# BIRLA INSTITUTE OF TECHNOLOGY & SCIENCE, PILANI

## WORK INTEGRATED LEARNING PROGRAMMES

### Part A: Content Design

| | |
|---|---|
| **Course Title** | Introduction to Data Science |
| **Course No** | DSE ZG523 |
| **Credit Units** | 3 |
| **Last Revised by** | Pravin Y Pawar |
| **Version / Date** | 1.2 , 15/04/2020 |

### Course Objectives

| # | Course Objectives |
|---|---|
| 1 | Gain basic understanding of the role of Data Science in various scenarios in the real-world of business, industry and government |
| 2 | Understand various roles and stages in a Data Science Project and ethical issues to be considered. |
| 3 | Explore the processes, tools and technologies for collection and analysis of structured and unstructured data |
| 4 | Appreciate the importance of techniques like data visualization, storytelling with data for the effective presentations of the outcomes with the stakeholders |

### Text Books/References

| ID | Text Book |
|---|---|
| T1 | Introducing Data Science by Cielen, Meysman and Ali |
| T2 | Storytelling with Data, A data visualization guide for business professionals, by Cole Nussbaumer Knaflic; Wiley |
| T3 | Introduction to Data Mining, by Tan, Steinbach and Vipin Kumar |
| | **References** |
| R1 | The Art of Data Science by Roger D Peng and Elizabeth Matsui |
| R2 | Ethics and Data Science by DJ Patil, Hilary Mason, Mike Loukides |
| R3 | Python Data Science Handbook: Essential tools for working with data by Jake VanderPlas |
| R4 | KDD, SEMMA and CRISP-DM: A Parallel Overview , Ana Azevedo and M.F. Santos , IADS-DM, 2008 |

\* The above materials are reference only and are neither conclusive nor exhaustive. However, the student is advised to refer latest content from online sources or instructor supplied materials for more thorough understanding of the topics

**Contents & Session delivery**

| Session (2 hrs) | Topics to cover | Content Reference |
|---|---|---|
| 1. | **Introduction to Data Science**<br>• Fundamentals of Data Science<br>• Real World applications<br>• Data Science vs BI<br>• Data Science vs Statistics<br>• Roles and responsibilities of a Data Scientist<br>• Software Engineering for Data Science<br>• Data Scientists Toolbox<br>• Data Science Challenges | T1 – Chapter 1<br><br>Class Room Discussion<br>Class Notes<br><br>Additional Reading (AR) material provided on Canvas LMS<br><br>Class Room Discussion |
| 2.<br><br><br><br><br>3. | **Data Analytics**<br>• Defining Analytics<br>• Types of data analytics<br>   o Descriptive, Diagnostic<br>   o Predictive, Prescriptive<br>• Data Analytics – methodologies<br>   o CRISP-DM Methodology<br>   o SEMMA<br>   o BIG DATA LIFE CYCLE<br>   o SMAM<br>• Analytics Capacity Building<br>• Challenges in Data-driven decision making | Class Notes<br>Beginners Guide to Analytics<br>Data Analytics Tutorial<br><br>R4 and Class Notes<br><br><br><br>Class Room Discussions |
| 4.<br><br><br><br><br><br><br><br>5. | **Data Science Process**<br>• Data Science methodology<br>   o Business understanding<br>   o Data Requirements<br>   o Data Acquisition<br>   o Data Understanding<br>   o Data preparation<br>   o Modelling<br>   o Model Evaluation<br>   o Deployment and feedback<br>• Case Study<br>• Data Science Proposal<br>   o Samples<br>   o Evaluation<br>   o Review Guide | T1 - Chapter 2<br>R1 – Chapter 2<br><br><br><br><br><br><br><br>Class Room Discussion<br>Class Notes and AR |
| 6. | **Data Science Teams**<br>• Defining Data Team<br>• Roles in a Data Science Team<br>   o Data Scientists<br>   o Data Engineers<br>• Managing Data Team<br>   o On boarding and evaluating the success of team<br>   o Working with other teams<br>   o Common difficulties | Class Room Discussion<br><br><br><br>Class Notes and AR |

| | | |
|---|---|---|
| 7.<br><br><br><br><br><br><br>8. | **Data and Data Models**<br>• Types of Data and Datasets<br>• Data Quality<br>• Epicycles of Data Analysis<br>• Data Models<br>   o Model as expectation<br>   o Comparing models to reality<br>   o Reactions to Data<br>   o Refining our expectations<br>• Six Types of the Questions<br>• Characteristics of Good Question<br>• Formal modelling<br>   o General Framework<br>   o Associational Analyses<br>   o Prediction Analyses | T3 – Chapter 2<br><br><br>R1 – Chapter 2<br>R1 – Chapter 5<br><br><br><br>R1 – Chapter 3<br><br><br>R1 – Chapter 7 |
| 9.<br><br><br><br><br><br><br><br>10.<br><br><br><br><br><br><br><br><br>11. | **Data wrangling and Feature Engineering**<br>• Data cleaning<br>• Data Aggregation, Sampling,<br>• Handling Numeric Data<br>   o Discretization, Binarization<br>   o Normalization<br>   o Data Smoothening<br>• Dealing with textual Data<br>• Managing Categorical Attributes<br>   o Transforming Categorical to Numerical Values<br>   o Encoding techniques<br>• Feature Engineering<br>   o Feature Extraction (Dimensionality Reduction)<br>   o Feature Construction<br>   o Feature Subset selection<br>      ▪ Filter methods<br>      ▪ Wrapper methods<br>      ▪ Embedded methods<br>   o Feature Learning<br><br>• Case Study involving FE tasks | T3 – Chapter 2<br><br><br>T3 – Chapter 2<br><br><br>Class Notes and AR<br><br><br>R3 – Chapter 5.4<br><br><br><br>Class Notes and AR<br>T3 – Appendix B<br><br><br>AR<br>Class Room Discussions<br><br><br><br>Class Room Discussions |
| 12.<br><br><br><br>13. | **Data Visualizations**<br>• Data Need for visualization<br>• Exploratory vs Explanatory Analysis<br><br>• Tables , Axis based Visualization and Statistical Plots<br>• Lessons in Data Visualization Design<br>• The Data Visualization Design Process<br>• Stories and Dashboards | Why Visual Analytics?<br>Visual Analysis for Everyone<br>T2 – Chapter 1<br><br>T2 – Chapter 2<br><br>T2 – Chapter 8<br>Class Notes<br>Class Room Discussion |
| 14. | **Storytelling with Data**<br>• The final deliverable<br>• The Narrative - report / presentation structure<br>• Building narrative with Data | AR<br>Class Room Discussion<br><br>T2 – Chapter 10 |

| | | | |
|---|---|---|---|
| | • Effective storytelling | R1 – Chapter 10 | |
| 15. | **Ethics for Data Science**<br>• Bias and Fairness<br> ○ Types of Bias<br> ○ Identifying Bias<br> ○ Evaluating Bias<br>• Being a data skeptic – examples of misuse of Data | AR<br>Hidden Biases in Big Data<br><br><br>On Being a Data Skeptic | |
| 16. | • Doing Good Science<br>• Five C's<br>• Ethical guidelines for Data Scientist<br>• Ethics of data scraping and storage<br><br>• Case Study: IBM AI Fairness 360 | R2 – Chapter 1<br>R2 – Chapter 3<br><br>Class Room Discussion<br><br>Credit Decisions<br>Medical Expenditures | |

**Evaluation Scheme**:
Legend: EC = Evaluation Component; AN = After Noon Session; FN = Fore Noon Session

| | Name | Type | Duration | Wt. | Date/Deadline |
|---|---|---|---|---|---|
| EC1 | Quiz-I (Pre-Mid)<br>Quiz-II (Post-Mid)<br>Assignment | Online (30 mins)<br>Online (30 mins)<br>Take-home | 5 days open<br>5 days open<br>10 days open | 5%<br>5%<br>20% | TBA |
| EC2 | Mid-Semester Exam | Closed Book | 1.5 hours | 30% | Per Schedule |
| EC3 | Comprehensive Exam | Open Book | 2.5 hours | 40% | Per Schedule |

**Notes:**
➔ The release dates of Quiz-1/2 and assignments will be 5 days (for Quiz) and 10 days (for assignments) before the completion/submission deadline.
➔ **Deadlines will NOT be extended for whatever reason** and the student is requested not to wait for the deadline to start working on Quiz/Assignment
➔ Syllabus for Quiz-I: Sessions: 1 to 3 / Quiz-II : Session 9 to 12
➔ Syllabus for Assignment: Hands-on Python-based Exercise (real-world problem, for individual group of 3 / 4 students). Group formation procedure will be announced before Assignment release
➔ All Quiz/Assignments will be released and to be answered/submitted in Canvas LMS
➔ Syllabus for Mid-Semester Test (Closed Book): Topics in Session Nos. 1 to 8
➔ Syllabus for Comprehensive Exam (Open Book): All topics (Session Nos. 1 to 16)
➔ The student is strictly advised to stick to regular schedule of Mid-Sem and Compre examinations, and Makeup examinations will be only for those students with business-related absence/health related issues.
➔ Strictly NO MAKEUPS for Quiz and Assignments and all submissions after the above stated deadlines will not be considered/evaluated.
➔ All students should conform to BITS students' ethical code-of-conduct and all assignments will be subjected to plagiarism check, and if violated will be subject to disciplinary action apart from nullifying all the marks/grades assigned.

**Important links and information:**
<u>Canvas LMS:</u> All materials/announcements/discussions forums/Online Quizs/Assignment submissions will be via Canvas LMS portal. Students are expected to monitor this portal regularly for any content or announcements.

Contact sessions: Students should attend the online lectures as per the schedule provided in the Course Handout (posted on Canvas LMS)

Evaluation Guidelines:
1. EC-1 consists of 2 Quizzes and 1 Assignments. Students will attempt them through the course pages on the Canvas portal. Announcements will be made on the portal, in a timely manner.
2. For Closed Book tests: No books or reference material of any kind will be permitted.
3. For Open Book exams: Use of books and any printed / written reference material (filed or bound) is permitted. However, loose sheets of paper will not be allowed. Use of calculators is permitted in all exams. Laptops/Mobiles of any kind are not allowed. Exchange of any material is not allowed.
4. If a student is unable to appear for the Regular Test/Exam due to genuine exigencies, the student should follow the procedure to apply for the Make-Up Test/Exam which will be made available on the Elearn portal. The Make-Up Test/Exam will be conducted only at selected exam centers.

It shall be the responsibility of the individual student to be regular in attending the contact-session schedule as given in the course handout, and take all the prescribed evaluation components such as Assignment/Quiz, Mid-Semester Test and Comprehensive Exam according to the evaluation scheme provided in the handout