

INTRODUCTION TO BIG DATA ANALYTICS

(OPEN ELECTIVE)

(Effective from the academic year 2018 -2019)

SEMESTER – VII

Course Code	18CS751	CIE Marks	40
Number of Contact Hours/Week	3:0:0	SEE Marks	60
Total Number of Contact Hours	40	Exam Hours	03

CREDITS –3

Course Learning Objectives: This course (18CS751) will enable students to:

- Interpret the data in the context of the business.
- Identify an appropriate method to analyze the data
- Show analytical model of a system

Module – 1

Teaching
Hours

08

Introduction to Data Analytics and Decision Making: Introduction, Overview of the Book, The Methods, The Software, Modeling and Models, Graphical Models, Algebraic Models, Spreadsheet Models, Seven-Step Modeling Process. **Describing the Distribution of a Single Variable:** Introduction, Basic Concepts, Populations and Samples, Data Sets, Variables, and Observations, Types of Data, Descriptive Measures for Categorical Variables, Descriptive Measures for Numerical Variables, Numerical Summary Measures, Numerical Summary Measures with StatTools, Charts for Numerical Variables, Time Series Data, Outliers and Missing Values, Outliers, Missing Values, Excel Tables for Filtering, Sorting, and Summarizing.

Finding Relationships among Variables: Introduction, Relationships among Categorical Variables, Relationships among Categorical Variables and a Numerical Variable, Stacked and Unstacked Formats, Relationships among Numerical Variables, Scatterplots, Correlation and Covariance, Pivot Tables.

Textbook 1: Ch. 1,2,3

RBT: L1, L2, L3

Module – 2

08

Probability and Probability Distributions: Introduction, Probability Essentials, Rule of Complements, Addition Rule, Conditional Probability and the Multiplication Rule, Probabilistic Independence, Equally Likely Events, Courseive Versus Objective Probabilities, Probability Distribution of a Single Random Variable, Summary Measures of a Probability Distribution, Conditional Mean and Variance, Introduction to Simulation.

Normal, Binormal, Poisson, and Exponential Distributions: Introduction, The Normal Distribution, Continuous Distributions and Density Functions, The Normal Density, Standardizing: Z-Values, Normal Tables and Z-Values, Normal Calculations in Excel, Empirical Rules Revisited, Weighted Sums of Normal Random Variables, Applications of the Normal Random Distribution, The Binomial Distribution, Mean and Standard Deviation of the Binomial Distribution, The Binomial Distribution in the Context of Sampling, The Normal Approximation to the Binomial, Applications of the Binomial Distribution, The Poisson and Exponential Distributions, The Poisson Distribution, The Exponential Distribution.

Textbook 1: Ch. 4,5

RBT: L1, L2, L3

Module – 3

Decision Making under Uncertainty: Introduction, Elements of Decision Analysis, Payoff

08

Tables, Possible Decision Criteria, Expected Monetary Value(EMY),Sensitivity Analysis, Decision Trees, Risk Profiles, The Precision Tree Add-In,Bayes' Rule, Multistage Decision Problems and the Value of Information, The Value of Information, Risk Aversion and Expected Utility, Utility Functions, Exponential Utility, Certainty Equivalents, Is Expected Utility Maximization Used?

Sampling and Sampling Distributions: Introduction, Sampling Terminology, Methods for Selecting Random Samples, Simple Random Sampling, Systematic Sampling, Stratified Sampling, Cluster Sampling, Multistage Sampling Schemes, Introduction to Estimation, Sources of Estimation Error, Key Terms in Sampling, Sampling Distribution of the Sample Mean, The Central Limit Theorem, Sample Size Selection, Summary of Key Ideas for Simple Random Sampling.

Textbook 1: Ch. 6,7

RBT: L1, L2, L3

Module – 4

Confidence Interval Estimation: Introduction, Sampling Distributions, The t Distribution, Other Sampling Distributions, Confidence Interval for a Mean, Confidence Interval for a Total, Confidence Interval for a Proportion, Confidence Interval for a Standard Deviation, Confidence Interval for the Difference between Means, Independent Samples, Paired Samples, Confidence Interval for the Difference between Proportions, Sample Size Selection, Sample Size Selection for Estimation of the Mean, Sample Size Selection for Estimation of Other Parameters.

Hypothesis Testing:Introduction,Concepts in Hypothesis Testing, Null and Alternative Hypothesis, One-Tailed Versus Two-Tailed Tests, Types of Errors, Significance Level and Rejection Region, Significance from p-values, Type II Errors and Power, Hypothesis Tests and Confidence Intervals, Practical versus Statistical Significance, Hypothesis Tests for a Population Mean, Hypothesis Tests for Other Parameters, Hypothesis Tests for a Population Proportion, Hypothesis Tests for Differences between Population Means, Hypothesis Test for Equal Population Variances, Hypothesis Tests for Difference between Population Proportions, Tests for Normality, Chi-Square Test for Independence.

Textbook 1: Ch. 8,9

RBT: L1, L2, L3

Module – 5

Regression Analysis: Estimating Relationships: Introduction, Scatterplots : Graphing Relationships, Linear versus Nonlinear Relationships,Outliers,Unequal Variance, No Relationship,Correlations:Indications of Linear Relationships, Simple Linear Regression, Least Squares Estimation, Standard Error of Estimate, The Percentage of Variation Explained:R-Square,Multiple Regression, Interpretation of Regression Coefficients, Interpretation of Standard Error of Estimate and R-Square, Modeling Possibilities, Dummy Variables, Interaction Variables, Nonlinear Transformations, Validation of the Fit.

Regression Analysis: Statistical Inference:Introduction,The Statistical Model, Inferences About the Regression Coefficients, Sampling Distribution of the Regression Coefficients, Hypothesis Tests for the Regression Coefficients and p-Values, A Test for the Overall Fit: The ANOVA Table,Multicollinearity,Include/Exclude Decisions, Stepwise Regression,Outliers,Violations of Regression Assumptions,Nonconstant Error Variance,Nonnormality of Residuals,Autocorrelated Residuals ,Prediction.

Textbook 1: Ch. 10,11

RBT: L1, L2, L3

Course outcomes: The students should be able to:

- Explain the importance of data and data analysis
- Interpret the probabilistic models for data

- Define hypothesis, uncertainty principle
- Evaluate regression analysis

Question Paper Pattern:

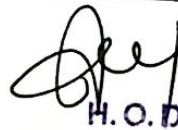
- The question paper will have ten questions.
- Each full Question consisting of 20 marks
- There will be 2 full questions (with a maximum of four sub questions) from each module.
- Each full question will have sub questions covering all the topics under a module.
- The students will have to answer 5 full questions, selecting one full question from each module.

Text Books:

1. S C Albright and W L Winston, Business analytics: data analysis and decision making, 5/e Cenage Learning

Reference Books:

1. ArshdeepBahga, Vijay Madiseti, "Big Data Analytics: A Hands-On Approach", 1st Edition, VPT Publications, 2018. ISBN-13: 978-0996025577
2. Raj Kamal and Preeti Saxena, "Big Data Analytics Introduction to Hadoop, Spark, and Machine-Learning", McGraw Hill Education, 2018 ISBN: 9789353164966, 9353164966



H. O. D.

Dept. Of Computer Science & Engineering
Alva's Institute of Engg. & Technology
Mijar, MOODBISRI - 574 225