

NIRMA UNIVERSITY	
<b>Institute:</b>	Institute of Technology
<b>Name of Programme:</b>	BTech All (Other than CSE)
<b>Course Code:</b>	4CS101ME25
<b>Course Title:</b>	Big Data Systems
<b>Course Type:</b>	Interdisciplinary Minor-Elective
<b>Year of Introduction:</b>	2025-26

L	T	Practical Component				C
		LPW	PW	W	S	
3	0	2	-	-	-	4

### Course Learning Outcomes (CLO):

At the end of the course, the students will be able to –

1. outline the significance and challenges of big data (BL2)
2. model big data applications using various platforms (BL3)
3. utilise big data systems for practical business analytics (BL3)
4. compare data mining algorithms for extracting knowledge from extensive datasets. (BL5)

Unit	Contents	Teaching Hours (Total 45)
Unit-I	<b>Introduction to Big Data and Big Data Storage Platforms:</b> Evolution of Big Data, Types of Big Data, Definition of Big Data, Importance of Big data analytics, Challenges of Conventional Systems, Big data platforms and data storage	06
Unit-II	<b>Hadoop and HDFS:</b> Hadoop Ecosystem, Comparisons of RDBMS and Hadoop, Distributed Computing Challenges, Hadoop Overview, Processing Data with Hadoop, Hadoop YARN	06
Unit-III	<b>MapReduce:</b> working with Map Reduce, Anatomy of a Map Reduce Job Run, Failures, Job Scheduling, Shuffle, and Sort, Task Execution, Map Reduce Types and Formats, Map Reduce Features.	08
Unit-IV	<b>Big data Machine Learning with Spark:</b> Basic concepts of Apache Spark, Spark - RDDs, DataFrames, PySpark, NumPy, SciPy, and Spark ML library, Big data Algorithms for Linear Regression, Clustering, Association rule mining, Decision tree	15
Unit-V	<b>Hadoop Eco-System:</b> CAP Theorem - BASE Concept, NoSQL, Types of No SQL databases, Introduction to MongoDB, Data Types in MongoDB, CRUD	10

### Self-Study:

The self-study contents will be declared at the commencement of the semester. Around 10% of the questions will be asked from self-study contents

### Suggested Readings/ References:

1. Bill Chambers & Matei Zaharis, Spark: The Definitive Guide: Big Data Processing Made Simple, O'Reilly Media, Inc.
2. Michael Berthold, David J. Hand, Intelligent Data Analysis, Springer

3. Chris Eaton, Dirk DeRoos, Tom Deutsch, George Lapis, Paul Zikopoulos, Understanding Big Data: Analytics for Enterprise Class Hadoop and Streaming Data, McGraw Hill
4. Anand Rajaraman and Jeffrey David Ullman, Mining of Massive Datasets, Cambridge University Press
5. Bill Franks, Taming the Big Data Tidal Wave: Finding Opportunities in Huge Data Streams with Advanced Analytics, John Wiley & sons
6. Glenn J. Myatt, Making Sense of Data, John Wiley & Sons
7. Da Ruan, Guoqing Chen, Etienne E.Kerre, GeertWets, Intelligent Data Mining, Springer
8. Paul Zikopoulos, Dirk deRoos, Krishnan Parasuraman, Thomas Deutsch, James Giles, David Corrigan, Harness the Power of Big Data the IBM Big Data Platform, Tata McGraw Hill
9. Michael Minelli, Michele Chambers, Ambiga Dhiraj, Big Data, Big Analytics: Emerging Business Intelligence and Analytic Trends for Today's Businesses, Wiley Publications
10. Zikopoulos, Paul, Chris Eaton, Understanding Big Data: Analytics for Enterprise Class Hadoop and Streaming Data, Tata McGraw Hill
11. Seema Acharya and Subhashini C, Big Data and Analytics, Wiley

### Suggested List of Experiments:

Sr. No.	Title	Hours
1	Study and explore various applications of big data in different domains. Choose one of them and study it in detail; also, write down the report on different types of digital data generated in selected applications. For example, Big Data in Retail, Big Data in Healthcare, Big Data in Education, Big Data in E-commerce, Big Data in Media and Entertainment, Big Data in Finance, Big Data in Travel Industry, Big Data in Telecom, Big Data in Automobile	04
2	Learning limitations of data analytics by applying Machine Learning Techniques to a large amount of data. Write a program to read data sets from any online website, excel file, and CSV file and to perform a) Linear regression and logistic regression on the iris dataset. b) K-means clustering. Students will learn the limitations of platforms and algorithms.	02
3	Set up a single-node Hadoop cluster and apply HDFS commands to the single-node Hadoop Cluster.	04
4	Design MapReduce algorithms to take a very large file of integers and produce as output: a) The largest integer b) The average of all the integers. c) The same set of integers, but each integer appears only once. * d) The count of the number of distinct integers in the input.*	04
5	Apply MapReduce algorithms to find phrase frequency from a given dataset. Prepare a report to guide the design of the mapper and reducer.	02
6	Analyze the impact of different numbers of mappers and reducers on the same definition as practical 4. Prepare a conclusive report on the analysis.	02
7	Implement regression or classification analytic algorithm using MapReduce by handling given datasets using PySpark.	04



- 8 Implement any one of the analytic algorithms of the Clustering application using PySpark. 02
- 9 Set up the MongoDB environment in your system. Import Restaurant Dataset and perform CRUD operation. 04
- 10 Case study: Use various open-source platforms to solve any big data analytic problem of your choice. 02

*Handwritten signature*