> <li>1. Cross-Validate data with other sources</li> <li>2. Regularly update Datasets</li> <li>3. Use data cleaning techniques to ensure accuracy</li> <li>4. Investigate Data origins for reliability</li> </ol>
<b><u>Frequency Data</u></b>	<ol style="list-style-type: none"> <li>1. High variability can obscure meaningful patterns</li> <li>2. Aggregation can lead to data loss</li> <li>3. difficult to handle with very large dataset</li> </ol>	<ol style="list-style-type: none"> <li>1. Use smoothing techniques like moving averages</li> <li>2. Keep raw data and provide both detailed and summarized views</li> <li>3. Implement robust database management systems for large datasets</li> </ol>
<b><u>Non-Frequency Data</u></b>	<ol style="list-style-type: none"> <li>1. Harder to apply statistical methods or quantitative analysis</li> <li>2. Requires categorization for analysis, which can introduce subjectivity</li> </ol>	<ol style="list-style-type: none"> <li>1. Convert non-frequency data into usable formats via classification techniques</li> <li>2. Apply consistent and transparent categorization criteria to minimize subjectivity</li> </ol>
<b><u>Quantitative Data</u></b>	<ol style="list-style-type: none"> <li>1. Outliers can skew the results</li> <li>2. Assumptions like normal distribution may not always hold true.</li> <li>3. Requires careful handling of missing or incomplete data</li> </ol>	<ol style="list-style-type: none"> <li>1. Use Robust statistics like the median for outlier resistance</li> <li>2. Test for distribution assumptions and apply appropriate methods</li> <li>3. Impute or remove missing values based on analysis needs</li> </ol>
<b><u>Qualitative Data</u></b>	<ol style="list-style-type: none"> <li>1. Difficult to analyze systematically.</li> <li>2. Requires thematic coding, which can be subjective</li> <li>3. Large volumes of data may require advanced tools like Natural Language Processing</li> </ol>	<ol style="list-style-type: none"> <li>1. Use qualitative analysis software for organization</li> <li>2. Maintain consistent coding practices</li> <li>3. Utilize AI or NLP tools to process large volumes efficiently</li> </ol>
<b><u>Continuous Data</u></b>	<ol style="list-style-type: none"> <li>1. Requires precision in measurement tools</li> <li>2. Small measurement errors can propagate in calculation</li> <li>3. Potential difficulty in data aggregation and summarization</li> </ol>	<ol style="list-style-type: none"> <li>1. Use high-quality measurements tools and calibrate regularly</li> <li>2. Implement error-checking algorithms</li> <li>3. Use appropriate aggregation techniques that account for precision loss</li> </ol>
<b><u>Discrete Data</u></b>	<ol style="list-style-type: none"> <li>1. Can lead to oversimplification of phenomena</li> <li>2. Limited granularity may miss important leads</li> <li>3. Not always suitable for predictive models</li> </ol>	<ol style="list-style-type: none"> <li>1. Use hybrid approaches, combining discrete and continuous data when possible</li> <li>2. Increase Data resolution if needed</li> <li>3. Choose models better suited for discrete data, like decision tree</li> </ol>

<u>Data Type</u>	<u>Common Issues</u>	<u>How to Overcome</u>
<b><u>Structured Data</u></b>	<ol style="list-style-type: none"> <li>1. Inflexible format, requires predefined schema</li> <li>2. Scaling can be difficult as data grows</li> <li>3. Cannot easily handle unanticipated variations of data</li> </ol>	<ol style="list-style-type: none"> <li>1. Implement Flexible database structures (NoSQL)</li> <li>2. Use Modular Data Models</li> <li>3. Regularly Update schema based on emerging data trends</li> </ol>
<b><u>Unstructured Data</u></b>	<ol style="list-style-type: none"> <li>1. Hard to process and analyze without advanced techniques (e.g., text mining, AI)</li> <li>2. May require extensive storage and processing resources.</li> <li>3. Quality of insights depends heavily on preprocessing and cleaning.</li> </ol>	<ol style="list-style-type: none"> <li>1. Use AIML Techniques for better processing</li> <li>2. Invest in scalable cloud storage and processing solutions</li> <li>3. Implement Rigorous data cleaning protocols before analysis</li> </ol>
<b><u>Categorical Data</u></b>	<ol style="list-style-type: none"> <li>1. Risk of oversimplification when grouping complex data</li> <li>2. May require subjective grouping (e.g., high or low)</li> <li>3. Harder to quantify relationship between categories</li> </ol>	<ol style="list-style-type: none"> <li>1. Use more nuanced category creation and maintain transparency in groupings</li> <li>2. Test Category groupings with different thresholds</li> <li>3. Apply statistical tests for categorical relationships like Chi-square</li> </ol>
<b><u>Numerical Data</u></b>	<ol style="list-style-type: none"> <li>1. Assumptions about data distribution might not hold</li> <li>2. Outliers can heavily influence metrics like the mean</li> <li>3. Missing values can complicate statistical analysis</li> </ol>	<ol style="list-style-type: none"> <li>1. Use non-parametric tests if assumptions fail</li> <li>2. Apply outlier detection and removal techniques</li> <li>3. Impute missing data based on statistical models</li> </ol>

## **Data Visualization:**

### **What is Data Visualization?**

It refers to the graphical representation of information and data. It uses elements like charts, graphs, maps, etc. It provides tools to provide an accessible way to understand trends, outliers and patterns in data.

### **Why is Data Visualization Needed?**

- Simplifies Complex Data
- Detect Patterns and Trends
- Communicate Results Effectively
- Quick Decision Making

### **How does Data Visualization Help us?**

- Improves Data Accessibility
- Reveals Hidden Insights
- Data Quality Validation
- Better Understanding of Relationships

### **Benefits of Data Visualization?**

1. Enhanced Comprehension
2. Increased Engagement
3. Decision Support
4. Accessibility for Non- Technical Users
5. Highlighting Key Metrics

## **Visualization for Exploratory Data Analysis:**

EDA or Exploratory Data Analysis is the process of investigating datasets to summarize their main characteristics, often with visual methods.

1. Initial Data Exploration
2. Pattern and Trend Identification
3. Understanding Distributions
4. Spotting Anomalies
5. Feature Selection and Engineering
6. Validation of Assumptions
7. Clustering and Segmentation
8. Comparing Groups or Categories

### Types of Charts and their uses:

<u>Chart Name</u>	<u>Definition</u>	<u>Description</u>	<u>Where to Apply</u>	<u>How to Infer:</u>
<b><u>Bar Chart</u></b>	Rectangular bars to represent data values	<ol style="list-style-type: none"> <li>Used to compare categorical data</li> <li>Categories on X axis and values on Y axis</li> <li>Can be plotted vertically or horizontally</li> </ol>	Categorical Comparisons (Sales by region, number of product sold)	Taller Bar represents higher values.
<b><u>Pie Chart</u></b>	represents data as slices of a circle with each slice representing portion of a whole	<ol style="list-style-type: none"> <li>Used for showing proportions or percentages</li> <li>Each slice size is proportional to the quantity it represents</li> </ol>	Representing Parts of a Whole (market share, budget distribution)	The larger the slice, the larger the proportion of that category within the dataset
<b><u>Histogram</u></b>	Distribution of a dataset by dividing it into bins and counting the numbers of data points that falls into each bin.	<ol style="list-style-type: none"> <li>Used to understand the distribution of continuous data</li> <li>X-axis represents intervals or bins and Y-axis represents frequency</li> </ol>	Examining the distribution of continuous variables like age distribution	The height of each bin represents the frequency of values within that range. it helps to assess weather
<b><u>Scatter Plot</u></b>	Uses dots to represent the values of two different variables	<ol style="list-style-type: none"> <li>Visualizes relationships or correlations between two continuous variables</li> <li>Each point represents an observation in the dataset</li> </ol>	Exploring relationship between two continuous variables (Height vs Weight, sales vs marketing spend)	A positive correlation shows an upward trend, and a negative correlation shows a downward trend. No visible pattern indicates no correlation
<b><u>Box Plot</u></b>	Summarizes data by showing its median, quartiles, and outliers in a graphical format	<ol style="list-style-type: none"> <li>Understand distribution of data, especially its spread, and any outliers</li> <li>Represents the IQR and the whiskers are the outliers showing the extent of the data</li> </ol>	Summarizing data distributions like test scores and stock prices	The line inside the box represents the median. Outliers are plotted as individual points. The Whiskers give a sense of the range
<b><u>Line Plot</u></b>	Uses line to connect individual data points, typically showing the change of a variable over time	<ol style="list-style-type: none"> <li>Often used to represent time series data</li> <li>X-axis represent time or sequential events and Y-axis represents the values</li> </ol>	Tracking changes over time (stock prices, sales trends)	The slope of the line shows whether the variables is increasing, decreasing, or staying constant over time
<b><u>Violin Plot</u></b>	Similar to boxplot but rotated kernel density plot on each side.	<ol style="list-style-type: none"> <li>Combines a Box-plot with density estimation</li> <li>Shows the probability density of the data at different values</li> </ol>	analyzing distribution and their shapes comparing distributions between categories	The thickness of the plot shows the density of the values, wider part indicates more datapoints

<b><u>Chart Name</u></b>	<b><u>Definition</u></b>	<b><u>Description</u></b>	<b><u>Where to Apply</u></b>	<b><u>How to Infer:</u></b>
Heat Map	Represents data in matrix form, where values are depicted using color gradients	1. Visualizes correlation Matrices or grid based data 2. Color represents the magnitude of values in the data	Visualizing correlation matrices or distribution of values over 2D space	Darker or more intense colors indicate higher values, lighter colors indicate lower values
Pair Plot	a matrix of scatter plot used to show pairwise relationships between multiple variables in a dataset	each scatter plot shows the relationship between two variables, allowing a quick overview of potential correlations, trends, and distributions	best for EDA to identify trends and correlations between multiple numerical variables	Look for patterns in the scatter plots, such as linear relationships or clustering of points
Stacked Bar Chart	A bar Chart where different categories within each bar are stacked on top of each other	Useful for comparing both the totals and the sub-components(categories) within the data	Use for data with hierarchical categories to compare the total values and the contribution of each sub-category	Compare the height of each bar assess the totals and see how each part contributes to the whole
Boxen Plot	Provides more detailed information about the data distribution especially the tails.	Focuses on large datasets and display more granular information about the distribution, for extremes specifically	Use for large datasets with skewed distributions where you want to visualize data in more details especially outliers and extremes	Similar to box plot but provides deeper insights into the tail
Density Plot	A smooth curve representing the distribution of a continuous variables	Used to estimate the probability density function of a continuous variable, showing where data points are concentrated	Best for visualizing the distribution of continuous data and comparing distribution across different groups	Look for peaks in the curve, which indicate areas of high concentration. Multiple peaks suggest bimodal or multimodal distribution
Area Chart	line chart where area under the line is filled with color to show magnitude	typically used to visualize cumulative frequency over time or between different categories	use when you want to emphasize the magnitude of the changes over time or across categories	look at the shape of the filled area to understand trends and cumulative totals.
Bubble Chart	A scatter plot where the size of each point represents a third variable	Allows visualization of three variables at once by encoding two variables in the axes and the third in the size of the bubble	Useful when you want to visualize relationships between three variables and show relative importance via bubble size	Larger Bubbles indicate higher values of the third variables, while the position on the axes shows the relationship between the other two variables
Dot Plot	Uses dots to represent individual data points	Useful for visualizing distributions and frequencies in small datasets. Similar to bar chart but more focused on individual points	Ideal for smaller datasets where you want to highlight individual data points and their frequency	The Spread of dots along the axis indicates the distribution and relative frequency of values