



BITS Pilani

Pilani | Dubai | Goa | Hyderabad

INTRODUCTION TO DATA SCIENCE

SESSION # 1 : INTRODUCTION

SANKARA NAYAKI K

sankaranayaki@wilp.bits-pilani.ac.in

The instructor is gratefully acknowledging the authors who made their course materials freely available online.

References:

- Introducing Data Science by Cielen, Meysman and Ali
- Storytelling with Data by Cole Nussbaumer Knaflic; Wiley
- Introduction to Data Mining by Tan, Steinbach and Vipin Kumar
- The Art of Data Science by Roger D Peng and Elizabeth Matsui
- Python Data Science Handbook: Essential tools for working with data by Jake VanderPlas

TABLE OF CONTENTS

1 COURSE HANDOUT

2 EVALUATION COMPONENTS

3 DATA SCIENCE



COURSE HANDOUT

- M1 Introduction to Data Science
- M2 Data Analytics
- M3 Data Science Process
- M4 Data Science Teams
- M5 Data and Data Models
- M6 Data wrangling and Feature Engineering
- M7 Data visualization
- M8 Storytelling with Data
- M9 Ethics for Data Science

LAB SESSIONS

The Lab capsules has to be practised. Details will be posted soon in Impartus.

L1 17-Oct-20

L2 16-Jan-21



PLATFORMS / TOOLS

1 Python

DATASET



1 Iris dataset

Any other dataset can be used.

TABLE OF CONTENTS

1 COURSE HANDOUT

2 EVALUATION COMPONENTS

3 DATA SCIENCE

EVALUATION COMPONENTS

Component	Weightage	Deadline
Quiz-I (Pre-Mid)	5 %	TBA
Quiz-II (Post-Mid)	5 %	TBA
Assignment	20 %	TBA
Mid-Semester Exam	30 %	Per Schedule
Comprehensive Exam	40 %	Per Schedule

Further announcements will be posted in Canvas.

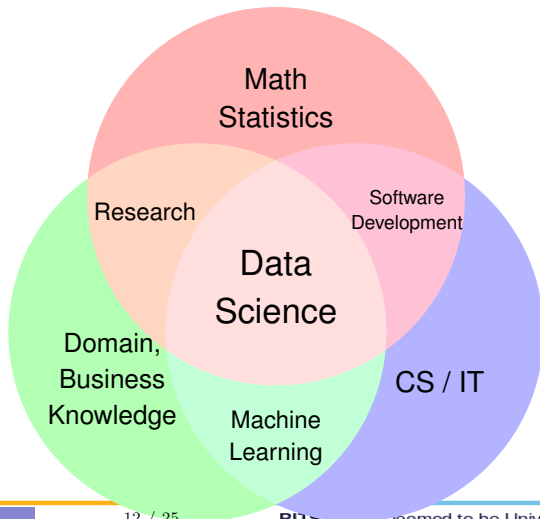
TABLE OF CONTENTS

- 1 COURSE HANDOUT
- 2 EVALUATION COMPONENTS
- 3 DATA SCIENCE



- Data Science is a study of data.
- Data Science is an art of uncovering insights and trends that are hiding behind the data.
- Data Science is the process of using data to understand different things.
- Data Science helps to translate data into a story. The story telling helps in uncovering insights. The insights help in making decision or strategic choices.

DATA SCIENCE – MULTIPLE DISCIPLINES



NEED OF DATA SCIENCE

- Data deluge, tons of data.
- Powerful algorithms.
- Open software and tools.
- Computational speed, accuracy and cost.
- Data storage in terms of capacity and cost.

USE CASES OF DATA SCIENCE



APPLICATIONS OF DATA SCIENCE



APPLICATIONS OF DATA SCIENCE

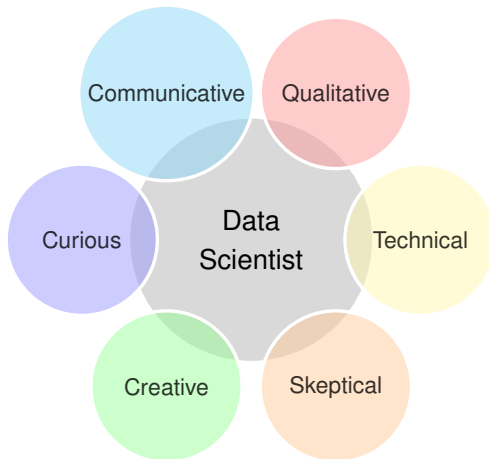




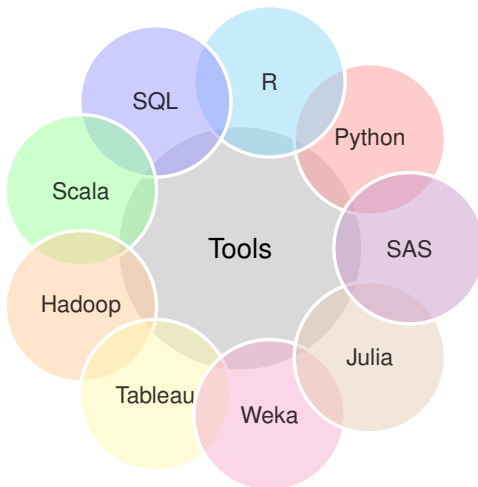
ROLE OF A DATA SCIENTIST

- ❶ Reframe business challenges as analytics challenges.
This is a skill to diagnose the problem, consider the core of a given problem, and determine which kinds of candidate analysis analytical method can be applied to solve it.
- ❷ Design, implement and deploy statistical models and data mining technique on data.
This activity is mainly the role of data scientist, applying complex or advanced analytical methods to a variety of business problem using data.
- ❸ Develop insights that lead to actionable recommendations.
Learn how to draw insights out of data and communicate them effectively.

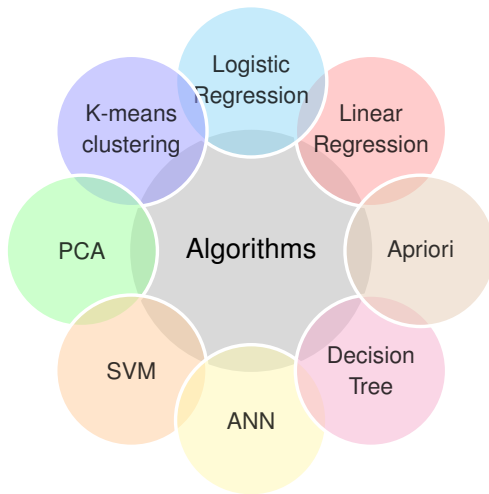
SKILLS REQUIRED FOR A DATA SCIENTIST



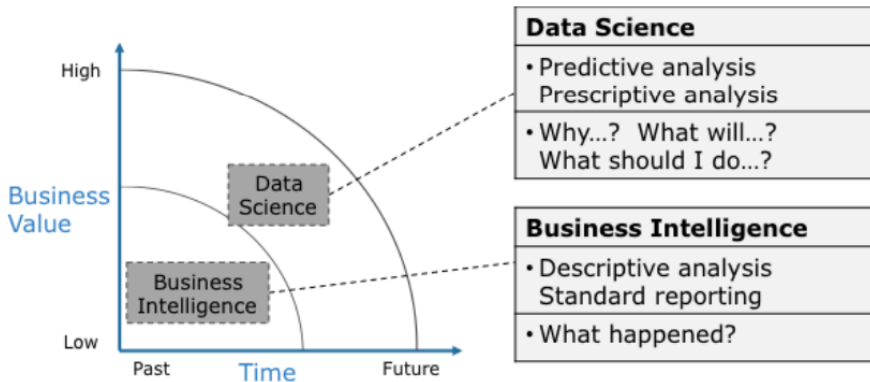
TOOLS AVAILABLE TO A DATA SCIENTIST



ALGORITHMS FOR A DATA SCIENTIST



DATA SCIENCE VS. BUSINESS INTELLIGENCE



DATA SCIENCE VS. BUSINESS INTELLIGENCE

	Data Science	Business Intelligence
Perspective	Looking forward	Looking backward
Analysis	Predictive Explorative	Descriptive Comparative
Data	Same data, New analysis Listens to data Distributed	New Data, Same analysis Speaks for data Warehoused
Scope	Specific to business question	Unlimited
Expertise	Data scientist	Business analyst
Deliverable	Insight or story	Table

DATA SCIENTIST VS. BUSINESS ANALYST

Area	BI Analyst	Data Scientist
Focus	Reports, KPIs, trends	Patterns, correlations, models
Process	Static, comparative	Exploratory, experimentation, visual
Data sources	Pre-planned, added slowly	On the fly, as-needed
Transform	Up front, carefully planned	In-database, on-demand, enrichment
Data quality	Single version of truth	"Good enough", probabilities
Data model	Schema on load	Schema on query
Analysis	Retrospective, Descriptive	Predictive, Prescriptive

DATA SCIENCE VS. STATISTICS

	Data Science	Statistics
Type of problem	Semi structured or un-structured	Well structured
Analysis Objective	Need not be well formed	Well formed objective
Type of Analysis	Explorative	Confirmative
Data collection	Data collection is not linked to the objective	Data collected based on the objective
Size of dataset	Large Heterogeneous	Small Homogeneous
Paradigm	Theory and heuristic (deductive & inductive)	Theory based (deductive)



SOFTWARE ENGINEERING FOR DATA SCIENCE

- ❶ For data scientists, software is the generalization of a specific aspect of a data analysis. Software allows for the systematizing and the standardizing of a procedure, so that different people can use it and understand what it's doing, at any given time. Software encompasses all required tools into a specific module or procedure that can be repeatedly applied in a variety of settings
- ❷ Software will have an interface, or a set of inputs and a set of outputs that are well understood.
3 Levels of S/W : Code, Write Function S/W Package or API
- ❸ Application of Skill Sets differ in
Methodologies, Objectives, Approaches Tools.

DATA SCIENCE CHALLENGES

- 1 Complexity of Data Reality
- 2 Identifying the problem
- 3 Access to right data – Data quantity
- 4 Data Cleansing – Data quality - Data Security
- 5 Granularity, Consistency Availability of Data
- 6 Lack of domain expertise
- 7 Cognitive Bias Content Bias

THANK YOU