

## Session -1

What is data science :-

A multi-disciplinary field that uses scientific methods, processes & algorithms and systems to extract knowledge & insights from structured & unstructured data.

Components of datascience

Domain Expertise & scientific methods

- Analysis
- Mathematical & statistical models
- Scientific tools & methods

Data Science

Technology

- Python language
- operating system
- Application design
- library
- data processing tools

Data science in sectors

- Data science is multidisciplinary
- In healthcare etc

## Data science applications

- prediction. → classification → Recommendations
- pattern detection & grouping
- Anomaly detection.
- Recognition. → Actionable processes
- Scoring & ranking → Segmentation
- optimization → Forecasts

## Bigdata overview

- Huge volume of data
- Complexity of datatypes & structures
- Speed of new data creation & growth

## Sources of big data

- Sensors → Audio
- social → server /
- Text Network log.
- location
- image
- video

## Industry:

Big data analytics in industry Verticals  
Banking & securities, Communication media & enter-  
tainment, Health care providers, Education,  
manufacturing & material resources, Government,  
Insurance, Retail & wholesale traders, Transporta-  
tion, Energy & utilities.

## Session - 2

### Data analytics life cycle

#### Life cycle overview

- The data analytics life cycle is designed to fit big data problems & data science projects
- with six phases. The project work can occur in several phases simultaneously.
- The cycle is iterative.
- Work can return to earlier phases as new information is uncovered.

## Description of key Roles in analytics project

- Business user - understands The domain area.
- Project sponsor - provides requirements.
- project manager - ensures meeting objectives
- Business Intelligence Analyst - provides business domain expertise based on deep understanding of the data.
- Database Administration (DBA) - creates DB environment.
- Data engineer - provides technical skills assists data management & extraction, supports analytic sandbox.
- Data scientist - provides analytic techniques & modeling.



## Description of key Roles in analytics project

- Business user - understands The domain area.
- Project sponsor - provides requirements.
- Project manager - ensures meeting objectives.
- Business Intelligence Analyst - provides business domain expertise based on deep understanding of the data.
- Database Administration (DBA) - creates DB environment.
- Data engineer - provides technical skills assists data management & extraction, supports analytic Sandbox.
- Data scientist - provides analytic techniques & modeling.



## The Analytics DELTA

Data breadth, integration, quality

Enterprise wide approach to managing analytics  
leadership, passion & commitment

Targets first deep, then broad

Analysts professionals & amateurs

## Data analytics life cycle phases

phase 1: Discovery

phase 2: Data preparation

phase 3: Model planning

phase 4: Model building

phase 5: Communicate results

phase 6: Operationalize

## phase 1: Discovery - Activities

1. Learning the Business Domain

2. Acquiring Resources

3. Framing the problem Definition

4. Identifying key stakeholders.
5. Interviewing the Analytics sponsor
6. Developing Initial Hypotheses.
7. Identifying potential Data sources.

## phase 2: Data preparation.

it has the following activities

1. preparing The Analytic Sand box.
2. performing ETLT.
3. Learning about the data.
4. Data conditioning.
5. Survey & visualize.
6. common Tools for data preparation.

→ 2.1 preparing The analytic Sand box.

→ 2.2 performing ETLT (Extract , Transform,  
Load , Transform)

2.3 Learning about the data

2.4 Data conditioning

2.5 Survey & visualize

2.6 Common tools for data preparation.

### Session-3

phase 3 model planning: Where the team

determines the methods, techniques & workflow it intends to follow for the subsequent, mod 4 building phase. The team explores the data to learn about the relationships between variables and subsequently selects key variables of the most suitable models.

phase 4: model building: The team develops

data sets for testing, training & production purposes.

In this phase the team builds & executes models based on the work done in the model planning phase.

## Phase 5: Communicate results:

The team, in collaboration with major stakeholders, determines if the results of the project are success (or) failure based on the criteria developed in phase 1.

## Phase 6: Operationalize:

The team delivers final reports, briefings, code & technical documents. In addition, the team may run a pilot project to implement the models in a production environment.

## Tools: Model building phase

### Commercial Tools:

- SAS Enterprise Miner
- SPSS Modeler → Alpine Miner
- Matlab → STATISTICA & Mathem atica

## free & open source tools:

- SQL → python, → WEKA → Octave
- R & PL/R

## benefits & use of data science

Example 1: Google AdSense, which collects data from internet users so relevant commercial messages can be matched to the person browsing the internet.

Example 2: thinkpoint Benefits & uses of data science & big data analytic is another example of real time personalized advertising.

## Facets of data

- Structured → Audio, video & images
- Unstructured → streaming.
- Natural language.
- Machine-generated.
- Graph-based.

## 1. Structured data

Structured data is data that depends on a data model and resides in a fixed field within a record. As such it's often easy to store structured data in tables within databases (e.g. Excel files shown in figure 1). Structured Query Language is the preferred way to manage and query data that resides in database.

## 2. Unstructured data

- Unstructured data is data that isn't easy to fit into a data model because the content is context specific (or varying).
- The thousands of different languages & dialects out there further complicate this. A human written email, as shown in figure 2, is also a perfect example of natural language data.

### 3. Natural language

→ Natural language is a special type of unstructured data; it's challenging to process because it requires knowledge of specific data science techniques & linguistics.

### 4. Machine generated data

→ Machine-generated data is information that's automatically created by a computer, process, application (or) other machine without human intervention. Machine-generated data is becoming a major data resource and will continue to do so.

### 5. Graph based (or) network data.

→ Graph data can be a confusing term because any data can be shown in graph in this case points to mathematical graph theory.

### 6. Audio, image & video

→ Audio, image & video all data types

that pose specific challenges to a data scientist tasks that are trivial for humans, such as recognizing objects in pictures, turn out to be challenging for computers.

## 7. Streaming data

→ while streaming data can take almost any of the previous forms, it has an extra property. The data flows into the system when an event happens instead of being loaded into a datastore in a batch.

## Session-4

### Data science process & Data wrangling



# Data science process

- 1. Setting The research goal.
- 2. Retrieving data.
- 3. Data preparation.
- 4. Data exploration.
- 5. Data modeling.
- 6. presentation & automation.

## 1. setting the research goal.

- Data science is mostly applied in the context of an organization, when the business asks to perform a data science project. This charter contains info such as what is the work going to research.

## 2. Retrieving data

## 2. Retrieving data

- The second step is to collect data. stated in the project charter which data is need & where can find it.

### 3. Data preparation

→ Data collection is an error prone process in this phase enhance the quality of the data and prepare it for use in subsequent steps.

### 4. Data exploration

Data exploration is concerned with building a deeper understanding of data. Try to understand how variable interact with each other, The distribution of data, & whether there are outliers.

### 5. Data modeling (o) model building.

→ In This phase use models, domain knowledge and insights about the data from in previous steps to answer the research question.

## 6. presentation & automation

→ Finally, present the results to business.  
These results can take many forms, ranging from presentations to research reports.

These six are again subdivided:

1. setting the research goal

↓  
define research goal      Create project charter  
                                ↓  
                                Internal data      External data  
                                ↓  
                                Data retrieval      Data ownership

2. retrieving data

↓  
Data entry errors  
Physical impossible values  
Missing values  
Outliers  
Spaces, typos

3. Data preparation

↓  
Combining data      Data transformation  
                                ↓  
                                Aggregating data  
                                Extrapolating data  
                                Derived measures  
                                Creating dummies  
                                Reducing number of variables  
                                ↓  
                                Joining/merging datasets  
                                Set operations  
                                Creating views

4. Data exploration

- simple graphs
- combined graphs
- link & brush
- Non graphical techniques

5. Data modeling

- model and variable selection
- model execution.
- model diagnostic & model comparison.

6. presentation & automation

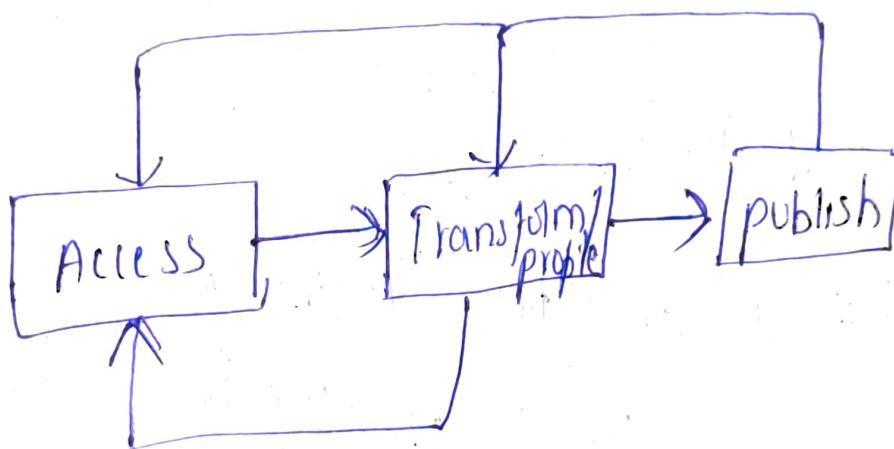
- presenting data
- Automating data analysis.

## Data wrangling

→ Data wrangling is a generic phrase capturing the range of the tasks involved in preparing the data for analysis; it begins with accessing data, sometimes, access

is gated on getting appropriate permission & making the corresponding changes in data infrastructure. Access also involves manipulating the locations & relationships between dataset.

## A simple diagram



## Session 5

### Data analytic process

#### Why data analytics?

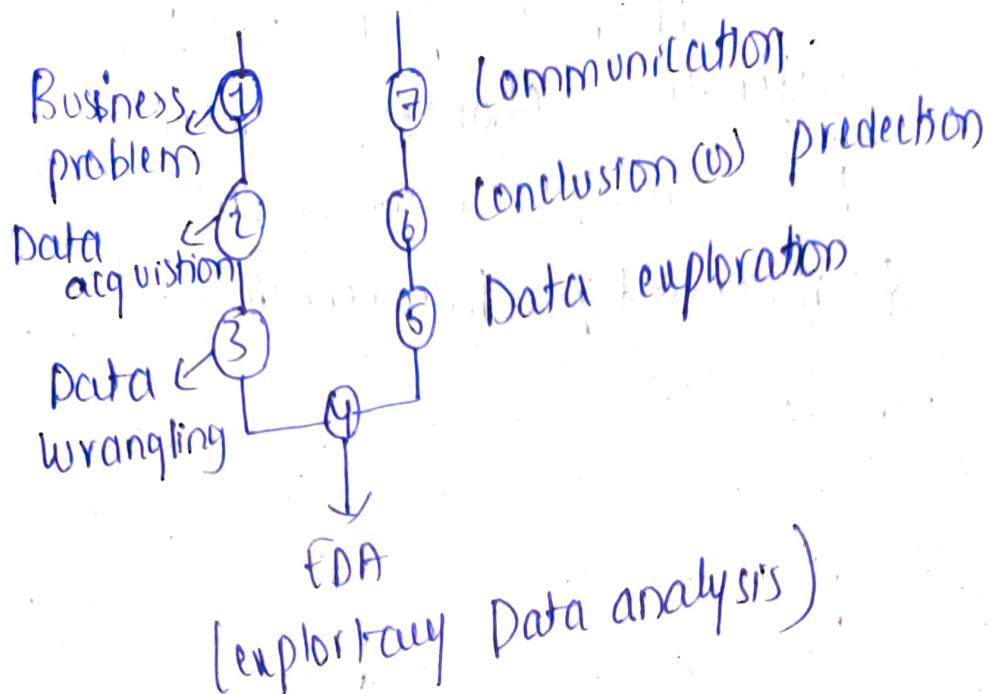
Data by itself is just an information source

But unless you understand it, you will be not be able to use it effectively.

Introduction

#### \* Data analytics:

It is a combination of process to extract data from datasets.



## 1. Business problem

The process of analytics begins with questions of stakeholders.

- (1) business problems of stakeholders.

## 2. Data acquisition

Data scientists Expertise

- File handling
- File formats
- Web scrapping

## ③ Data wrangling

Data cleansing

Data manipulation

↓

Data wrangling.

## 4. exploratory Data analysis (EDA)

- EDA approach studies the data to recommend suitable models that best fit the data.
  - The focus is on data: its structure, outliers, and models suggested by the data.
- Significance of EDA, Steps in EDA, EDA making sense of data, EDA: measurement scales, EDA quantitative technique, EDA graphical techniques all in PR (course material 2/2).

## 5. Data Exploration: model selection

### model selection

- Based on the overall data analysis process.

- Should be accurate to avoid iterations
- Depends on pattern identification & algorithms
- Depends on pattern identification & algorithms
- Depends on hypothesis building & testing

## 7. Communication

Forms of data analysis presentations:

- visual graphs
- plotting maps.
- Reports
- white paper reports
- powerpoint presentations

### Session-6

#### EDA: Graphical Techniques

##### Shape based EDA : Skewness

- The literal meaning of "skew" is a bias, dragging (or) distortion toward some particular value, group, subjects (or) direction.
- A measure of skewness is a numeric metric to concisely summarize the degree of asymmetry of a unimodal distribution that can be compared with other similar numbers.

## Calculating Skewness

$$g_1 = \frac{\sum_{i=1}^N (y_i - \bar{y})^3}{s^3}$$

where  $\bar{y}$  is the mean

$s$  is the standard deviation

$N$  is the number of datapoints.

The above formula for skewness is referred to as the Fisher-Pearson coefficient of skewness.

Skewed graphs in pg. 24.

## Shape based EDA & kurtosis

→ Kurtosis is a statistical measure that defines how heavily the tails of distribution differ from the tails of normal distribution.

→ In other words, kurtosis identifies whether the tails of a given distribution contain extreme values.

### Types of kurtosis

→ mesokurtic → platykurtic

→ leptokurtic

## Mesokurtic

- Data that follows a mesokurtic distribution shows an excess kurtosis of zero or close to zero.

## Leptokurtic

- leptokurtic indicates a positive excess kurtosis.

## Platykurtic

- A platykurtic distribution shows a negative excess kurtosis.

## Histograms

- Histogram graphically summarizes the distribution of univariate dataset. It shows:
  - The centre (or) location of data.
  - The Skewness of data.
  - The presence of outliers.
  - The presence of multiple modes in the data.

## Scatter plots

→ A scatterplot represents relationships between two variables.

CO<sub>2</sub>

## Data visualization

## Session 7

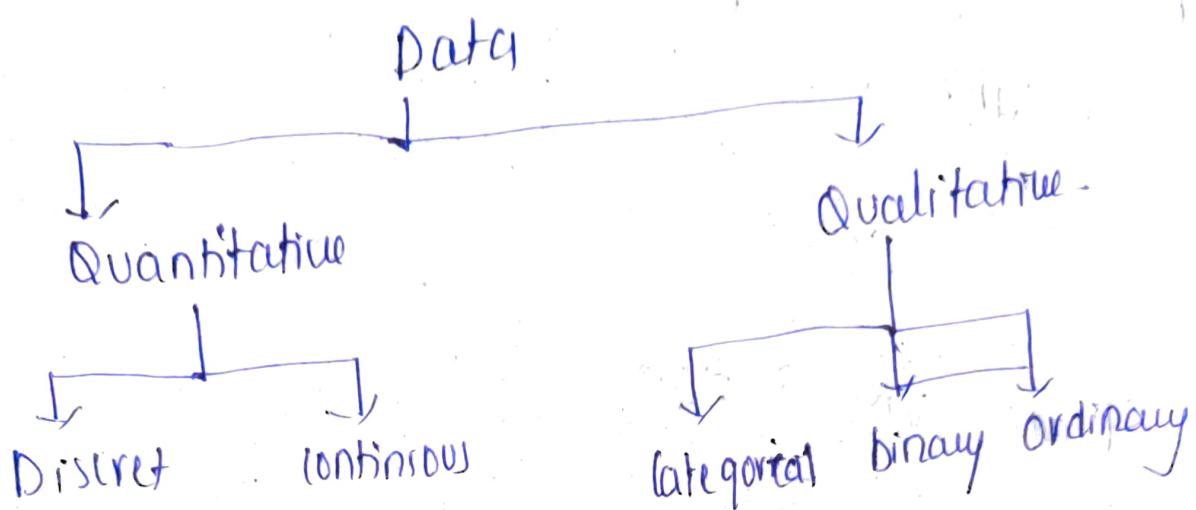
### Data analysis

It is the most crucial steps that deal with descriptive statistics and analysis of the data.

### Data analysis techniques

- Data summarizations are summary tables
- Graphs.
- Descriptive statistics, inferential statistics, correlation statistics
- Searching, grouping & mathematical models.

# Data types for plotting



## Quantitative data

The information is recorded as numbers.  
it represents an objective measurement (or) a count

### continuous data

continuous variables can take any numeric value  
if can meaningfully divided into smaller increments  
including fractional & decimal values.

### Discrete Data

discrete data has a limited number of  
possible values.