



Exploratory Data Analysis

Dr.G.Malathi
Associate Professor, SCSE
Image Processing and Vision - Head
Vellore Institute of Technology, Chennai

Syllabus

- 0 CSE5007 Exploratory Data Analysis L,T,P,J,C
2,0,0,4,3
- 0 Objective:
- 0 This course introduces the methods for data preparation and data understanding. It covers essential exploratory techniques for understanding multivariate data by summarizing it through statistical methods and graphical methods

Syllabus

0 Expected Outcomes

0 After successfully completing the course the student should be able to

0 1. Handle missing data in the real world data sets by choosing appropriate methods

0 Summarize the data using basic statistics

0 Visualize the data using basic graphs and plots.

0 Identify the outliers if any in the data set.

0 Choose appropriate feature selection and dimensionality reduction techniques for handling multi-dimensional data.

Syllabus

0 **Module1 INTRODUCTION TO EXPLORATORY DATA ANALYSIS**

- 0 Data Analytics life cycle,
- 0 Exploratory Data Analysis (EDA) -Definition,
- 0 Motivation,
- 0 Steps in data exploration,
- 0 The basic data types,
- 0 Data Type Portability

Syllabus

0 **Module2 PREPROCESSING - TRADITIONAL METHODS AND MAXIMUM LIKELIHOODESTIMATION**

- 0 Introduction to Missing data,
- 0 Traditional methods for dealing with missing data,
- 0 Maximum Likelihood Estimation – Basics,
- 0 Missing data handling,
- 0 Improving the accuracy of analysis

Syllabus

0 **Module3 PREPROCESSING - BAYESIAN ESTIMATION**

- 0 Introduction to Bayesian Estimation,
- 0 Multiple Imputation - Imputation Phase, Analysis and Pooling Phase,
- 0 Practical Issues in Multiple Imputation,
- 0 Models for Missing Not at Random Data

Syllabus

0 **Module4 DATA SUMMARIZATION & VISUALIZATION**

- 0 Statistical data elaboration,
- 0 1-D Statistical data analysis,
- 0 2-D Statistical data Analysis,
- 0 N-D Statistical data analysis

Syllabus

0 **Module5** OUTLIERANALYSIS

0 Introduction,

0 Extreme Value Analysis,

0 Clustering based, Distance Based and
Density Based outlier analysis,

0 Outlier Detection in Categorical Data

Syllabus

0 **Module6 FEATURE SUBSET SELECTION**

0 Feature selection algorithms:

0 filter methods,

0 wrapper methods and

0 embedded methods,

0 Forward selection,

0 backward elimination,

0 Relief,

0 greedy selection,

0 genetic algorithms for feature selection

Syllabus

0 **Module7** DIMENSIONALITY REDUCTION

0 Introduction,

0 Principal Component Analysis (PCA),

0 Kernel PCA,

0 Canonical Correlation Analysis,

0 Factor Analysis,

0 Multidimensional scaling,

0 Correspondence Analysis

Reference Books

Reference Books

- 0 1. Charu C. Aggarwal , “Data Mining The Text book”, Springer, 2015.
- 0 2. Craig K. Enders, “Applied Missing Data Analysis”, The Guilford Press, 2010.
- 0 3. Inge Koch, “Analysis of Multivariate and High dimensional data”, Cambridge University Press, 2014.
- 0 4. Michael Jambu, “Exploratory and multivariate data analysis”, Academic Press Inc. , 1990.
- 0 5. Charu C. Aggarwal, “Data Classification Algorithms and Applications”, CRC press, 2015

Projects

0 **Team: 5 members per team**

- 0 1. Exploring the data sets for Data Science problems from Kaggle website
- 0 2. Applying exploratory data analysis in the field of biometrics for reliable and robust identification of humans from their personal traits, mainly for security and authentication purposes
- 0 3. Analyze the dataset for Fraud Detection, Customer segmentation etc.
- 0 Note: Students can down load real-time data sets for different Machine Learning Tasks from <https://archive.ics.uci.edu/ml/datasets.html> and <http://sci2s.ugr.es/keel/datasets.php#sub1> and do the projects

Data Analytics

- 0 It's a process of examining data sets to draw conclusion with the aid of specialized systems and softwares
- 0 Data Analytics initiatives can help businesses
 - 0 increase revenues,
 - 0 improve operational efficiency,
 - 0 optimize marketing campaigns,
 - 0 customer service efforts,
 - 0 quick response to emerging market trends,
 - 0 gain information over rivals

Data Analytics

Applications of Data Analytics

- 0 It supports variety of business uses.
- 0 Example1: banks and credit card companies analyze withdrawal and spending patterns to prevent fraud and identity theft.
- 0 Example 2: E-commerce companies and marketing services providers do clickstream analysis to identify website visitors who are more likely to buy a particular product or service based on navigation and page-viewing patterns

Applications of Data Analytics

- 0 Example 3: Healthcare organizations mine patient data to evaluate the effectiveness of treatments for cancer and other diseases

Data Science Domain

Domain	Usage
	Predicting flight delay
	Predicting life time value of a customer Cross selling Up selling
	Disease Prediction Medicine effectiveness

Data Science Domain

Domain	Usage
	Sentiment Analysis Digital Marketing
	Discount Offering Demand Forecasting
	Self driving cars Pilotless aircrafts Drones

Data Analytics Process

- 0 It starts with data collection, in which data scientist identify the information they need for a particular analytics application and then assemble it for use
- 0 Data from different source systems need to be combined via data integration routines, transformed into a common format and loaded into the analytics system such as Hadoop or NoSQL Database

Data Analytics Process

- 0 Alternatively, the collection process may consist of pulling a relevant subset out of a stream of raw data that flows into Hadoop and moving it to a separate partition so that it can be analyzed without affecting the overall data
- 0 The quality issues in the data may affect the accuracy of analytics application. So it is subjected to data profiling and data cleansing.


Data Analytics Process

- 0 This is required to make sure that the information in a data set is consistent, errors and duplicate entries are eliminated
- 0 Data Governance policies are applied to ensure that the data matches the corporate standards
- 0 The data scientist builds an analytical model using predictive modelling tools and or other programming languages such as python, Scala, R etc..

Data Analytics Process

- 0 The model is initially run against a partial data set to test its accuracy and tested as a process known as training.
- 0 Finally the model is run in production mode against the full data set
- 0 At last the results generated by analytical models are communicated to business executives and other end users to aid in their decision making using visualization techniques

Machine Learning Methodology



Identify Problem: Define the problem statement and the end outcome expected

Gather data: Identify, Collect and prepare data available for the use case

Perform EDA , Build Features: Explore , Analyze and Study the length and depth of data

Build Machine Learning Models: Train and develop machine learning models for the use case

Productionize solution: Develop data products, deploy automated solutions

Exploratory Data Analysis

- 0 It is the process of studying the data by leveraging various statistical and visualization techniques
- 0 It aims to find patterns and relationships in data.

Data Analytics

- 0 Analyzing big data is the process of
 - 0 examining large data sets
 - 0 to uncover hidden patterns,
 - 0 show changes over time, and
 - 0 confirm or challenge theories

What is Data Analytics?

Contd..

Examples:

0 **BigData in Child Welfare**

- 0 Dr. John Snow was skeptical of the theory that foul air was the cause of a significant cholera outbreak.
- 0 He collected interview data on those infected and identified patterns from the data he collected, concluding that the water pump on the street was the source of the outbreak.
- 0 The pump was turned off and the outbreak stopped.

What is Data Analytics?

Contd..

Examples:

0 **BigData in Healthcare**

- 0 Dr. John Halamka, one of the foremost health care professional in the world shares a personal situation demonstrating how data and analytics can benefit patients and catalyze positive changes in health care delivery.
- 0 His wife, of Asian descent, was diagnosed with stage IIIA breast cancer in December 2011.

What is Data Analytics?

Contd..

- 0 He and his team queried data from all the Harvard hospitals on treatment and outcomes for the last 10,000 Asian females with a tumor similar to his wife's.
- 0 Data revealed a medication that would be most effective for her specific case.
- 0 Halamka credits his wife's full recovery to the findings that he was able to draw from this analysis.

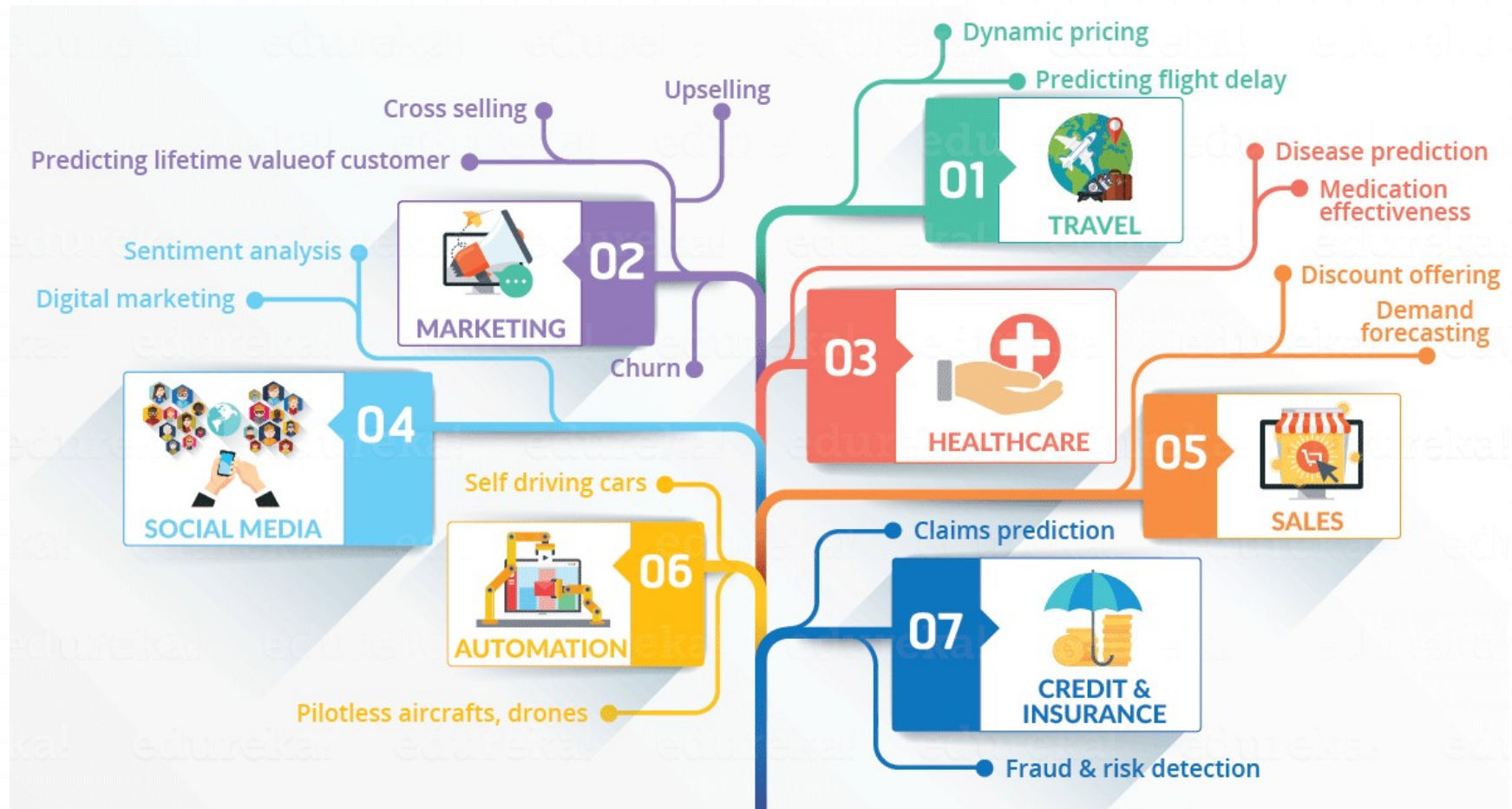
Data Science - Definition

- 0 It is an interdisciplinary field about
 - 0 scientific methods
 - 0 processes, and
 - 0 systems to extract knowledge or
 - 0 insights from data in various forms, either
 - 0 structured or unstructured

Data Science

- 0 Data Science is a field that encompasses anything related to
- 0 data cleansing
- 0 preparation and
- 0 Analysis
- 0 To **get knowledge from the data**

Applications of Data Science



Why we need Data Science

- 0 Tradition data was structure and small in size.
- 0 Business Intelligence Tools was enough for analyzing it
- 0 Current trend, data is either unstructured or semi-structured.

Sources of data

- 0 It is generated from financial logs, multimedia forms, sensors and instruments
- 0 BI tools are not capable of processing huge data
- 0 So advanced analytical tools and algorithms for processing, analyzing and drawing insights

What makes Data Scientist Special?

0 Traditional Data Analyst explains **what is going** on by processing the history of the data.

0 Data Scientist

0 does an Exploratory Data Analysis

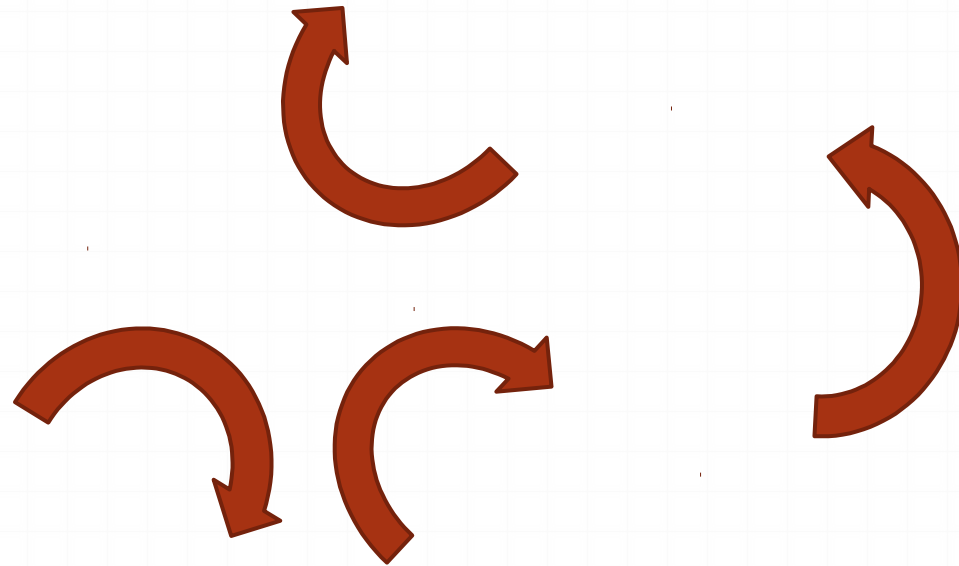
0 Makes use of machine learning algorithms to identify the **occurrence** of a particular **event** in the **future**

Life cycle of Data Science

Myth:

0 Directly go for data collection and analysis

Life cycle of Data Science



Life Cycle of Data Science

Phase1 – Discovery: Frame the business problem and formulate initial hypothesis

- 0 Understand:

- 0 Specifications

- 0 Requirements

- 0 Priorities

- 0 Required budget

Life Cycle of Data Science

Phase1 – Discovery: Frame the business problem and formulate initial hypothesis

- 0 Ability to ask the right question
- 0 Assess required resources present in terms of
 - 0 People, Technology, Time , Data to support the project

Life Cycle of Data Science

Phase1 – Discovery: Frame the business problem and formulate initial hypothesis

- 0 Identify people and key stakeholders

- 0 Identify the data sources

Life Cycle of Data Science

Phase2 - Data Preparation:

- 0 Prepare Analytic Sandbox/workspace
- 0 Perform ELT(Extract, Load, Transform)
- 0 Learning about the data
- 0 Data Conditioning
- 0 Survey & Visualize

Life Cycle of Data Science

Phase2 - Data Preparation:

It includes steps:

- 0 to explore,
- 0 preprocess and
- 0 condition data

Most labour intensive step in the analytical life cycle

Its generally the most iterative phase

Life Cycle of Data Science

Phase2 - Data Preparation:

Require analytics sandbox to perform analytics for the entire duration of the project.



Life Cycle of Data Science

- 0 Phase2 - Data Preparation: Require analytics sandbox to perform analytics for the entire duration of the project.
- 0 A secure private cloud is set up in an organization behind the firewall
- 0 It is interconnected to the IT production environment for data sharing
- 0 The users create simple functions using a language built for specific purpose
- 0 Sandbox gives the analyst a safe place to build models and run experiments to analyse how the company works under such conditions

Life Cycle of Data Science

Phase2 - Data Preparation:

Life Cycle of Data Science

Phase2 - Data Preparation:

- 0 Perform ELT to get the data into sandbox:
 - 0 Extract
 - 0 Load
 - 0 Transform
- 0 Data cleaning, transformation and visualization can be done using R which will help to
 - 0 Spot outliers
 - 0 Establish a relationship between the variables

Life Cycle of Data Science

Phase2 - Data Preparation: Learning about the data

Its important to get familiar with the data

- 0 List your data sources
- 0 List what is needed vs. what is available
- 0 Highlight gaps - Identifies data not currently available
- 0 Identifies the data outside the organization that might be useful

Life Cycle of Data Science

Phase2 - Data Preparation: Data Conditioning

Clean and normalize the data

Discern what you keep and what you discard

Life Cycle of Data Science

Phase2 - Data Preparation: Survey & Visualize

- 0 Overview, filter then maintain data of interest
- 0 Descriptive statistics
- 0 Use data visualization tools to gain an overview of data
 - 0 Does the data contains unexpected values?

Life Cycle of Data Science

□ Learning about the Data: Sample Dataset Inventory

Dataset	Data Available and Accessible	Data Available, but not Accessible	Data to Collect	Data to Obtain from Third Party Sources
Products shipped	●			
Product Financials		●		
Product Call Center Data		●		
Live Product Feedback Surveys			●	
Product Sentiment from Social Media				●

Life Cycle of Data Science

Phase3 – Model Planning:

- 0 Data Exploration

- 0 Variable Selection

- 0 Model Selection

Life Cycle of Data Science

Phase3 – Model Planning: Data Exploration

- 0 Assess the data to understand the relationship between the variables
- 0 Assess the structure of the data – this dictates the tools and analytic techniques for the next phase

Life Cycle of Data Science

Phase3 – Model Planning: Variable Selection

- 0 Explore the data to select variables and methods using visualization tools
- 0 Inputs from stakeholders and domain experts
- 0 Iterative testing to confirm the most significant variables

Life Cycle of Data Science

Phase3 – Model Planning: Model Selection

- 0 Choose an analytical technique based on the end goal of the project

Life Cycle of Data Science

Phase3 – Model Planning:

- 0 Determine the methods and techniques to draw the relationships between the variables
- 0 relationships will set the base for choosing the algorithms which you will implement in the next phase.
- 0 insights into the nature of your data and have decided the algorithms to be used

Life Cycle of Data Science

Phase3 – Model Planning:

The problem to Solve	The category of techniques
I want to group items by similarity. I want to find structures(commonalities) in the data	Clustering
I want to discover relationships between actions or items	Association Rules
I want to determine relationship between the outcome and the input variables	Regression
I want to assign (known) labels to objects	Classification

Life Cycle of Data Science

Phase3 – Model Planning:

The problem to Solve	The category of techniques
I want to find the structure in a temporal process. I want to forecast the behavior of a temporal process	Time series analysis
I want to analyze my text data	Text analysis

Life Cycle of Data Science

Phase4 – Model Building:

- 0develop datasets for training and testing purposes.
- 0consider whether your existing tools will suffice for running the models
- 0analyze various learning techniques like classification, association and clustering to build the model.

Life Cycle of Data Science

Phase5 – Operationalize:

- 0 deliver final reports, briefings, code and technical documents
- 0 a pilot project is also implemented in a real-time production environment
- 0 To get a clear picture of the performance and other related constraints on a small scale before full deployment

Life Cycle of Data Science

Phase6- Communicate Results:

- 0 Evaluate the results achieved with the goals mentioned in phase1.
- 0 By communicating the key findings with the stake holder

Case study: Diabetes Prevention

Aim: **Predict the occurrence of diabetes.**

Step1: Data collection based on medical history of the patient

Attributes:

1. ngravid - No. of times pregnant
2. glu - Plasma glucose concentration
- 3.bp - Blood pressure
4. skin - Triceps skinfold thickness
5. bmi - Body mass index
6. ped - Diabetes pedigree function
7. age - Age
8. income - Income

Case study: Diabetes Prevention

Patient id	pregn	glu	bp	skin	bmi	ped	age	income
1	6	148	72	35	33.6	0.627	50	
2	one	85	66	29	26.6	0.351	31	
3	1	97	6600	15	23.2	0.487	22	

Case study: Diabetes Prevention

Step2: Clean and prepare the data for data analysis

1. Data consists of inconsistencies like missing values, blank columns and abrupt values
2. Data cleaning required

Case study: Diabetes Prevention

This data has a lot of inconsistencies.

1. In the column ngravid, 1 is written as 'one'
2. In the column bp, one of the values is 6600 which is a huge value
3. Income column is blank. Does not give sense in contributing to diabetes. It has to be removed from the table.
4. Clean and preprocess the data by removing outliers so that it can be used for analysis.

Case study: Diabetes Prevention

Step3: Analysis

1. load the data into the analytical sandbox
2. apply various statistical functions to know about the number of missing values,
3. Visualize the data to get an idea of distribution of the data

Case study: Diabetes Prevention

Step4: Analysis the problem to fix the algorithm

1. we already have the major attributes for analysis like *npreg*, *bmi*, etc., so we will use supervised learning technique to build a model here.
2. decision tree because it takes all attributes into consideration in one go, like the ones which have a linear relationship as well as those which have a non-linear relationship.
3. linear relationship between *ngravid* and

Case study: Diabetes Prevention

Step5:Run a small pilot project

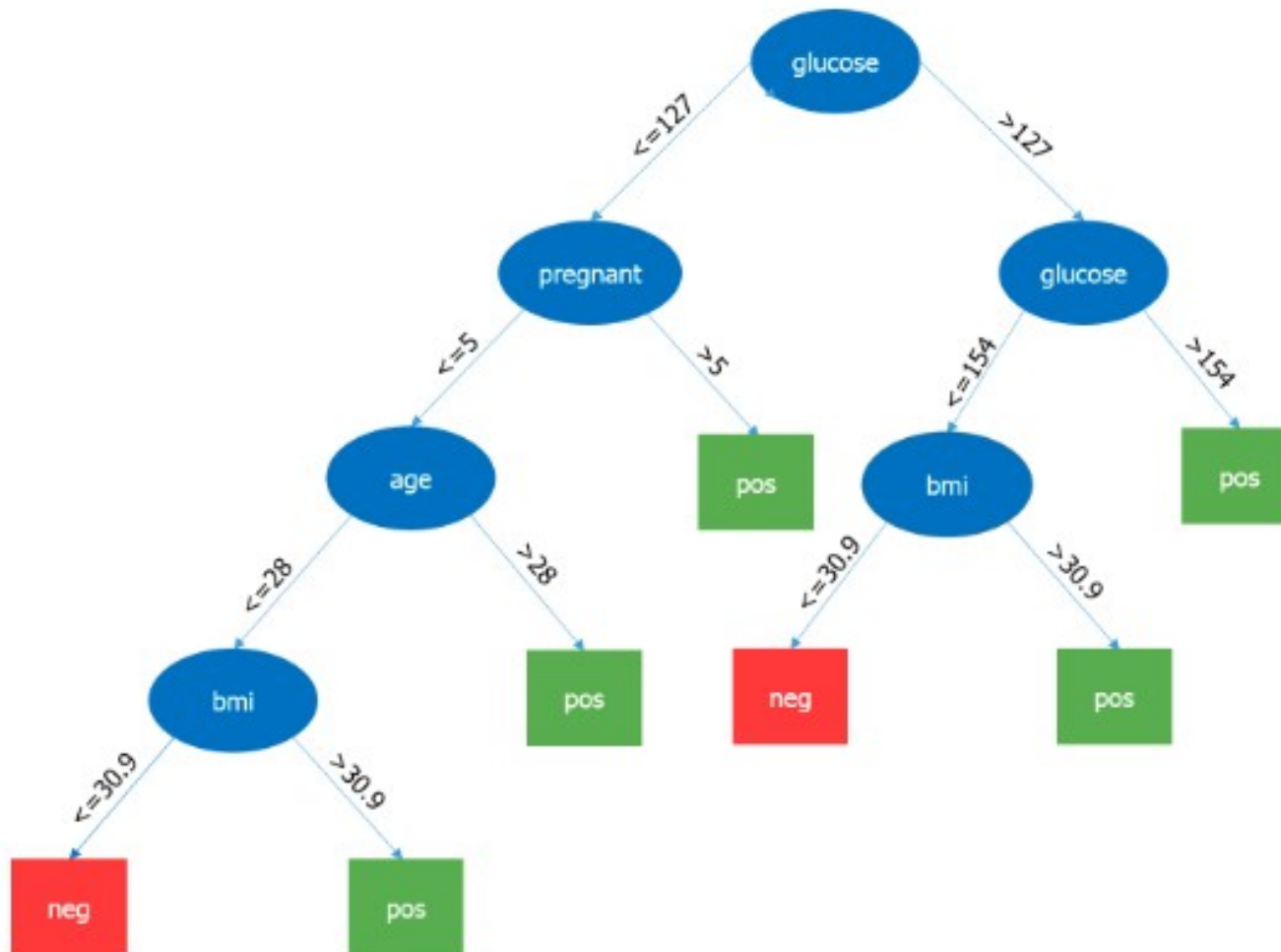
1. to check if our results are appropriate
2. look for performance constraints
3. Re-plan and rebuild the model if the results are not accurate

Case study: Diabetes Prevention

Step6: Full deployment

1. Once we have executed the project successfully, we will share the output for full deployment.

Case study: Diabetes



Applications: Healthcare

0 Example: Reducing hospital readmissions

- 0 In a pilot study, Manhattan's Mount Sinai Medical Center cut readmissions with the help of a computer model that predicts which patients have the highest chances of returning to the hospital.
- 0 The model draws information on factors like disease and past hospital visits from hospital claims data.
- 0 Caregivers give follow-up calls and other assistance to those likely to come back to the hospital.

Applications: Healthcare

contd...

Medical Exams by Bathroom Mirrors:

0 Bathroom mirrors, for example, could read one's skin temperature, pulse and blood pressure, alerting doctors to early indicators of health problems.

Applications: Urban Living

- 0 Urban Informatics:

- 0 System that allows police officers to remotely monitor traffic flows at key junctions via a network of cameras and sensors.

- 0 Passing real-time information about traffic conditions to each other, re-routing themselves to keep traffic moving.

Applications: Crime Prevention

Predictive Policing:

0 uses a constantly-calibrated feed of data on criminal incidents to tell officers where and when future crimes may occur.

Case Study: Marketing Analytics

- 0 Cross sell: It involves the sale of multiple products offered by a single product/service provider to a new or existing customer.
- 0 Up sell: selling higher value products/services to an existing customer.

Case Study: Marketing Analytics

0 Cross Selling:



0 Up Selling.



Next Best Product to Recommend Model Framework

0 This model provide answers to the following process:

0 What – choice of product

0 Whom – selection of customers

0 When – timing

0 How – contact strategy

Next Best Product to Recommend Model Framework



Target variable is response = 1 if the customer has expressed intent to buy, 0 otherwise. Get the response rate.

Response rate is no. of responders divided by customers contacted for the offer.

Next Best Product to Recommend Model Framework

0 Response Model: The following table shows the variables description and relationship to cross-sell response in descending order of importance

Next Best Product to Recommend Model Framework

Variable Description	Relationship with new product response
% of times customer has reacted positively when contacted for an offer	positive
Total credit card limit	positive
Ratio of international spend on the card	positive
Is a premium card holder	positive
Average cash usage on credit card	positive

Next Best Product to Recommend Model Framework

0Implementation:

- 0Based on the response model, a cut-off of the score can be decided
- 0Customers exceeding the cut-off should be considered for marketing.

Machine Learning

- 0 Machine learning is a set of algorithms that train
 - 0 on a data set to make predictions or
 - 0 take actions in order to optimize some systems

Machine Learning Algorithms

Supervised Learning



Dr.G.Malathi, Associate Professor and Coordinator for Image
Processing Research Group, School of Computing Science and
Engineering, VIT University, Chennai Campus

09/26/
2019

Supervised Learning

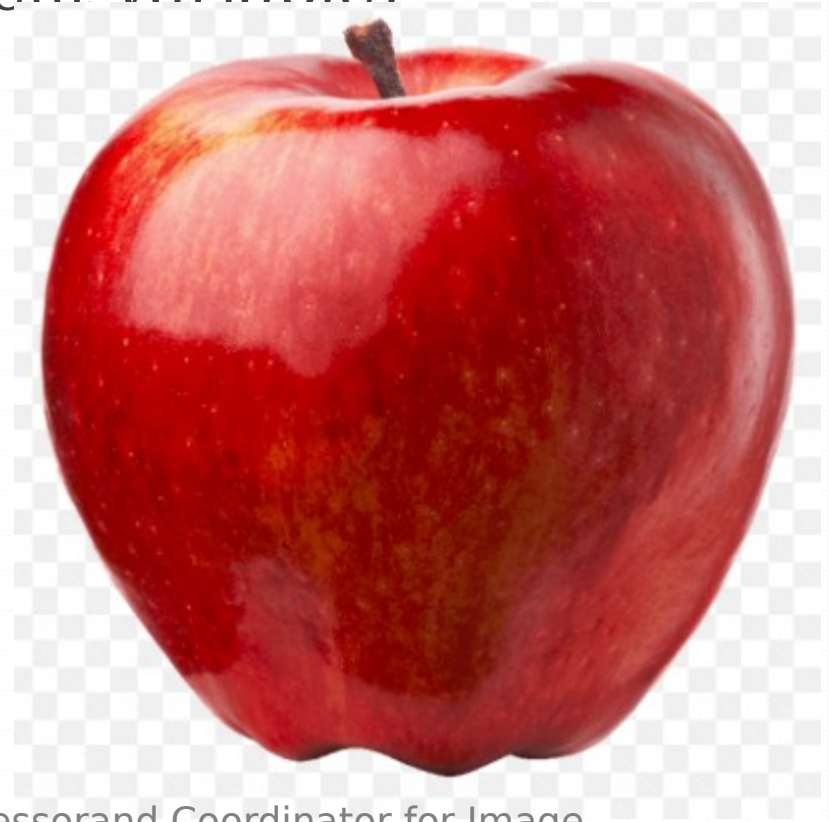


Supervised learning



Supervised Learning

- 0 This algorithm consist of a target / outcome variable (or dependent variable)



Supervised Learning

- 0 It is to be predicted from a given set of predictors (independent variables)



Supervised Learning

- 0 Using these set of variables, we generate a function that map inputs to desired outputs



Supervised learning

The training process continues until the model achieves a desired level of accuracy on the training data



Unsupervised Learning



Unsupervised Learning



we do not have
any target or
outcome variable
to predict /
estimate

Unsupervised Learning

It is used for clustering population in different groups



Unsupervised Learning

0 widely used for segmenting customers in different groups for specific intervention



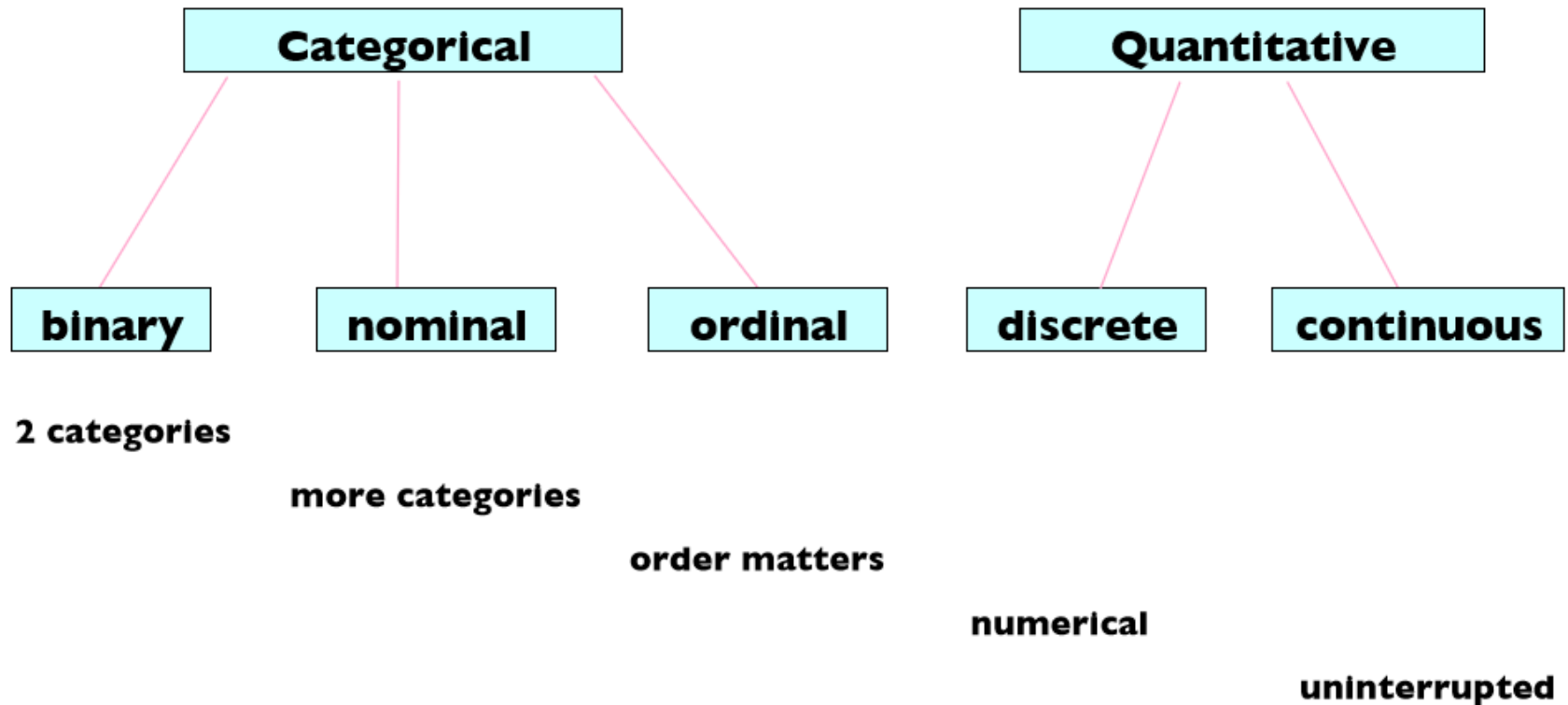
Reinforcement Learning

- 0the machine is trained to make specific decisions
- 0the machine is exposed to an environment where it trains itself continually using trial and error
- 0This machine learns from past experience and tries to capture the best possible knowledge to make accurate business decisions

EDA

- 0 Before making inferences from data it is essential to examine all your variables.
- 0 Why?
 - 0 to catch mistakes
 - 0 to see patterns in the data
 - 0 to find violations of statistical assumptions
 - 0 to generate hypotheses

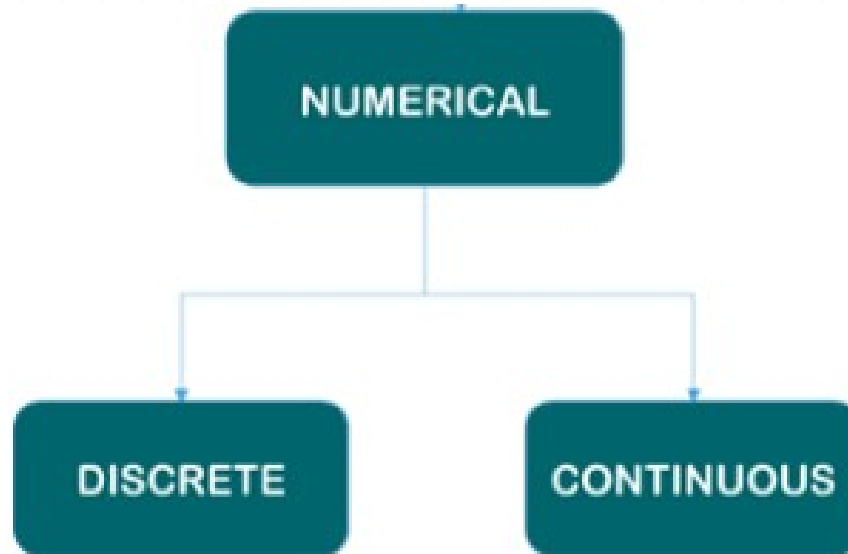
Types of Data



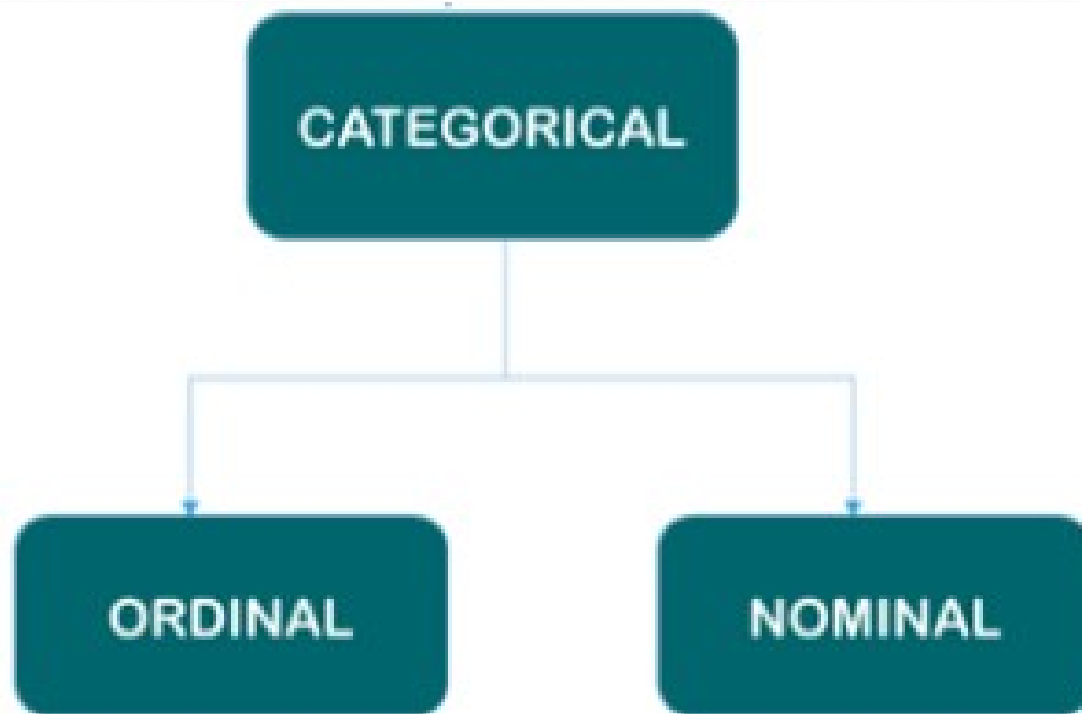
Exploring Data Types



Exploring Data Types



Exploring Data Types



Exploring Data Types

DISCRETE

**Whole Numbers.
Example: No of
students in a class**

CONTINUOUS

**Any value within a
range.
Example: Annual
Income**

Exploring Data Types

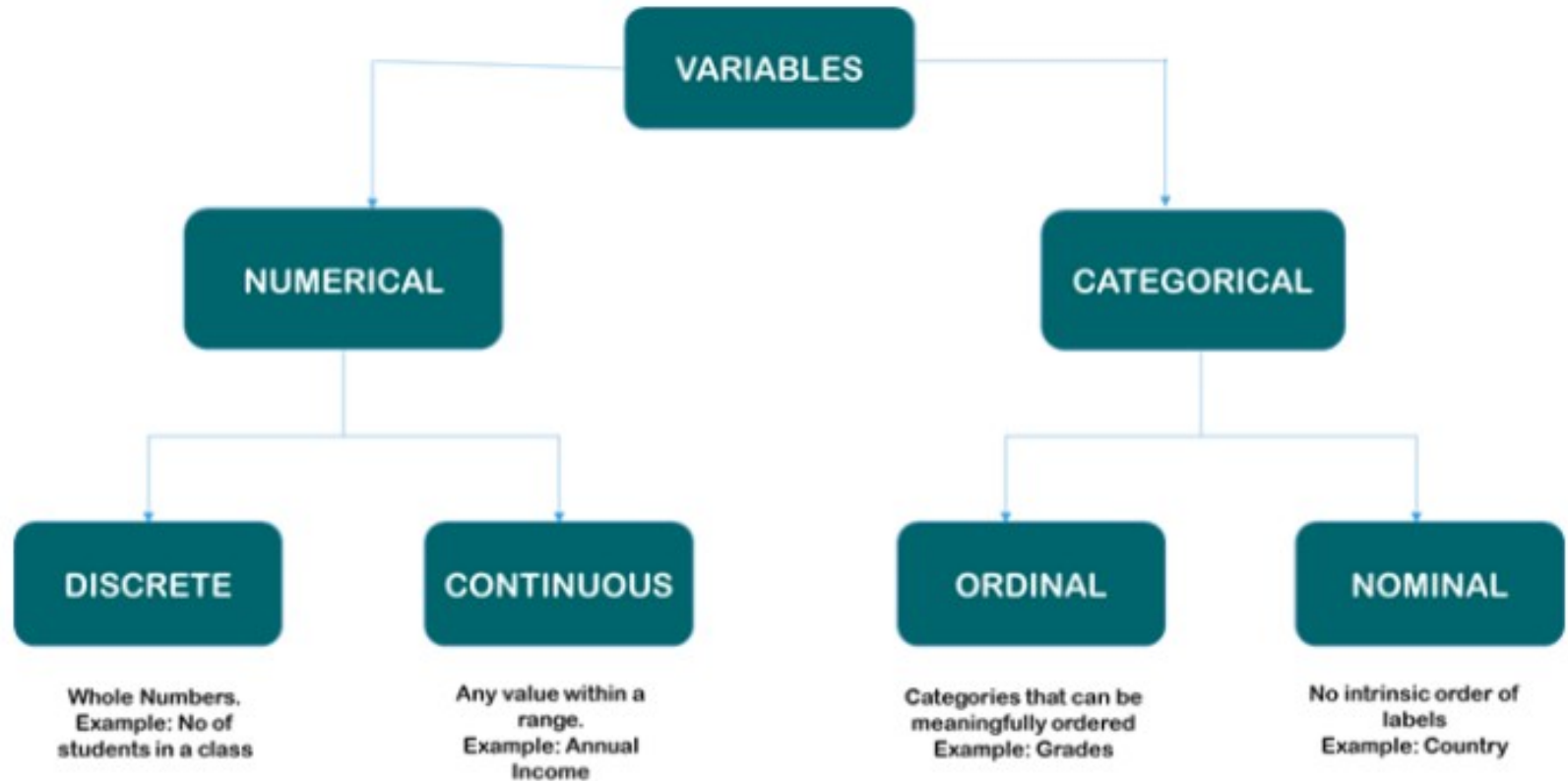
ORDINAL

**Categories that can be
meaningfully ordered
Example: Grades**

NOMINAL

**No intrinsic order of
labels
Example: Country**

Classification of Variables



Dimensionality of Data Sets

- 0 Univariate: Measurement made on one variable per subject
- 0 Bivariate: Measurement made on two variables per subject
- 0 Multivariate: Measurement made on many variables per subject

Univariate, Bivariate, Multivariate Data and its analysis

- 0 Univariate Data:
- 0 This type of data consists of **only one variable**
- 0 It does not deal with causes or relationships
- 0 the main purpose of the analysis is to describe the data and find patterns that exist within it

Univariate, Bivariate, Multivariate Data and its analysis

0 The example of a univariate data can be height.

Heights (in cm)	164	167.3	170	174.2	178	180	186
----------------------------	------------	--------------	------------	--------------	------------	------------	------------

- 0 Conclusions can be drawn using central tendency measures (mean, median, mode)
- 0 Dispersion or spread of data (range, min, max, quantile, variance and SD)
- 0 Graphical Representation using histogram, pie-chart, bar chart

Univariate, Bivariate, Multivariate Data and its analysis

- 0 Bivariate Data:
- 0 This type of data involves **two different variables**
- 0 the analysis is done to find out the relationship among the two variables
- 0 Example of bivariate data can be temperature and ice cream sales in summer season.
- 0 one of these variables is independent while the other is dependent.

Univariate, Bivariate, Multivariate Data and its analysis

- 0 the relationship is visible from the table that temperature and sales are directly proportional to each other
- 0 as the temperature increases, the sales also increase

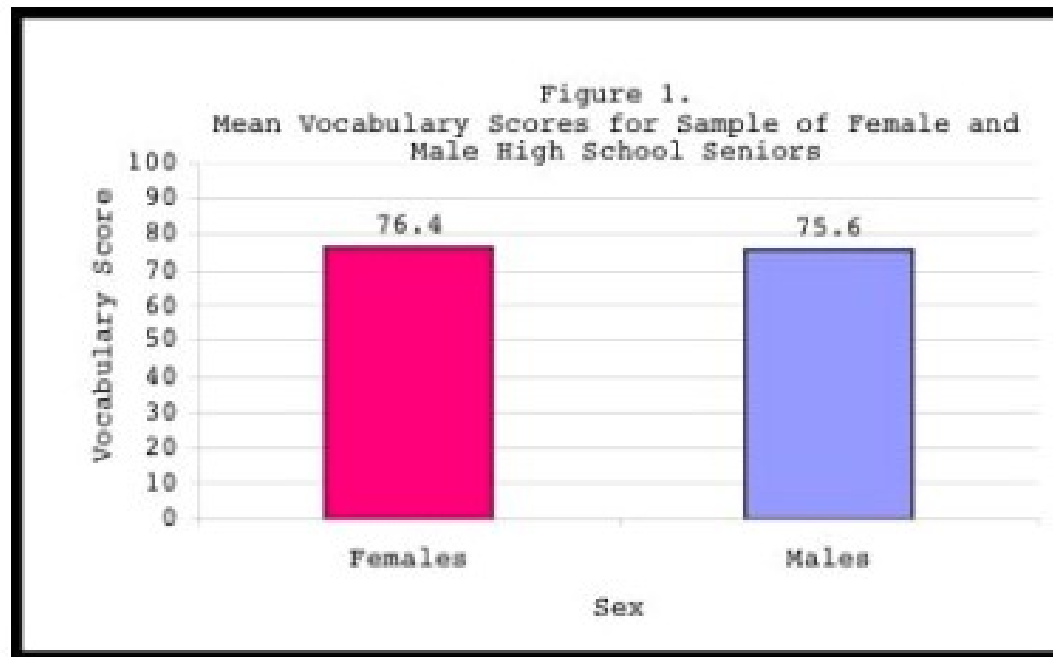
TEMPERATURE(IN CELSIUS)	ICE CREAM SALES
20	2000
25	2500
35	5000
43	7800

Bivariate Analysis

- 0 Relationship between two variables
- 0 Independent variable and dependent variable
- 0 Independent variable: are not affected by anything
- 0 Dependent variable: That can be changed by the outside factors
- 0 Control variable: It may alter either a dependant variable or independent variable

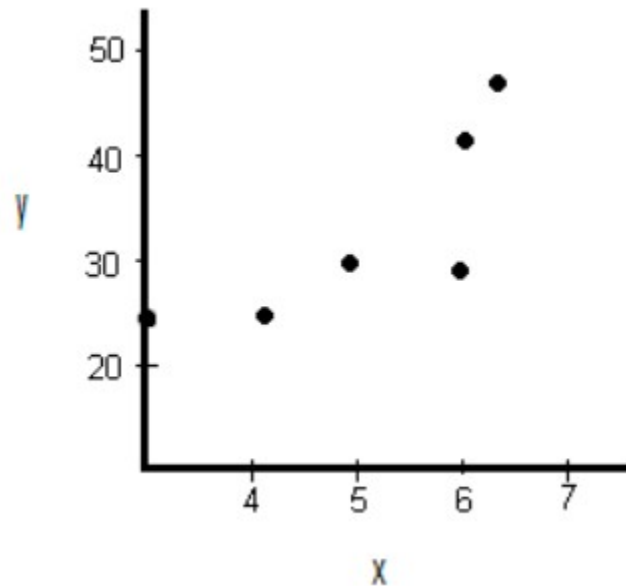
Bivariate Analysis - Example

- 0 Consider the given graph where gender is independent variable and mean vocabulary score is a dependant variable
- 0 the mean vocabulary scores depend on the independent variable. The dependent variable is male or female



Types of Bivariate Analysis

- 0 Scatter Plots: These give you a visual idea of the pattern that your variables follow.



Types of Bivariate Analysis

- 0 Correlation coefficient
- 0 This coefficient tells you if the variables are related
- 0 zero means they aren't correlated
- 0 while a 1 (either positive or negative) means that the variables are perfectly correlated
- 0 Example:
- 0 P

$$r = \frac{\text{mean}(XY) - \text{mean}(X) \times \text{mean}(Y)}{\text{SD}(X) \times \text{SD}(Y)}$$

Attribute Selection- Example

X	Y
10	30
15	45
20	60
25	65
30	80

x	y	xy	x ²	y ²
10	30			
15	45			
20	60			
25	65			
30	80			

$$= [(Xy/n) - (x/n)*(y/n)]/\text{sqrt}(((x^2/n) - \text{sqr}(x/n)) * ((y^2/n) - \text{sqr}(y/n)))$$

Attribute Selection - Example

- 0 If correlation is ≤ 0.5 it is interpreted as very low.
- 0 If correlation is between 0.51 to 0.79 it is interpreted as low.
- 0 If correlation is between 0.80 to 0.89 it is interpreted as moderate.
- 0 If correlation is ≥ 0.90 it is interpreted as high.
- 0 Inference: Variables x,y are highly correlated then any one can be removed

Univariate, Bivariate, Multivariate Data and its analysis

- 0 Multivariate Data:
- 0 involves **three or more variables**
- 0 Example of this type of data is suppose an advertiser wants to compare the popularity of four advertisements on a website, then their click rates could be measured for both men and women and relationships between variables can then be examined.
- 0 It contains more than one dependent variable
- 0 Technique: PCA

Numerical Summaries of Data

- 0 Central Tendency measures: They are computed to give a “center” around which the measurements in the data are distributed.
- 0 Variation or Variability measures. They describe “data spread” or how far away the measurements are from the center.
- 0 Relative Standing measures. They describe the relative position of specific measurements in the data.

Estimates of Location

- 0 Variables with measured or count data might have thousands of distinct values.
- 0 A basic step in exploring your data is getting a
 - 0 “typical value” for each feature (variable):
 - 0 an estimate of where most of the data is located (i.e., its central tendency).

Key terms for estimates of location

0 Mean:

0 The most basic estimate of location is the mean, or *average* value

0 The mean is the sum of all the values divided by the number of values

0 Consider the following set of numbers: {3 5 1 2}.

0 The mean is $(3 + 5 + 1 + 2) / 4 = 11 / 4 = 2.75$.

0 Mean or x-bar represents the mean of a sample from a population.

Estimates of location

0 Mean:

To calculate the average \bar{x} of a set of observations, add their value and divide by the number of observations:

$$\bar{x} = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i$$

Estimates of location

- 0 Trimmed mean:
- 0 Trimmed means are widely used, and in many cases, are preferable to use instead of the ordinary mean
- 0 A trimmed mean eliminates the influence of extreme values
- 0 *trimmed mean* is calculated by dropping a fixed number of sorted values at each end and then taking an average of the remaining values.
- 0 the sorted values
by $x(1), x(2), \dots, x(n)$ where $x(1)$ is the smallest value and $x(n)$ the largest, trimmed mean =
$$\frac{\text{sum}(x(2)+x(3)+\dots+x(n-1))}{(n-2)}$$

Estimates of location

- 0 Example: Consider an international diving competition in which 5 judges are employed. The top and bottom scores from five judges are dropped, and the final score is the average of the three remaining judges.
- 0 This makes it difficult for a single judge to manipulate the score, perhaps to favor his country's contestant.

Estimates of location

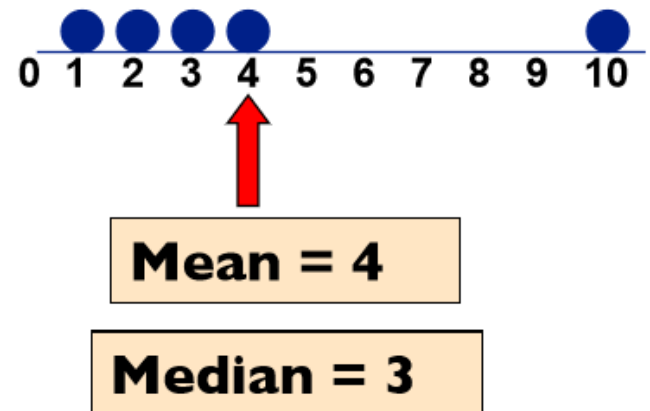
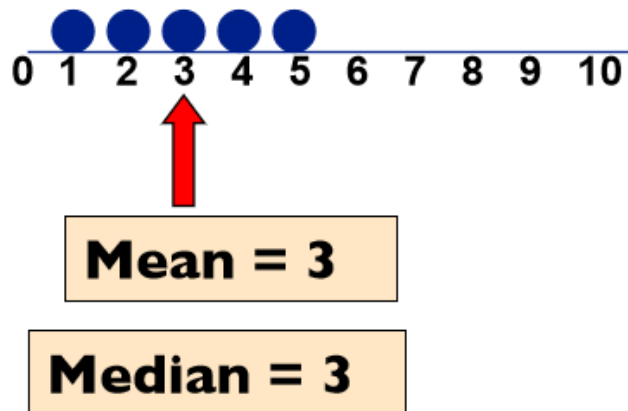
- 0 Weighted mean is calculated by multiplying each data value (x_i) by a weight (w_i) and dividing their sum by the sum of weights

Estimates of location

- 0 Median is the middle number on a sorted list of the data
- 0 If there are an odd number of observations, find the middle value
- 0 - If there are an even number of observations, find the middle two values and average them
- 0 Age of participants: 17 19 21 22 23 23
23 38
- 0 Median = $(22+23)/2 = 22.5$

Which Location is best?

- 0 Mean is best for symmetric distributions without outliers
- 0 Median is useful for skewed distributions or data with outliers
- 0 Assume that 1 to 10 refers to the salary of a person in x place where a celebrity is also in

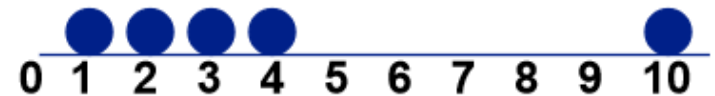


Which Location is best?



Mean = 3

Median = 3



Mean = 4

Median = 3

Outliers

- 0 An outlier is any value that is very distant from the other values in a data set
- 0 When outliers are the result of bad data, the mean will result in a poor estimate of location, while the median will be still be valid
- 0 In anomaly detection, the points of interest are the outliers

Standard Deviation

- 0 *deviations*, between the estimate of location and the observed data.
- 0 For a set of data $\{1, 4, 4\}$, the mean is 3 and the median is 4
- 0 The deviations from the mean are the differences: $1 - 3 = -2$, $4 - 3 = 1$, $4 - 3 = 1$
- 0 The variance is an average of the squared deviations
- 0 standard deviation is the square root of the variance.

Variance

0 Average of squared deviations of values from the mean

$$\hat{\sigma}^2 = \frac{\sum_i^n (x_i - \bar{x})^2}{n - 1}$$

Standard Deviation

0 Standard Deviation is the square root of variance

$$\hat{\sigma} = \sqrt{\frac{\sum_i^n (x_i - \bar{x})^2}{n - 1}}$$

Estimates based on percentile

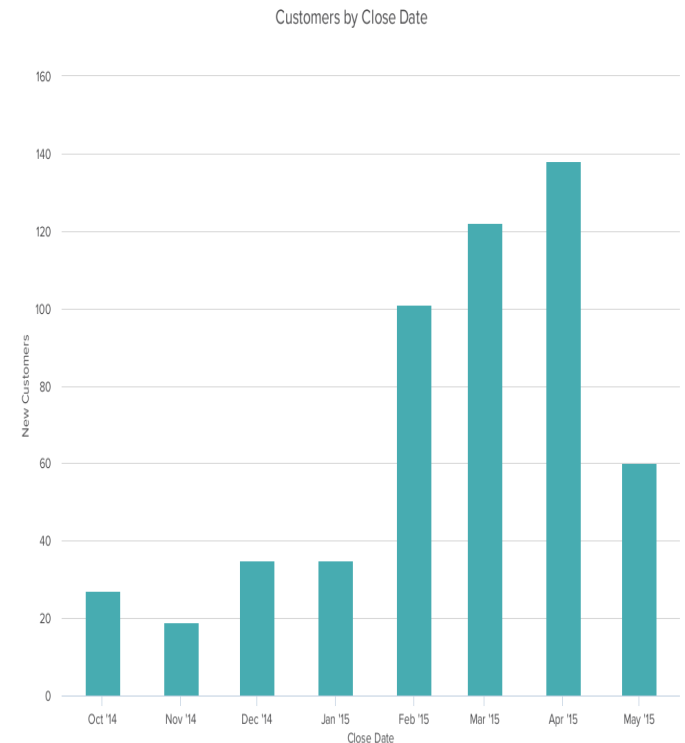
- 0 Statistics based on sorted (ranked) data are referred to as *order statistics*.
- 0 In a data set, the P th percentile is a value such that at least P percent of the values take on this value or less
- 0 For example, to find the 80th percentile, sort the data. Then, starting with the smallest value, proceed 80 percent of the way to the largest value.

IQR

- 0 Difference between the 25th percentile and the 75th percentile, called the *interquartile range* (or IQR)
- 0 Example:
- 0 Consider 3,1,5,3,6,7,2,9
- 0 Sort 1,2,3,3,5,6,7,9
- 0 25th percentile is at 2.5 (1,2,3,**3,5**,6,7,9) median (50%)
- 0 So 25th percentile is (1,**2,3**,3,5,6,7,9) = 2.5
- 0 So 75th percentile is (1,2,3,3,5,**6,7**,9) = 6.5
- 0 Interquantile range is $6.5 - 2.5 = 4$

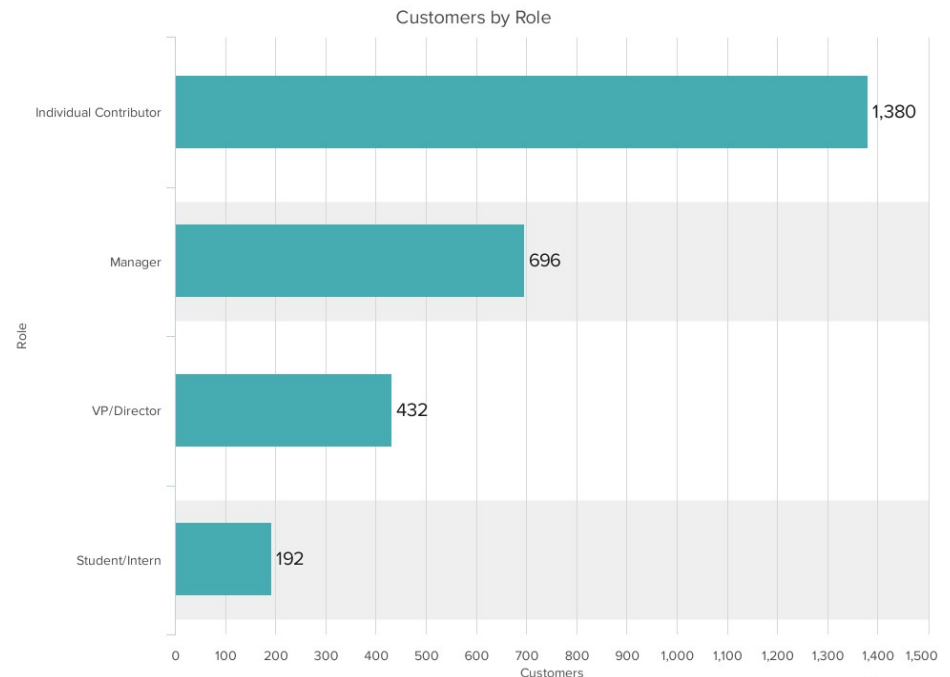
Exploring Data Distribution

- 0 Column Chart
- 0 A column chart is used to show a comparison among different items
- 0 Example: Customers by close



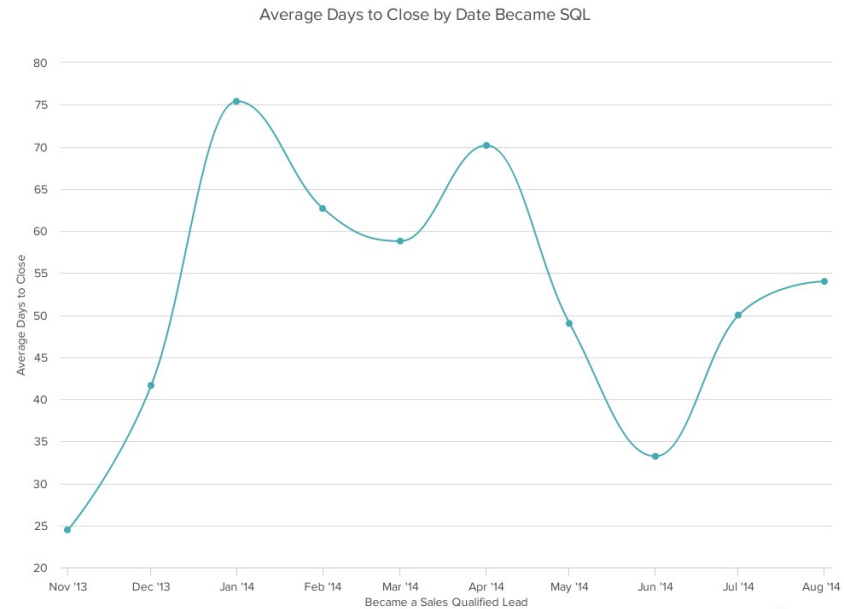
Exploring Data Distribution

- 0 Bar Graph
- 0 basically a horizontal column chart
- 0 To avoid clutter when one data label is long and another is shc



Exploring Data Distribution

- 0 Line Graph
- 0 continuous data set.
- 0 A line graph reveals trends or progress over time

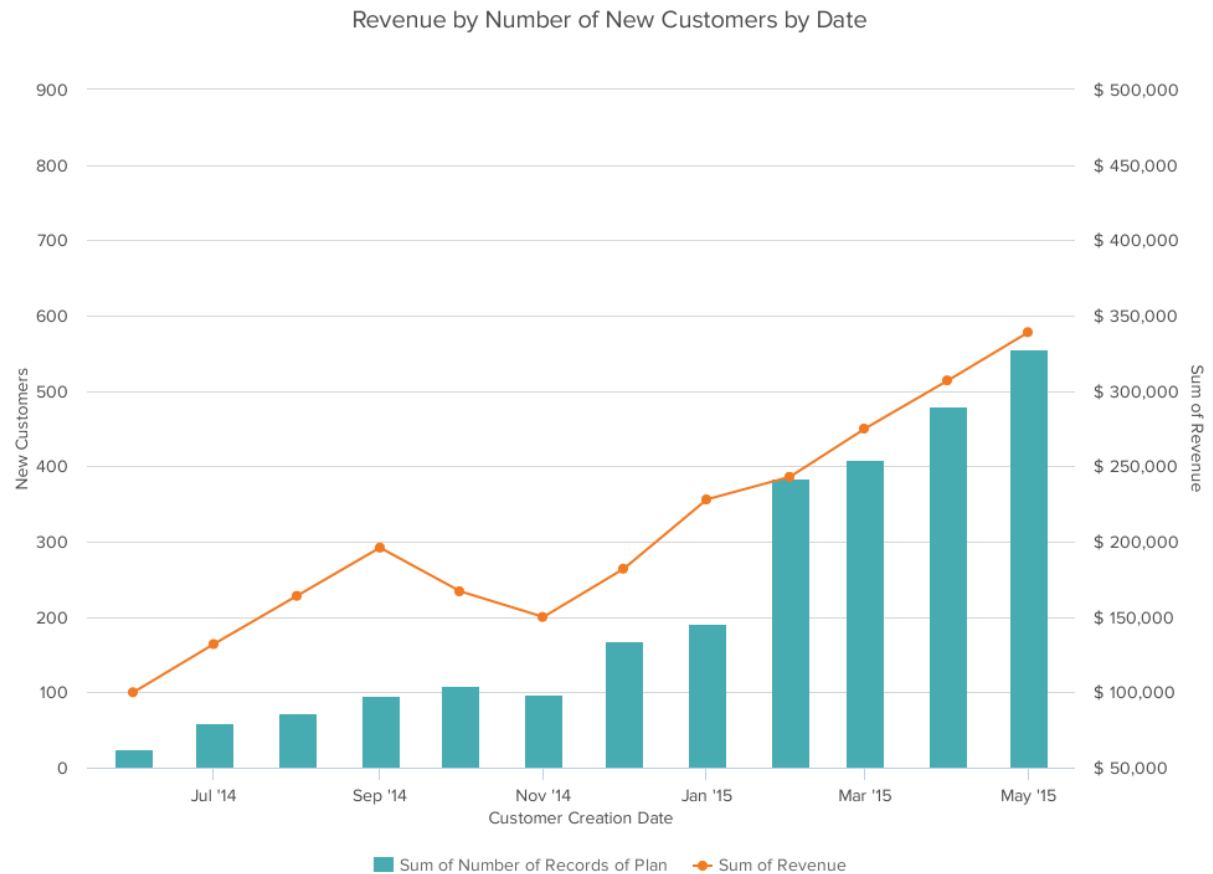


Exploring Data Distribution

- 0 Dual Axis Chart
- 0 allows you to plot data using two y-axes and a shared x-axis
- 0 It's used with three data sets, one of which is based on a continuous set of data and another which is better suited to being grouped by category
- 0 his should be used to visualize a correlation

Exploring Data Distribution

0 Dual Axis Chart



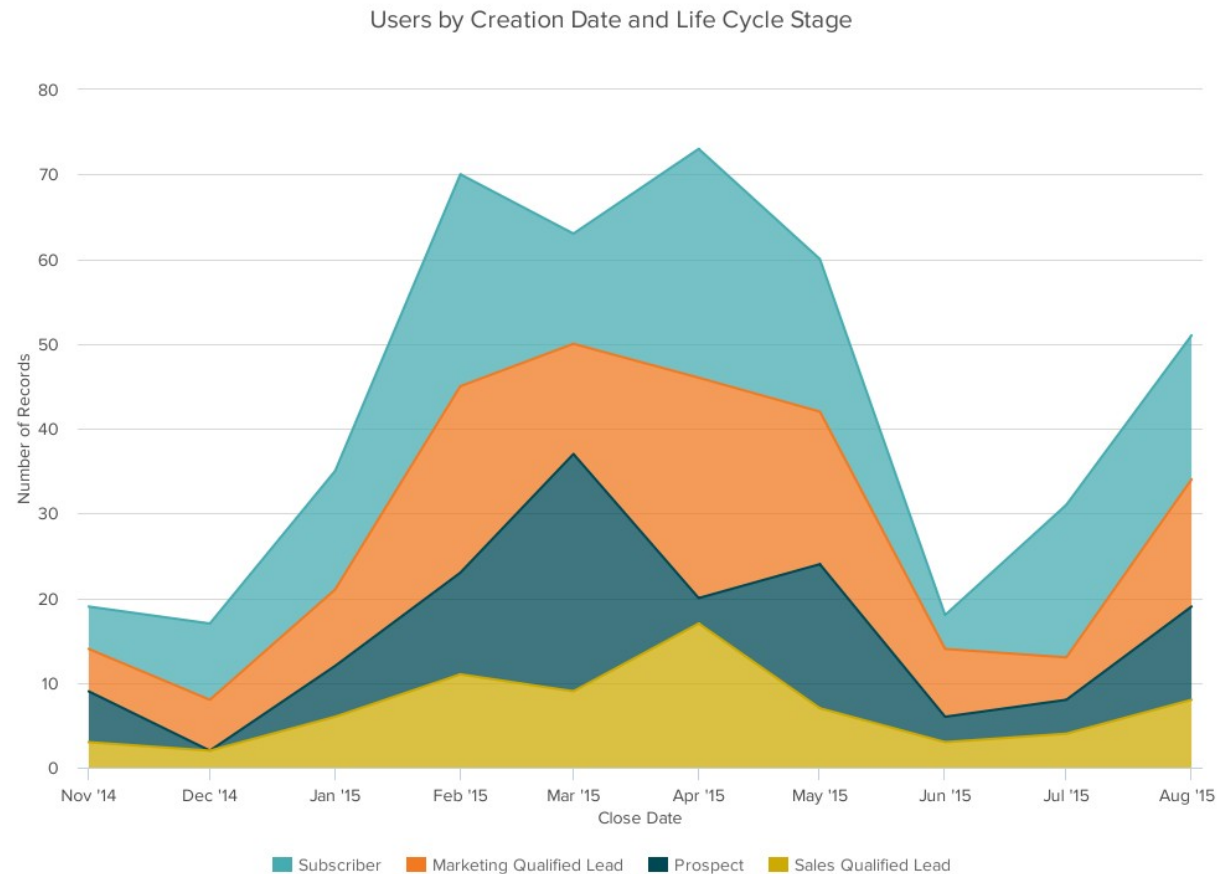
Exploring Data Distribution

0 Area Chart

- 0 An area chart is basically a line chart, but the space between the x-axis and the line is filled with a color or pattern
- 0 It is useful for showing part-to-whole relations
- 0 Example: individual sales reps' contribution to total sales for a year

Exploring Data Distribution

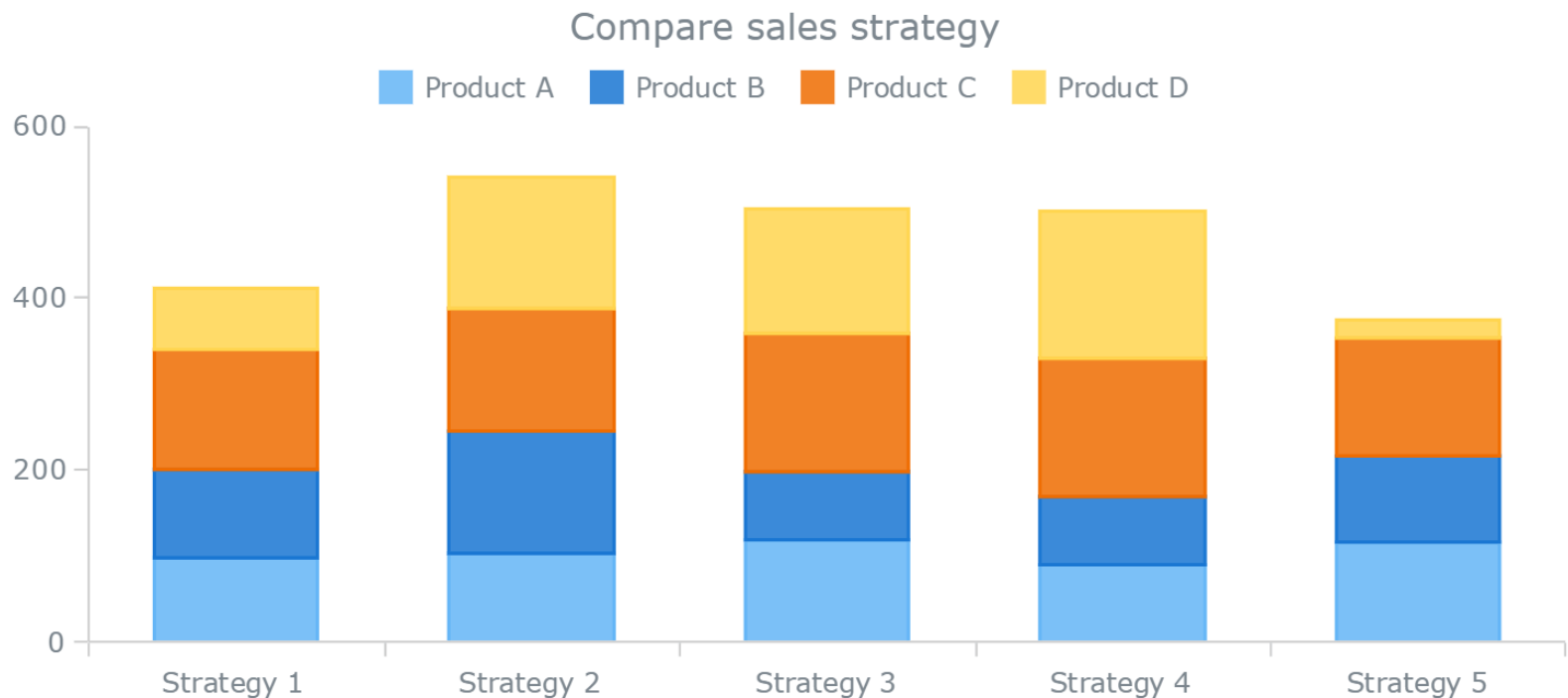
0 Area Char



Exploring Data Distribution

0 Stacked Bar Chart

0 This should be used to compare many different items



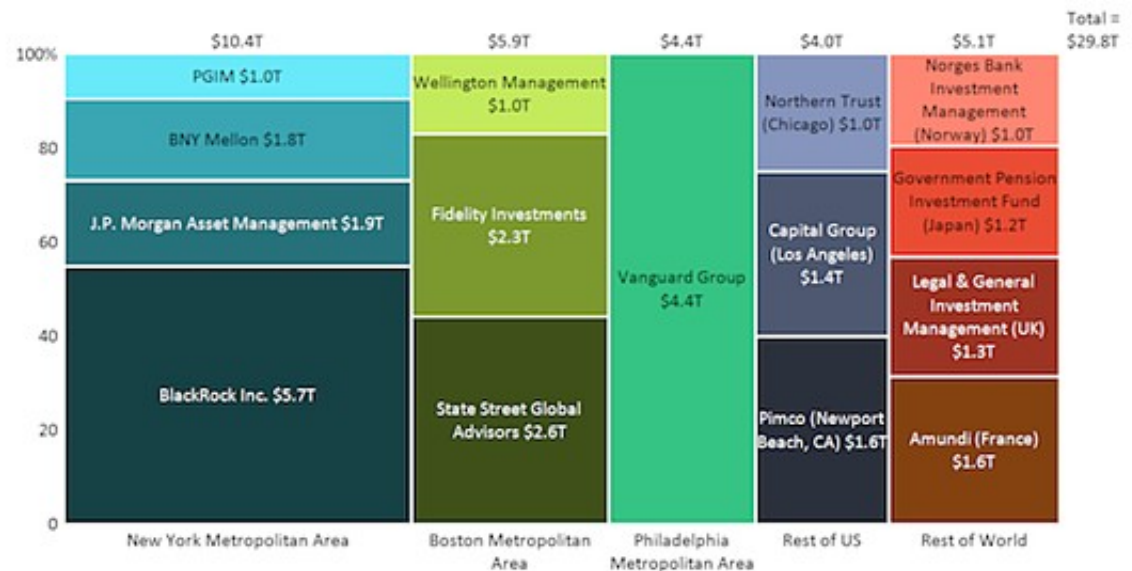
Exploring Data

Introduction

- 0 Mekko Chart
- 0 This type of graph can compare values, measure each one's composition
- 0 To show how your data is distributed across each one.

World's Largest Asset Managers

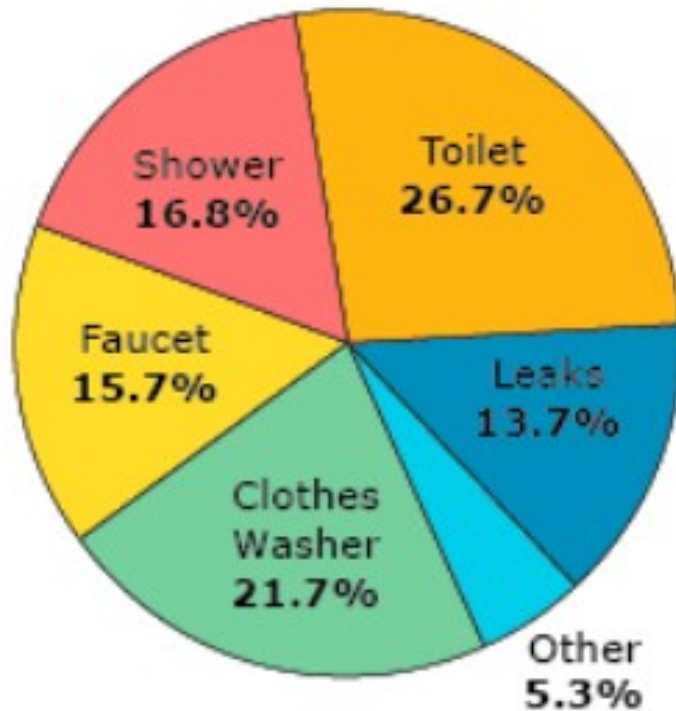
Most of the world's largest asset managers are grouped in the Northeast US. Eight of the 14 firms that manage \$1T or more are in the NY, Boston or Philadelphia areas.



Exploring Data Distribution

- 0 Pie Chart
- 0 shows a static number
- 0 how categories represent part of a whole
- 0 the total sum of all segments needs to equal 100%

How Much Water Do We Use?



Exploring Data Distribution

- 0 Scatter Plot Chart
- 0 Show the relationship between two different variables
- 0 It should be used when there are many different data points, and you want to highlight similarities in the data set.
- 0 This is useful when looking for outliers

Exploring Data Distribution

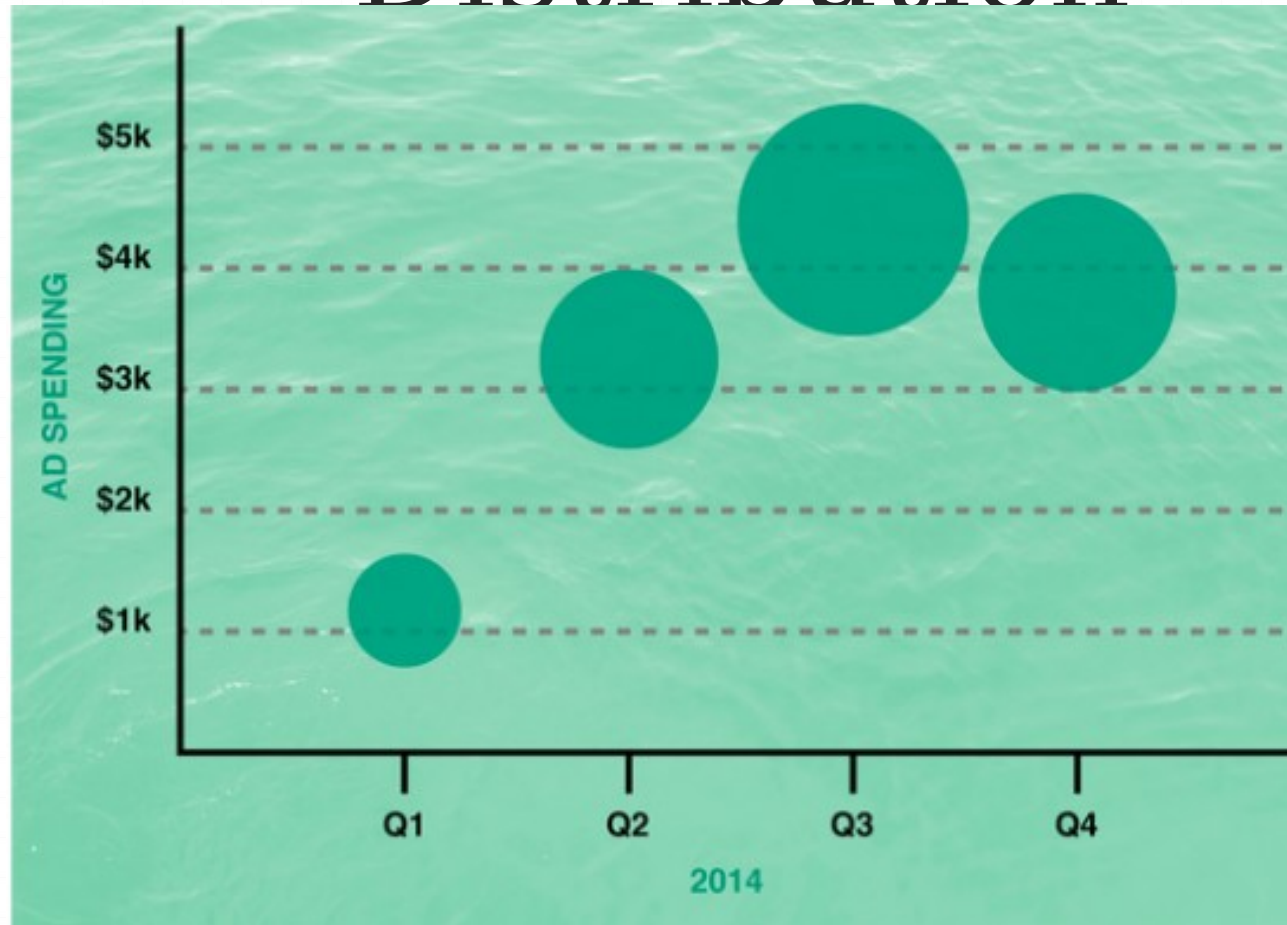


Exploring Data Distribution

0 Bubble Chart

- 0 is similar to a scatter plot in that it can show distribution or relationship.
- 0 There is a third data set, which is indicated by the size of the bubble or circle.

Exploring Data Distribution



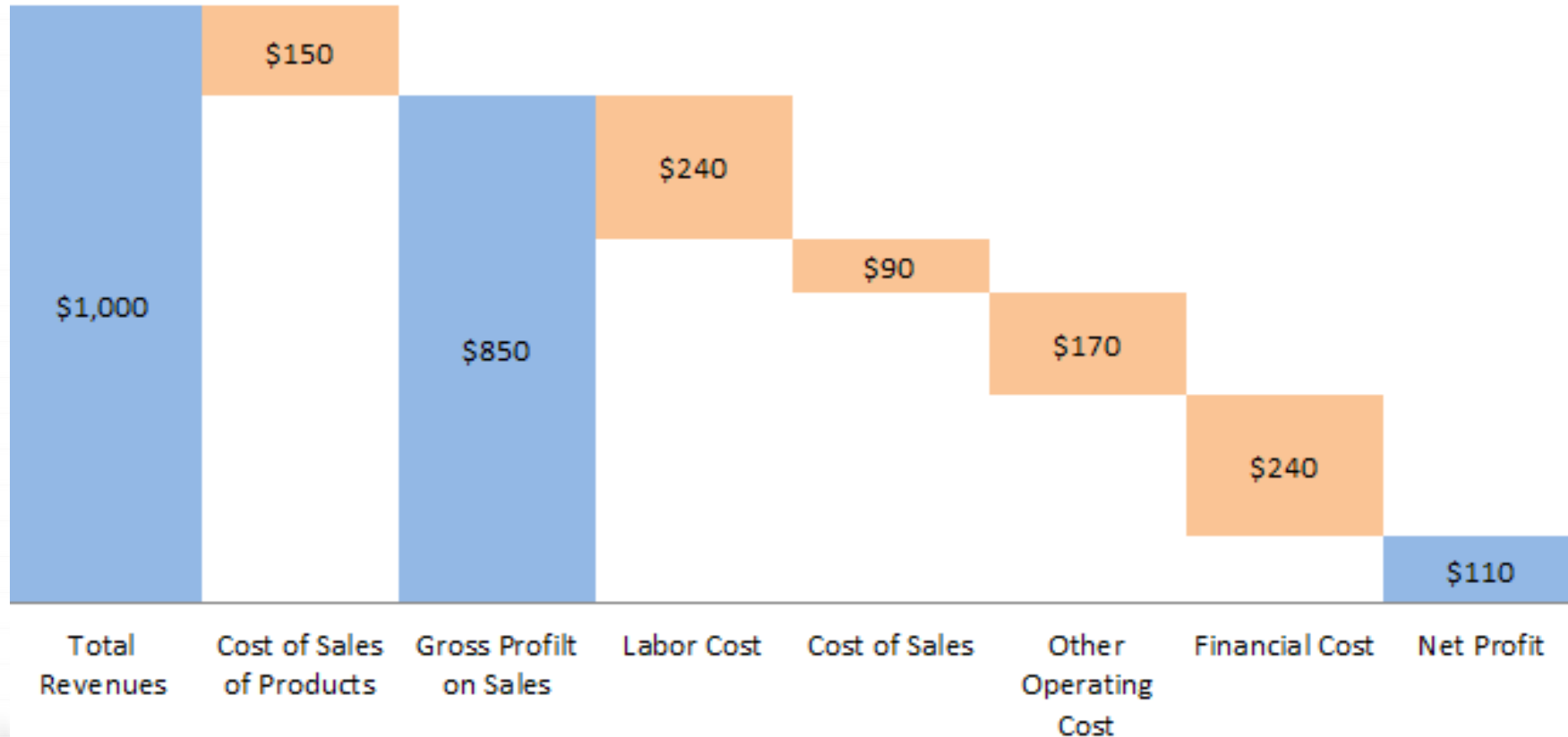
Exploring Data Distribution

0 Waterfall Chart

- 0 A waterfall chart should be used to show how an initial value is affected by intermediate values -- either positive or negative -- and resulted in a final value.

Exploring Data Distribution

Product Profit Analysis



Exploring Data Distribution

0 Funnel Chart

- 0 shows a series of steps and the completion rate for each step
- 0 This can be used to track the sales process or the conversion rate across a series of pages or steps.

Exploring Data Distribution

Sales Analysis - June 2016



Exploring Data Distribution

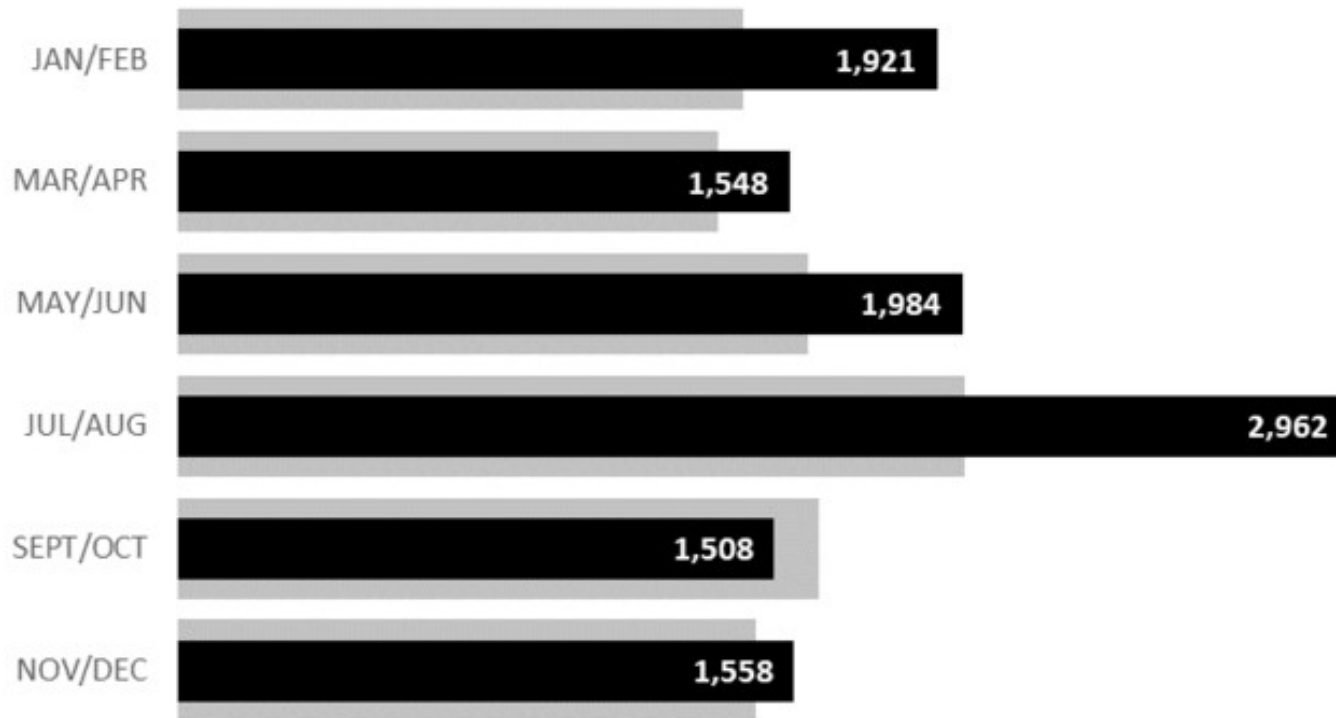
0 Bullet Chart

- 0 A bullet graph reveals progress toward a goal, compares this to another measure,
- 0 provides context in the form of a rating or performance.

Exploring Data Distribution

Water Usage Chart By BI Monthly Billing Cycle
(Usage In Cubic Feet)

■ Previous Year Consumption ■ Current Year Consumption

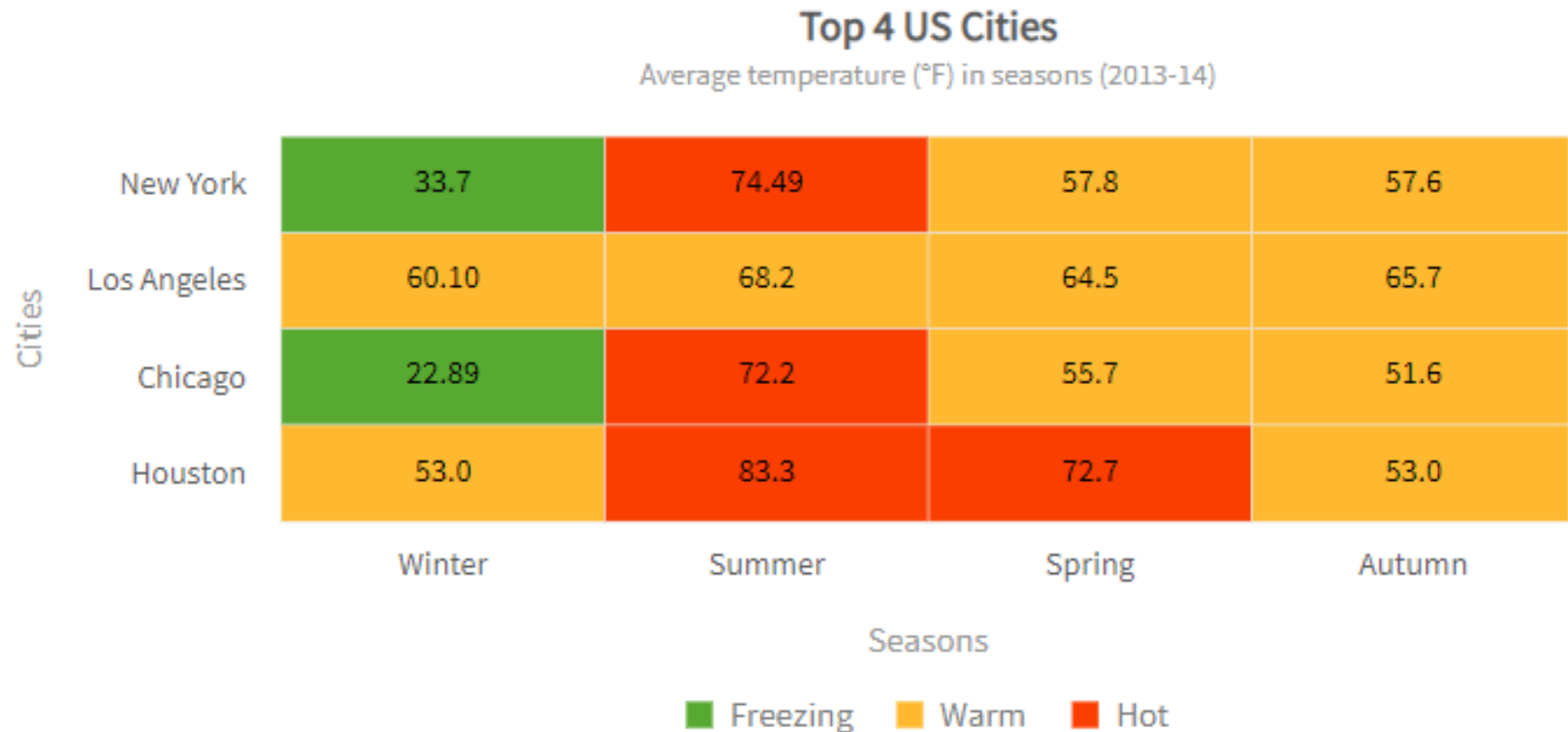


Exploring Data Distribution

0 Heat Map

- 0 A heat map shows the relationship between two items and provides rating information, such as high to low or poor to excellent.
- 0 to plot data like employee attendance records
- 0 you can use colors like red, yellow, blue and green to indicate a bad, average, good, and excellent grade

Exploring Data Distribution



Example1

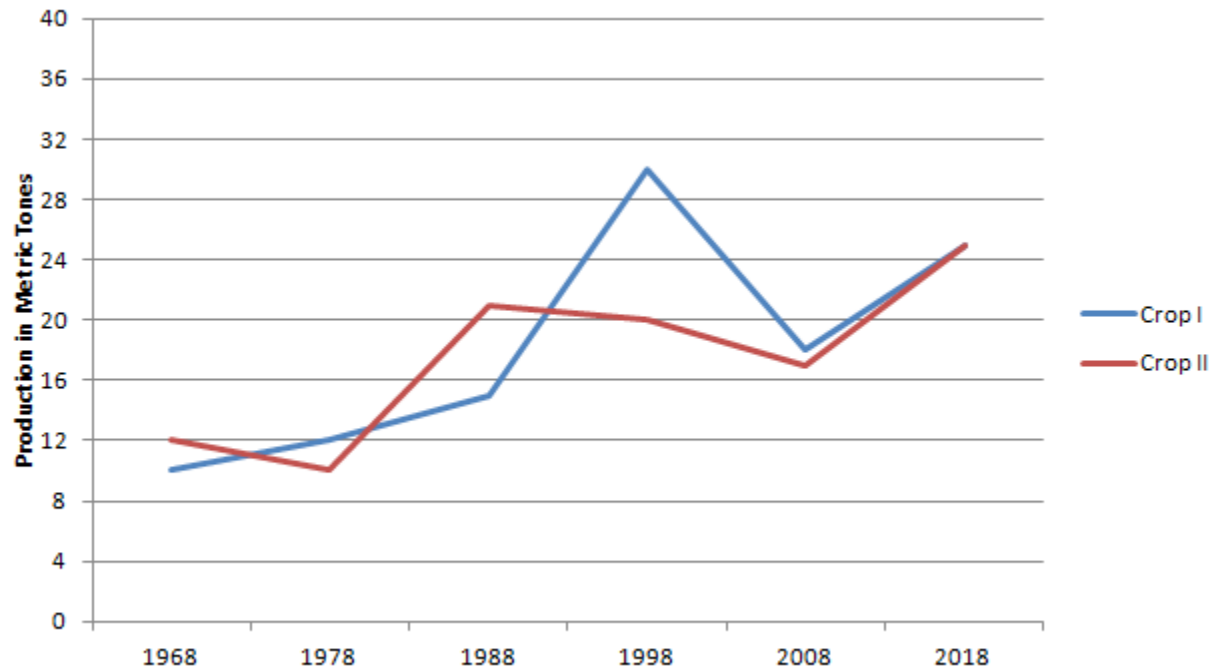
0 Problem: Draw a line graph for the production of two types of crops for the given years.

Production in metric tones

Year	Crop I	Crop II
1968	10	12
1978	12	10
1988	15	21
1998	30	20
2008	18	17
2018	25	25

Example 1

0 Solution: The required graph is



Example 2

0 Draw the histogram for the given data.

Marks	No. of Students
15 - 18	7
19 - 22	12
23 - 26	56
27 - 30	40
31 - 34	11
35 - 38	54
39 - 42	26
43 - 46	37
47 - 50	7
Total	250

Example2

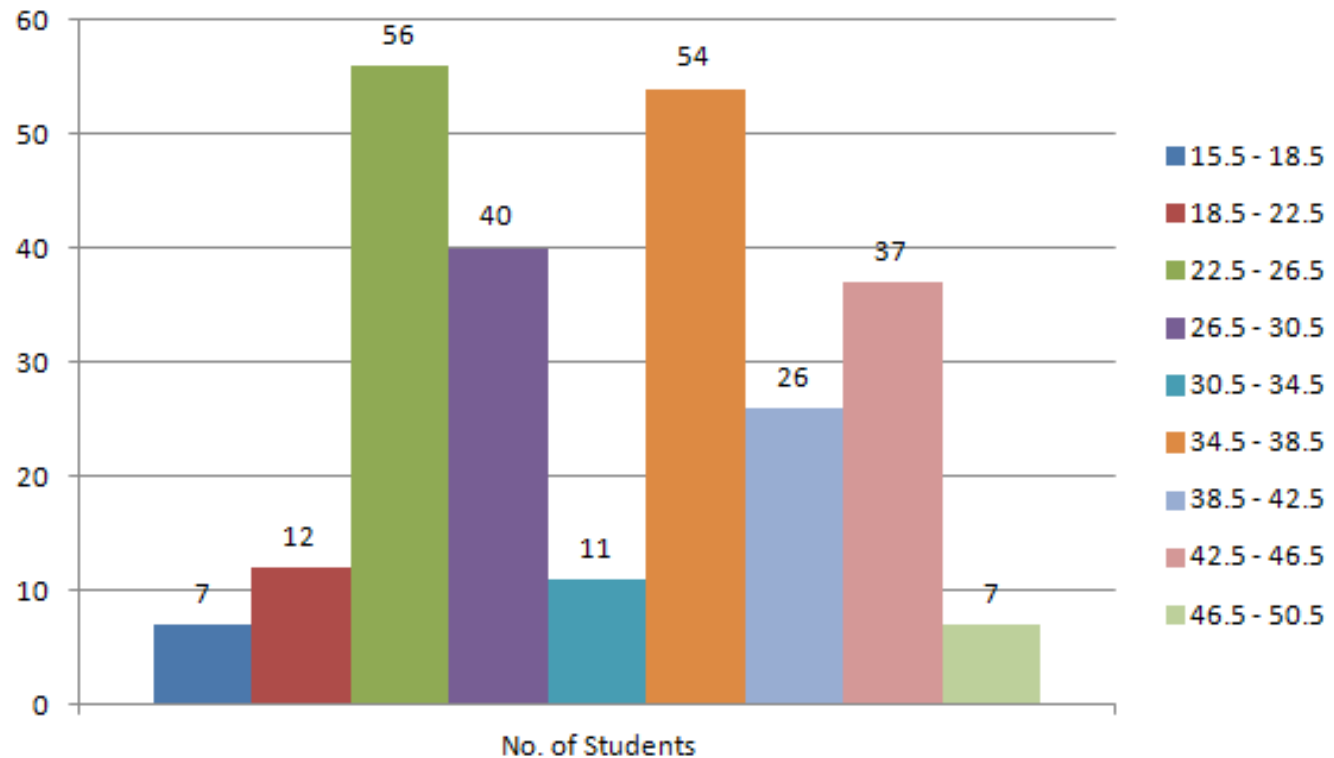
0 Solution: This grouped frequency distribution is not continuous. We need to convert it into a continuous distribution with exclusive type classes. This is done by averaging the difference of the lower limit of one class and the upper limit of the preceding class. Here, $d = \frac{1}{2} (19 - 18) = \frac{1}{2} = 0.5$. We add 0.5 to all the upper limits and we subtract 0.5 from all the lower limits.

Example2

Marks	No. of Students
14.5 - 18.5	7
18.5 - 22.5	12
22.5 - 26.5	56
26.5 - 30.5	40
30.5 - 34.5	11
34.5 - 38.5	54
38.5 - 42.5	26
42.5 - 46.5	37
46.5 - 50.6	7
Total	250

Example2

0 The corresponding histogram is:



Take Home:

Case Study: XYZ MNC Company

- 0 Case: Possibilities of an enquiry being converted into a buyer for sure.
- 0 Software Design Requirements:
- 0 How many people in family
- 0 Per month salary
- 0 Loans if any
- 0 Two wheeler/four wheeler possession if any
- 0 Residing in rental or own
- 0 Has life insurance policies
- 0 Has dogs at home?
- 0 Dog's breed?

T
H
A
N
K

Y
O
U

CONTACT DETAILS

Dr. G. Malathi

Associate Professor and Coordinator for
Image Processing Research Group
School of Computing Sciences and
Engineering

VIT University

Chennai Campus

malathi.g@vit.ac.in

9840833337

Take Home:

Case Study: XYZ MNC Company

- 0 Case: Possibilities of an enquiry being converted into a buyer for sure.
- 0 Knows to drive
- 0 Is married or lone
- 0 Does he/she has children
- 0 Frequent place of visits
- 0 Which type of hotel preferred for stay?
- 0 School in which the kids are studying?
- 0 Restaurants they visit
- 0 Frequency of visit to restaurant

Take Home:

Case Study: XYZ MNC Company

- 0 Case: Possibilities of an enquiry being converted into a buyer for sure.
- 0 Surprising Inferences from the above data obtained from the enquiry:
- 0 If a person has insurance policies, he cares his family a lot and does not spend on higher end cars
- 0 If a family has dogs at home, they are trustable. They will surely pay the lone amount
- 0 The types of hotels or resorts they stay has 5 stars, then they can be pushed to choose higher end if they opt for sedan

Take Home:

Case Study: XYZ MNC Company

- 0 Case: Possibilities of an enquiry being converted into a buyer for sure.
- 0 Surprising Inferences from the above data obtained from the enquiry:
- 0 If the person comes for enquiry more than once, then the chance of conversion to be a buyer is more if negotiation is initiated or offer value added services
- 0 People with kids have 90% chances of buying

Why EDA?

- 0 Detection of mistakes
- 0 Checking of assumptions
- 0 Preliminary selection of appropriate models
- 0 Determining relationship among the variables

Exploratory Data Analysis (EDA)

– An introduction

- 0 Through EDA

- 0 get to know about the data

- 0 distributions,

- 0 Data quality problems

- 0 Outliers

- 0 Correlations and inter-relationships

Steps to Understand, clean and prepare

- 0 Variables identification
- 0 Univariate Analysis
- 0 Bi-variate Analysis
- 0 Missing values treatment
- 0 Outliers treatment

Variable Identification

- 0 Identify
 - 0 Predictor(Input) variables
 - 0 Target(Output) variables
 - 0 Data type of the variables
 - 0 Category of the variables

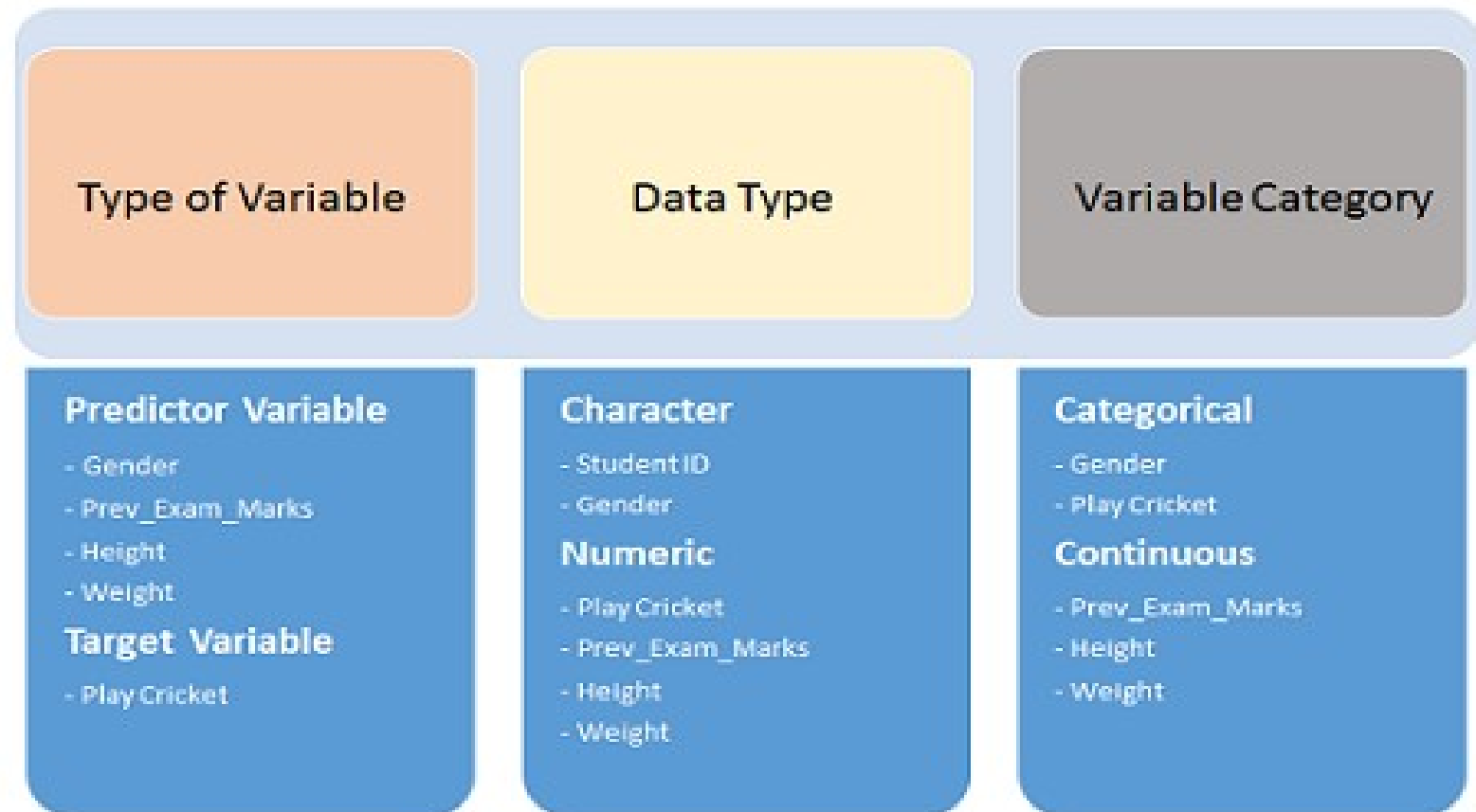
Variable Identification cont...

- 0 To predict, whether the students will play cricket or not. Need to identify predictor variables, target variable, data type of variables and category

Target Variable

Student_ID	Gender	Prev_Exam_Marks	Height (cm)	Weight Caregory (kgs)	Play Cricket
S001	M	65	178	61	1
S002	F	75	174	56	0
S003	M	45	163	62	1
S004	M	57	175	70	0
S005	F	59	162	67	0

Variable Identification cont...



Types of Data

- 0 Categorical data
 - 0 Nominal data
 - 0 Ordinal data
 - 0 Binary data
- 0 Measurement data
 - 0 Discrete
 - 0 Continuous

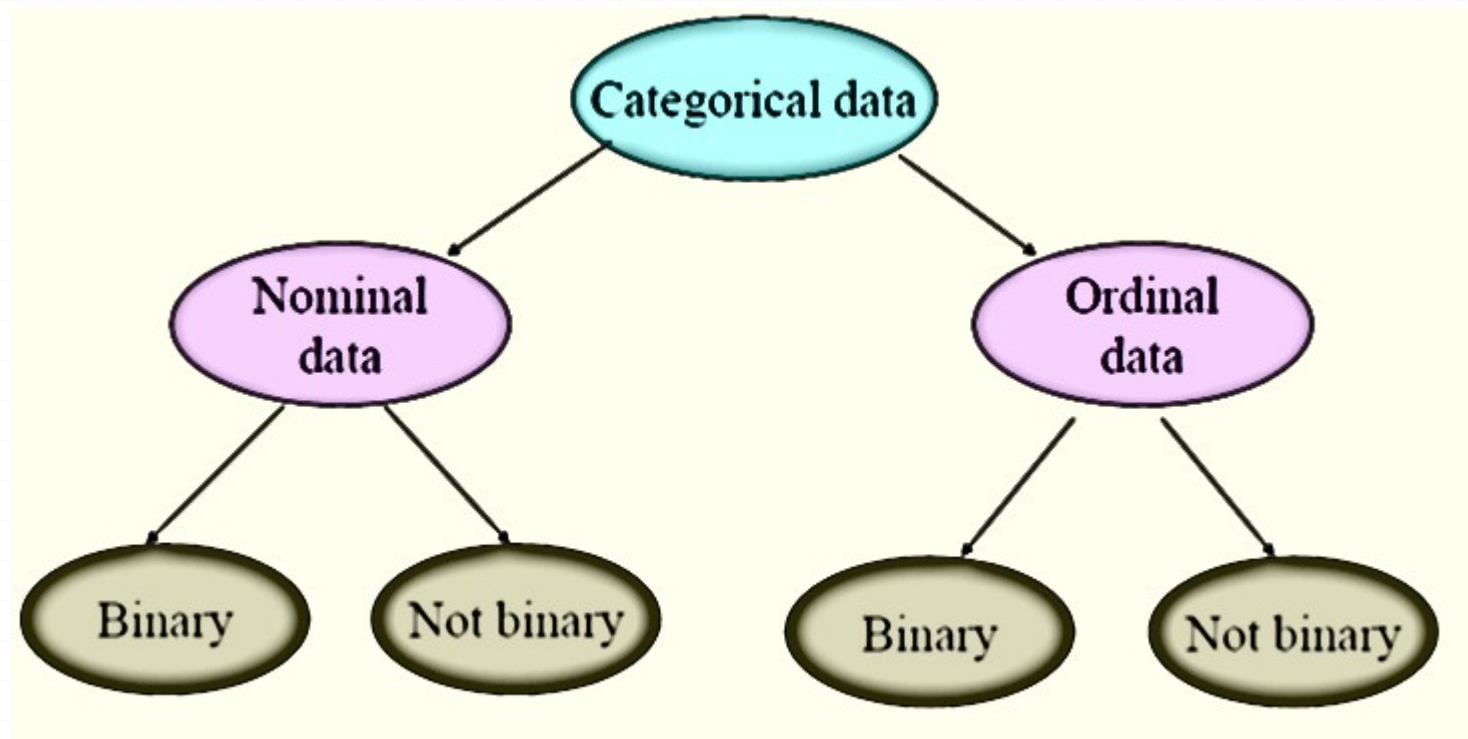
Categorical Data

- 0 The objects being studied are grouped into categories based on some **qualitative trait**
- 0 The resulting data are merely labels or categories

Examples: Categorical Data

- Hair color
 - brown, red, black, white, etc.
- Nationality
 - Indian, non-Indian

Categorical data classified as Nominal, Ordinal, and/or Binary



Nominal Data

- A type of categorical data in which objects fall into **unordered** categories.

0 Hair color

- brown, red, black, etc.

- Race

- Caucasian, African-American, Asian, etc.

- Nationality

- Indian, non-Indian

Ordinal Data

0 A type of categorical data in which **order** is important.

0 Class

- fresh, junior, senior, super senior
- Degree of illness
 - none, mild, moderate, severe, ..., going, going, gone

Binary Data

0 A type of categorical data in which there are **only two categories**.

- Binary data can either be nominal or ordinal.

0 Nationality

- Indian, non-Non Indian

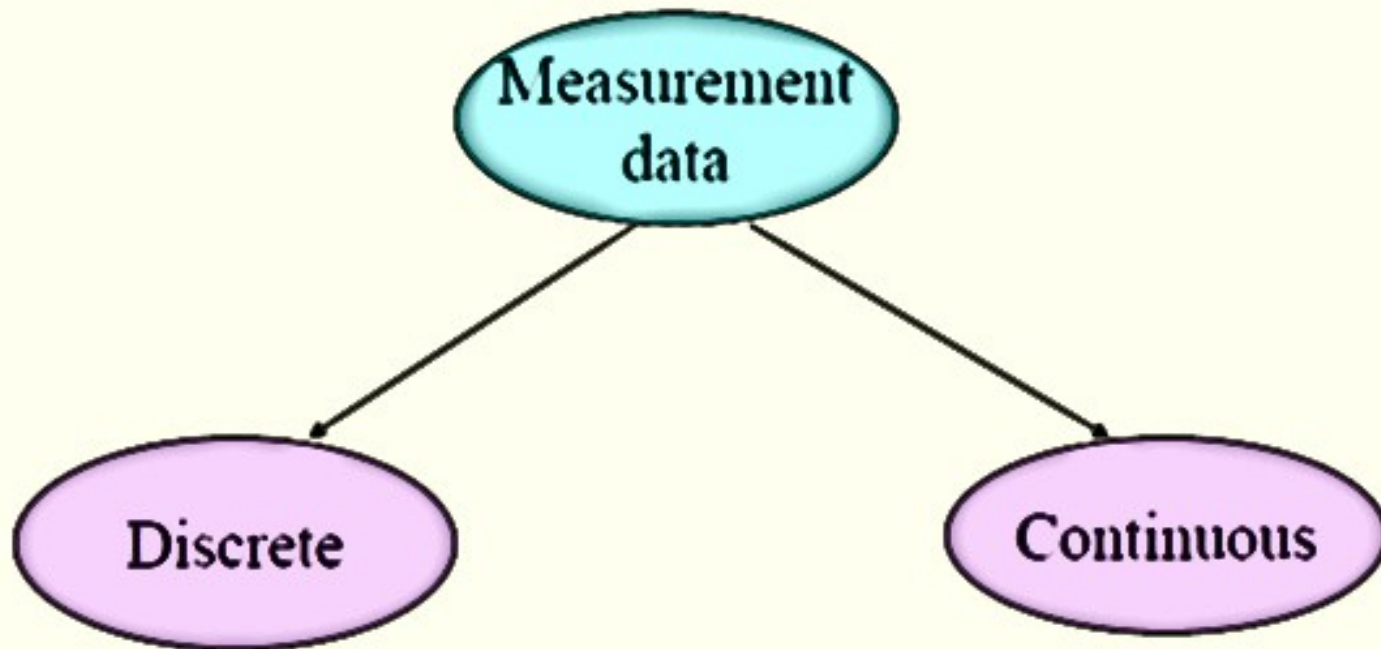
- Attendance

- present, absent

Measurement Data

- The objects being studied are **measured** based on some **quantitative** trait.
- The resulting data are set of numbers.
 - 0 Bio markers like cholesterol level, insulin level etc
 - 0 Height
 - 0 Age
 - 0 No. of students absent for the exam
 - 0 No. of parents not turned up for parents meeting

Classification of Measurement Data



Classification of Measurement Data

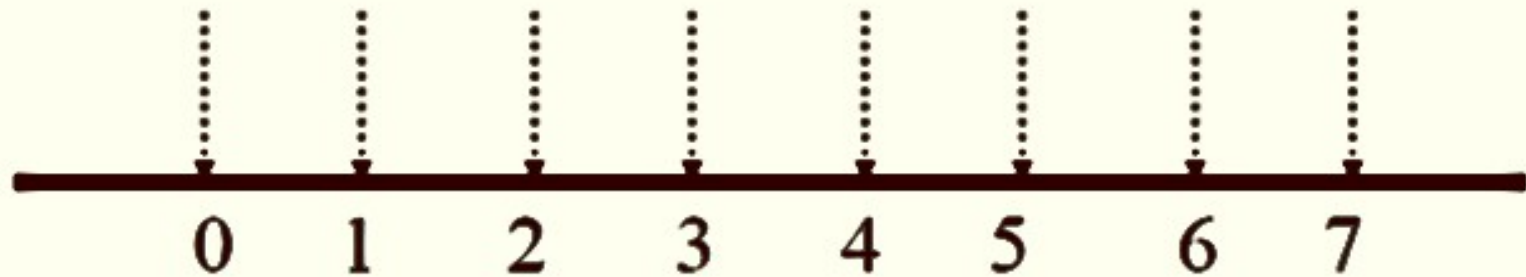
0 Discrete Measurement Data

0 Only certain values are possible (there are gaps between the possible values).

0 Continuous Measurement Data

0 Theoretically, any value within an interval is possible with a fine enough measuring device.

Discrete data -- Gaps between possible values



Continuous data -- *Theoretically,*
no gaps between possible values



Examples: Discrete Measurement Data

- 0 No. of students absent for the exam
- 0 No. of parents not turned up for the parents meeting
- 0 Generally, discrete data are counts.

Examples: Continuous Measurement Data

- 0 Bio markers like cholesterol level, insulin level etc
- 0 Height
- 0 Age
- 0 Generally, continuous data come from measurements.

Statistical Analysis

0 Categorical data are commonly summarized

using “**percentages**” (or “**proportions**”).

– 10% of students cleared exam was boys

Statistical Analysis

0 Measurement data are typically summarized

using **averages** (or **means**).

- Average weight of male students of II Year BE CSE is 70 KG.
- Average weight of female students of II Year BE CSE is 50 KG.

Missing Value Treatment

- 0 Missing data in the training data set can reduce the power/fit of a model
- 0 Lead to a biased model because the relationship between the other variables are not analyzed properly
- 0 May lead to wrong prediction or classification

Missing Value Treatment cont...

Name	Weight	Gender	Play Cricket/ Not
Mr. Amit	58	M	Y
Mr. Anil	61	M	Y
Miss Swati	58	F	N
Miss Richa	55		Y
Mr. Steve	55	M	N
Miss Reena	64	F	Y
Miss Rashmi	57		Y
Mr. Kunal	57	M	N

Gender	#Students	#Play Cricket	%Play Cricket
F	2	1	50%
M	4	2	50%
Missing	2	2	100%

Name	Weight	Gender	Play Cricket/ Not
Mr. Amit	58	M	Y
Mr. Anil	61	M	Y
Miss Swati	58	F	N
Miss Richa	55	F	Y
Mr. Steve	55	M	N
Miss Reena	64	F	Y
Miss Rashmi	57	F	Y
Mr. Kunal	57	M	N

Gender	#Students	#Play Cricket	%Play Cricket
F	4	3	75%
M	4	2	50%

chances of playing cricket by males is higher than females – in first table

chances of playing cricket by females is higher than males – in second table

Why data set has missing values?

- May occur at two stages:
 - Data Extraction
 - Data Collection
- **Data Extraction:**
 - Double-check for correct data with data guardians.
 - Some hashing procedures can be used to make sure data extraction is correct.
 - Typically easy to find and can be corrected easily as well.

Data collection

0 These errors occur at time of data collection and are harder to correct.

0 They can be categorized in four types:

0 Missing completely at random

0 Missing at random

0 Missing that depends on unobserved predictors

Data collection cont...

0 Missing completely at random:

0 This is a case when the probability of missing variable is same for all observations.

0 Example:

0 Respondents of data collection process decide that they will declare their earning after tossing a fair coin.

0 Missing at random:

0 This is a case when variable is missing at random and missing ratio varies for different values / level of other input variables.

0 Example:

0 Collecting data for age

0 Female has higher missing value compare to male.

Data collection cont...

0 Missing that depends on unobserved predictors:

0 This is a case when the missing values are not random and are related to the unobserved input variable.

0 Example:

0 In a medical study, if a particular diagnostic causes discomfort, then there is higher chance of drop out from the study.

0 This missing value is not at random unless we have included “discomfort” as an input variable for all patients.

0 Missing that depends on the missing value itself:

0 This is a case when the probability of missing value is directly correlated with missing value itself.

0 Example:

Which are the methods to treat missing values ?

0 Deletion

0 Imputation

Which are the methods to treat missing values ?

0 **Deletion:** Two types:

0 List Wise Deletion

0 Delete observations where any of the variable is missing.

0 This method reduces the power of model because it reduces the sample size.

0 Pair Wise Deletion.

0 Perform analysis with all cases in which the variables of interest are present.

0 Advantage of this method is, it keeps as many cases available for analysis.

0 One of the disadvantage of this method, it uses different sample size for different variables.

List Wise Deletion

- 0 listwise deletion will remove a case completely if it is missing a value for one of the variables included in the analysis.
- 0 you are conducting analyses using cumulative high school GPA, hours of study for first semester, SAT score, and first semester grade in college algebra.
- 0 Participant X is missing data for cumulative high school GPA, therefore, Participant X will be completely removed from the analyses because the participant does not have complete data for all the variables.

List-wise deletion

List wise deletion

Gender	Manpower	Sales
M	25	343
F	.	280
M	33	332
M	.	272
F	25	.
M	29	326
	26	259
M	32	297

Pair wise deletion

- 0 pairwise deletion will not omit a case completely from the analyses.
- 0 Pairwise deletion omits cases based on the variables included in the analysis
- 0 If you are conducting analyses using cumulative high school GPA, hours of study for first semester, SAT score, and first semester grade in college algebra
- 0 Participant X is missing data for cumulative high school GPA

Pair-wise deletion

- 0 Participant X will be omitted from any analyses using cumulative high school GPA
- 0 but they will not be omitted from analyses for which the participant has complete data.

pair-wise deletion

Pair wise deletion

Gender	Manpower	Sales
M	25	343
F	— .	280
M	33	332
M	— .	272
F	25	— .
M	29	326
— .	26	259
M	32	297

Imputation

- 0 Imputation is a method to fill in the missing values with estimated ones.
- 0 The objective is to employ known relationships that can be identified in the valid values of the data set to assist in estimating the missing values.

Mean/ Mode/ Median Imputation

- 0 Replacing the missing data for a given attribute by the mean or median (quantitative attribute) or mode (qualitative attribute) of all known values of that variable.
- 0 Two types:-
 - 0 **Generalized Imputation:**
 - 0 Calculate the mean or median for all non missing values of that variable then replace missing value with mean or median.
 - 0 Average of all non missing values of “**Manpower**” (28.33) and then replace missing value with it.
 - 0 **Similar case Imputation:**
 - 0 Calculate average for gender “**Male**” (29.75) and “**Female**” (25) individually of non missing values then replace the missing value based on gender.
 - 0 For “**Male**”, we will replace missing values of manpower with 29.75 and for “**Female**” with 25.

Maximum Likelihood Estimation

- 0 Maximum likelihood estimation is a method that will find the values of μ and σ that result in the curve that best fits the data.
- 0 The values that we find are called the maximum likelihood estimates (MLE).

- 0 *How do we calculate the maximum likelihood estimates of the parameter values using the log likelihood*
- 0 *Solution:*
- 0 *Consider the data points: 9, 9.5 and 11*

$$\ln(P(x; \mu, \sigma)) = -3 \ln(\sigma) - \frac{3}{2} \ln(2\pi) - \frac{1}{2\sigma^2} [(9 - \mu)^2 + (9.5 - \mu)^2 + (11 - \mu)^2]$$

- 0 This expression can be differentiated to find the maximum.
- 0 find the MLE of the mean, μ .
- 0 take the partial derivative of the function with respect to μ , giving

$$\frac{\partial \ln(P(x; \mu, \sigma))}{\partial \mu} = \frac{1}{\sigma^2} [9 + 9.5 + 11 - 3\mu] .$$

0 setting the left hand side of the equation to zero and then rearranging for μ gives:

$$\mu = \frac{9 + 9.5 + 11}{3} = 9.833$$

0 Similarly find σ

Multiple Imputation

0 Multiple Imputation follows:

0 Imputation Phase

0 Analysis Phase

0 Pooling Phase

Multiple Imputation

- 0 All multiple imputation methods follow three steps.
- 0 **Imputation** – missing values are imputed.
However, the imputed values are drawn m times from a distribution rather than just once. At the end of this step, there should be m completed datasets.
- 0 **Analysis** – Each of the m datasets is analyzed. At the end of this step there should be m analyses.
- 0 **Pooling** – The m results are consolidated into one result by calculating the mean, variance, and confidence interval of the variable of concern.

Multiple Imputation

- 0 In multiple imputation, the imputation process is repeated multiple times resulting in multiple imputed datasets
- 0 In this method the imputation uncertainty is accounted for by creating these multiple datasets
- 0 In the imputation model, the variables that are related to missingness, can be included.

Imputation Model

- 0 Imputation Phase
- 0 several copies of the data set are created each containing different imputed values
- 0 The imputed values are estimated using the means and covariance of the observed data
- 0 Regression equations are used to predict the incomplete values from the complete values and a normally distributed residual term is added to each value to restore variability

- 0 This process is iterated several times, updating the regression parameters after every iteration, to obtain different imputed values each time.
- 0 auxiliary variables can be included to improve the estimation of the imputed values.
- 0 imputed dataset is stored until the required number of imputed datasets is reached.
- 0 it is important to include the correct variables in the imputation process
- 0 imputation model should fits the distribution assumptions of the data.

- 0 when incomplete data are continuous and normally distributed, a multivariate normal distribution or linear regression can be used for the imputation
- 0 when data are not normal, or not continuous other imputation algorithms should be applied

Analysis Phase

- 0 the statistical analysis is carried out
- 0 On each imputed dataset, the analysis is carried out that would have been applied had the data been complete
- 0 That way as many sets of results are created as the number of imputed datasets created in the imputation phase

Pooling Phase

- 0 the multiple sets of results or parameter estimates are combined into a single set of results
- 0 When the estimates are pooled by Rubin's Rules, the parameter estimates are summarized by taking the average over the parameter estimates from all imputed datasets
- 0 The standard errors are pooled by combining the within imputation variance and the between imputation variance

Pooling Phase

$$Var_{within} = \frac{\sum_{i=1}^M SE^2_i}{M}$$

$$Var_{between} = \frac{\sum_{i=1}^M (\beta_i - \bar{\beta})^2}{M - 1}$$

$$Var_{total} = Var_{within} + Var_{between} + \frac{Var_{between}}{M}$$

Pooling formula's: Var is variance; SE is standard error; M is the number of imputed datasets; Beta is the parameter estimate.

Rubin's Rules

- 0 Rubin's Rules (RR) are designed to pool parameter estimates, such as mean differences, regression coefficients, standard errors and to derive confidence intervals and p-values.

Multiple Imputation

- 0 Multiple imputation can be used in cases where the data is
- 0 missing completely at random,
- 0 missing at random,
- 0 and even when the data is missing not at random.

Observed variables

- 0 Observed variables are actually measured by the researcher
- 0 Its classified into observed exogenous variable ie. It is not controlled by other variables similar to independent variable
- 0 Endogenous variable ie it is controlled by other variables similar to dependant variable
- 0 Example: job satisfaction scale, happiness measurement scale etc

Why Bayesian Estimation is used?

- 0 **Bayesian** inference is therefore just the process of deducing properties about a population or probability distribution from data using **Bayes'** theorem

Missing Completely at Random (MCAR)

- 0 Missing data are MCAR when the probability of missing data on a variable is unrelated to any other measured variable and is unrelated to the variable with missing values itself
- 0 Missingness on the variables are completely unsystematic

Missing Completely at Random (MCAR)

- 0 when data are missing for respondents for which their questionnaire was lost in the mail
- 0 separating the missing and the complete cases and examine the group characteristics. If characteristics are not equal for both groups, the MCAR assumption does not hold.

Examples of MCAR Data

Complete data	
Age	IQ score
25	133
26	121
29	91
30	105
30	110
31	98
44	118
46	93
48	141
51	104
51	116
54	97

Incomplete data	
Age	IQ score
25	
26	121
29	91
30	
30	110
31	
44	118
46	93
48	
51	
51	116
54	

Missing at Random (MAR)

- 0 When the probability of missing data on a variable is related to some other measured variable in the model, but not to the value of the variable with missing values itself.
- 0 For example, only younger people have missing values for IQ. In that case the probability of missing data on IQ is related to age.

0 It is recommended to incorporate correlates of missingness into the missing data handling procedure to diminish bias and improve the chances of satisfying the MAR assumption.

Example of Missing at Random Data

Complete data	
Age	IQ score
25	133
26	121
29	91
30	105
30	110
31	98
44	118
46	93
48	141
51	104
51	116
54	97

Incomplete data	
Age	IQ score
25	
26	
29	
30	
30	
31	
44	118
46	93
48	141
51	104
51	116
54	97

Missing not at Random

- 0 Data are missing not at random (MNAR) when the missing values on a variable are related to the values of that variable itself, even after controlling for other variables.
- 0 For example, when data are missing on IQ and only the people with low IQ values have missing observations for this variable.
- 0 A problem with the MNAR mechanism is that it is impossible to verify that scores are MNAR without knowing the missing values.

Missing not at Random (MNAR) Data

Complete data	
Age	IQ score
25	133
26	121
29	91
30	105
30	110
31	98
44	118
46	93
48	141
51	104
51	116
54	97

Incomplete data	
Age	IQ score
25	133
26	121
29	
30	
30	110
31	
44	118
46	
48	141
51	
51	116
54	

Introduction to Bayesian Estimation

- 0 It is based on probability of occurrence of an event.
- 0 Example:
- 0 Tossing a coin,
- 0 rolling a die, and
- 0 drawing a card out of a well-shuffled pack of cards

Bayes Theorem

- 0 Sample Space:
- 0 The result of an experiment is called an **outcome**.
- 0 **The set of all possible outcomes of an event is called the sample space.**
- 0 For example, if our experiment is throwing dice and recording its outcome,
- 0 the sample space will be: $S_1 = \{1, 2, 3, 4, 5, 6\}$
- 0 What will be the sample when we're tossing a coin?
- 0 $S_2 = \{H, T\}$

Bayes Theorem

0 Event:

0 **An event is a set of outcomes (i.e. a subset of the sample space) of an experiment.**

0 Example: Consider rolling of a dice

0 $E = \text{An even number is obtained} = \{2, 4, 6\}$

0 $F = \text{A number greater than 3 is obtained} = \{4, 5, 6\}$

0 The probability of these events:

0 $P(E) = \text{Number of favourable outcomes} / \text{Total number of possible outcomes} = 3 / 6 = 0.5$

0 $P(F) = 3 / 6 = 0.5$

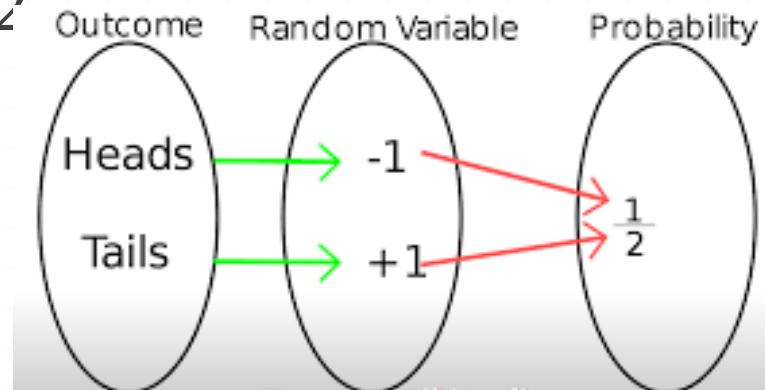
0 Then, $E \cup F = \{2, 4, 5, 6\}$ and $E \cap F = \{4, 6\}$

Bayes Theorem

- 0 Now consider an event G = An odd number is obtained: Then $E \cap G = \text{empty set} = \Phi$
- 0 Such events are called **disjoint** events
- 0 These are also called **mutually exclusive** events because only one out of the two events can occur at a time

Bayes Theorem

- 0 Random Variable:
- 0 Define a random variable X on the sample space of the experiment of tossing a coin.
- 0 It takes a value $+1$ if “Heads” is obtained and -1 if “Tails” is obtained.
- 0 Then, X takes on values $+1$ and -1 with equal probability of $1/2$



Bayes Theorem

0 Exhaustive Events

0 A set of events is said to be exhaustive if at least one of the events must occur at any time.

0 Thus, two events A and B are said to be exhaustive if $A \cup B = S$, the sample space.

0 Let's say that A is the event that a card drawn out of a pack is red and B is the event that the card drawn is black.

0 Here, A and B are exhaustive because the sample space $S = \{\text{red, black}\}$.

Bayes Theorem

0 Independent Events

0 If the occurrence of one event does not have any effect on the occurrence of another, then the two events are said to be independent

0 two events A and B are said to be independent if:

0 $P(A \cap B) = P(AB) = P(A) * P(B)$

0 if A is obtaining a 5 on throwing a die and B is drawing a king of hearts from a well-shuffled pack of cards, then A and B are independent just by their definition.

Bayes Theorem

- 0 **Conditional Probability**
- 0 **Conditional probability is defined as the probability of an event A, given that another event B has already occurred (i.e. A conditional B).**
- 0 This is represented by $P(A|B)$ and we can define it as:
- 0 $P(A|B) = P(A \cap B) / P(B)$
- 0 Let event A represent picking a king, and event B, picking a black card. Then, we find $P(A|B)$ using the above formula:

Bayes Theorem

0 $P(A|B) = P(A \cap B) / P(B)$

0 $P(A \cap B) = P(\text{Obtaining a black card which is a King}) = 2/52$
 $P(B) = P(\text{Picking a black card}) = 1/2$

0 Thus, $P(A|B) = 4/52.$

Bayes Theorem

- 0 **Marginal Probability**
- 0 **It is the probability of an event A occurring, independent of any other event B, i.e. marginalizing the event B.**
- 0 Marginal probability $P(A) = P(A|B)*P(B) + P(A|\sim B)*P(\sim B)$
- 0 $P(A) = P(A \cap B) + P(A \cap \sim B)$ //conditional probability
- 0 where $\sim B$ represents the event that B does not occur.

Bayes Theorem

- 0 the probability that a random card drawn out of a pack is red (event A) = $\frac{1}{2}$
- 0 calculate the same through marginal probability with event B as drawing a king.
- 0 $P(A \cap B) = \frac{2}{52}$ (because there are 2 kings in red suits)
- 0 $P(A \cap \sim B) = \frac{24}{52}$ (remaining cards from the red suit)
- 0 $P(A) = \frac{2}{52} + \frac{24}{52} = \frac{26}{52} = \frac{1}{2}$

Bayes Theorem

- 0 Consider that A and B are any two events from a sample space S where $P(B) \neq 0$.
- 0 Using conditional probability,
- 0 $P(A|B) = P(A \cap B) / P(B)$
- 0 $P(B|A) = P(A \cap B) / P(A)$
- 0 $P(A \cap B) = P(A|B) * P(B) = P(B|A) * P(A)$
- 0 $P(A|B) = P(B|A) * P(A) / P(B)$

Bayes Theorem

- 0 $P(A)$ and $P(B)$ are probabilities of observing A and B independently of each other (marginal probabilities)
- 0 $P(B|A)$ and $P(A|B)$ are conditional probabilities
- 0 $P(A)$ is called **Prior probability** and
- 0 $P(B)$ is called **Evidence**
- 0 $P(B) = P(B|A)*P(A) + P(B|\sim A)*P(\sim A)$
- 0 $P(B|A)$ is called **Likelihood** and $P(A|B)$ is called **Posterior probability**
- 0 Bayes Theorem can be written as:
- 0 $\text{posterior} = \text{likelihood} * \text{prior} / \text{evidence}$

Bayes Theorem

- 0 There are 3 boxes labeled A, B, and C:
- 0 Box A contains 2 red and 3 black balls
- 0 Box B contains 3 red and 1 black ball
- 0 And box C contains 1 red ball and 4 black balls
- 0 The three boxes are identical and have an equal probability of getting picked.
- 0 Consider that a red ball is chosen. Then what is the probability that this red ball was picked out of box A?

Bayes Theorem

- 0 Let E denote the event that a red ball is chosen and A, B, and C denote that the respective box is picked
- 0 We are required to calculate the conditional probability $P(A|E)$
- 0 We have prior probabilities $P(A) = P(B) = P(C) = 1/3$, since all boxes have equal probability of getting picked.
- 0 $P(E|A) = \text{Number of red balls in box A} / \text{Total number of balls in box A} = 2/5$
- 0 Similarly, $P(E|B) = 3/4$ and
- 0 $P(E|C) = 1/5$

Bayes Theorems

0 evidence $P(E) = P(E|A)*P(A) + P(E|B)*P(B) + P(E|C)*P(C)$

0 $= (2/5) * (1/3) + (3/4) * (1/3) + (1/5) * (1/3) = 0.45$

0 what is the probability that this red ball was picked out of box A?

0 $P(A|E) = P(E|A) * P(A) / P(E) = (2/5) * (1/3) / 0.45 = 0.296$

**T
H
A
N
K

Y
O
U**

CONTACT DETAILS

Dr. G. Malathi

Associate Professor

**Coordinator Image Processing Research
Group**

**School of Computing Sciences and
Engineering**

**VIT University
Chennai Campus**