



**Sion (East), Mumbai – 400022.  
(Autonomous)**

**Faculty: Science**

**Program: M.Sc.**

**Subject: DATA SCIENCE**

**PART I  
PROPOSED SYLLABUS**

**Academic Year: 2020 – 2021**

## M.Sc(Data Science) – Part I

SEMESTER – I			SEMESTER – II		
Subject Code	Subject Name	Credits	Subject Code	Subject Name	Credits
SIPSDS11	Statistical Methods and Linear Programming	4	SIPSDS21	Advanced Statistical Methods	4
SIPSDS12	Advanced Database Management Systems	4	SIPSDS22	Machine Learning	4
SIPSDS13	Data Mining for Business Intelligence	4	SIPSDS23	Linear Algebra	4
SIPSDS14	Data Science -I	4	SIPSDS24	Research Methodology	4
SIPSDSP11	Statistical Methods and Linear Programming Practical	2	SIPSDSP21	Advanced Statistical Methods practical	2
SIPSDSP12	Advanced Database Management Systems Practical	2	SIPSDSP22	Machine Learning practical	2
SIPSDSP13	Data Mining for Business Intelligence practical	2	SIPSDSP23	Linear Algebra practical	2
SIPSDSP14	Data Science -I practical	2	SIPSDSP24	Research Methodology Practical	2
<b>Total Credits</b>		<b>24</b>	<b>Total Credits</b>		<b>24</b>

# **SEMESTER – I**

## Statistical Methods and Linear Programming

### Learning Objective:

The purpose of this course is to familiarize students with basics of Statistics which is essential for prospective researchers and professionals.

### Learning Outcomes:

- Enable learners to know descriptive statistical concepts
- Enable learners to apply the various distribution methods to data.
- Demonstrate the competency on topics like basics of data science, data transformation, statistical methods, applied probability etc.
- Enable learners to know various statistical models concepts used for the study of Data Science.

### Theory Component:

M. Sc (Data Science)	Semester – I – SIPSDS11
Course Name	Statistical Methods and Linear Programming
Periods per week (1 Period is 60 minutes)	4
Credits (Theory + Internals)	4

Unit	Contents	No. of Lectures
I	<p><b>Data Presentation</b> :Data types : attribute, variable, discrete and continuous variable Data presentation : frequency distribution, histogram, ogive curves, stem and leaf display</p> <p><b>Data Aggregation</b> : Measures of Central tendency: Mean, Median, mode for raw data, discrete, grouped frequency distribution. Measures dispersion: Variance, standard deviation, coefficient of variation for raw data, discrete and grouped frequency distribution, quartiles, quantiles Real life examples</p> <p><b>Moments</b>: raw moments, central moments, relation between raw and central moments</p> <p><b>Measures of Skewness and Kurtosis</b>: based on moments, quartiles, relation between mean, median, mode for symmetric, asymmetric frequency curve.</p>	12
II	<p><b>Linear Regression</b> : fitting of linear regression using least square regression, coefficient of determination, properties of regression coefficients (only statement) Simple Linear Regression, Multiple Linear Regression,</p> <p><b>Classification</b>: logistic regression, Linear discriminant analysis, Quadratic discriminant analysis</p> <p><b>Resampling Methods</b> : Bootstrapping, cross validation,</p> <p><b>Subset Selection</b>: forward, backward, stepwise, best</p>	12

<b>III</b>	<p><b>Correlation and Regression:</b> bivariate data, scatter plot, correlation, nonsense correlation, Karl Pearson's coefficients of correlation, independence.</p> <p><b>Shrinkage:</b> Ridge regression</p> <p>Dimension Reduction: principal components regression, partial least squares.</p> <p><b>Nonlinear Models:</b> step function, piecewise function, splines, generalized additive model,</p> <p><b>Tree-Based Methods:</b> Bagging, Boosting, random forest.</p>	<b>12</b>
<b>IV</b>	<p>Introduction: linear programming, graphical method, simplex method, slack, surplus, artificial variables, Big M method, two Phase Method, conversion from simplex to dual and vice versa, dual simplex method, integer programming problem.</p>	<b>12</b>
<b>V</b>	<p><b>Transportation problem:</b> North west corner method, Least cost entry method, Vogel's approximation method, test for optimality.</p> <p><b>Assignment Problem:</b> mathematical models of assignment problem, Hungarian Method.</p> <p><b>Job sequencing Problem, Programme Evaluation and Review Technique and Critical Path Method (PERT AND CPM).</b></p>	<b>12</b>

### Books and References

Sr. No.	Title	Author/s	Publisher	Edition	Year
1	Probability, Statistics, Design of Experiments and Queuing theory, with applications of Computer Science	Trivedi, K.S.	Prentice Hall of India, New Delhi	2 <sup>nd</sup>	2009
2	Fundamentals of Mathematical Statistics	Gupta, S.C. and Kapoor, V.K.	S. Chand and Sons, New Delhi	11 <sup>th</sup>	2002
3	Applied Statistics	Gupta, S.C. and Kapoor, V.K.	S. Chand and Sons, New Delhi	7 <sup>th</sup>	1999
4	A First course in probability	Ross, S.M.	Pearson	6 <sup>th</sup>	2006

### Additional References :

1. "Probability and Statistics for Engineers", Dr. J. Ravichandran, 2010.
2. "Practical Statistics for Data Science", Peter Bruce, Andrew Bruce, O'Reilly, 2017.
3. "Statistics for Data Science", James D. Miller, Packt, 2017.
4. "Data Analysis with R", Tony Fischetti, 2015.
5. "R for data Science: Import, Tidy, Transform, Visualize and Model Data", Hadley Wickham, Garrett Grolemund.

## Advanced Database Management Systems

**Learning Objective:** To introduce students to the Extended Entity Relationship Model and Object Model, Object-Oriented Databases, Parallel and Distributed Databases and Client-Server Architecture and Databases on the Web and Semi Structured Data

**Learning Outcome:** Students will understand how to implement the Horizontal fragmentation of databases, Vertical fragmentation of database, Creating Replica of database., Create Temporal Database, Inserting and retrieving multimedia objects in database (Image / Audio /Video) and Implement Active database using Triggers.

### Theory Component:

<b>M. Sc (Data Science)</b>	<b>Semester – I – SIPSDS12</b>
<b>Course Name</b>	<b>Advanced Database Management Systems</b>
<b>Periods per week (1 Period is 60 minutes)</b>	<b>4</b>
<b>Credits (Theory + Internals)</b>	<b>4</b>

<b>Unit</b>	<b>Contents</b>	<b>No. of Lectures</b>
<b>I</b>	<b>The Extended Entity Relationship Model and Object Model:</b> The ER model revisited, Motivation for complex data types, User defined abstract data types and structured types, Subclasses, Super classes, Inheritance, Specialization and Generalization, Constraints and characteristics of specialization and Generalization, Relationship types of degree higher than two.	<b>12</b>
<b>II</b>	<b>Object-Oriented Databases:</b> Overview of Object-Oriented concepts, Object identity, Object structure, and type constructors, Encapsulation of operations, Methods, and Persistence, Type hierarchies and Inheritance, Type extents and queries, Complex objects; Database schema design for OODBMS; OQL, Persistent programming languages, OODBMS architecture and storage issues, Transactions and Concurrency control, Example of ODBMS <b>Object Relational and Extended Relational Databases:</b> Database design for an ORDBMS - Nested relations and collections, Storage and access methods, Query processing and Optimization, An overview of SQL3, Implementation issues for extended type, Systems comparison of RDBMS, OODBMS, ORDBMS	<b>12</b>
<b>III</b>	<b>Parallel and Distributed Databases and Client-Server Architecture:</b> Architectures for parallel databases, Parallel query evaluation, Parallelizing individual operations, Sorting, Joins, Distributed database concepts, Data fragmentation, Replication, and allocation techniques for distributed database design, Query processing in distributed databases, Concurrency	<b>12</b>

	control and Recovery in distributed databases. An overview of Client-Server architecture	
<b>IV</b>	<p><b>Databases on the Web and Semi Structured Data:</b> Web interfaces to the Web, Overview of XML, Structure of XML data, Document schema, Querying XML data, Storage of XML data, XML applications, The semi structured data model, Implementation issues, Indexes for text data</p> <p><b>Enhanced Data Models for Advanced Applications:</b> Active database Concepts. Temporal database Concepts, Spatial databases Concepts, Deductive databases and Query processing, Mobile databases, Geographic information systems.</p>	<b>12</b>
<b>V</b>	<p><b>Introduction and Getting Started :</b> Documents, Collections : Dynamic Schemas, Naming, Databases, Getting and Starting MongoDB, Introduction to the MongoDB Shell : Running the Shell, A MongoDB Client, Basic Operations with the Shell, Data Types : Basic Data Types, Dates, Arrays, Embedded Documents, _id and ObjectIds</p> <p><b>Creating, Updating, and Deleting Documents :</b> Inserting and Saving Documents : Batch Insert, Insert Validation, Removing Documents : Remove Speed, Updating Documents : Document Replacement, Using Modifiers, Upserts, Updating Multiple Documents, Returning Updated Documents</p> <p><b>Querying :</b> Introduction to find : Specifying Which Keys to Return, Limitations, Query Criteria : Query Conditionals, OR Queries, \$not , Conditional Semantics, Type-Specific Queries : null, Regular Expressions, Querying Arrays, Querying on Embedded Documents, \$where Queries : Server-Side Scripting, Cursors : Limits, Skips, and Sorts, Avoiding Large Skips, Advanced Query Options, Getting Consistent Results, Immortal Cursors, Database Commands : How Commands Work</p>	<b>12</b>

### Books and References

<b>Sr. No.</b>	<b>Title</b>	<b>Author/s</b>	<b>Publisher</b>	<b>Edition</b>	<b>Year</b>
1	Fundamentals of Database Systems	Elmasri and Navathe,	Pearson Education	4 <sup>th</sup>	2003
2	Database Management Systems	Raghu Ramakrishnan, Johannes Gehrke	McGraw-Hill	2 <sup>nd</sup>	2002
3	Database System Concepts	Korth, Silberchatz, Sudarshan	McGraw-Hill	7 <sup>th</sup>	2019
4	Database Systems, Design, Implementation and Management	Peter Rob and Coronel	Thomson Learning	9 <sup>th</sup>	2010
5	MongoDB: The Definitive Guide	Kristina Chodorow	O'Reilly Media	2 <sup>nd</sup>	2013

## **Data Mining for Business Intelligence**

### **Learning Objective:**

As Business Intelligence is a technology driven process, students will be exposed to various activities like Online Analytical Processing, Data Mining, Querying and Reporting which is prime requisite in business world.

### **Learning Outcome:**

The student becomes an expert to do analysis of complex data. The Business Intelligence concepts helps in accelerating and improving decision making.

### **Theory Component:**

<b>M. Sc (Data Science)</b>	<b>Semester – I – SIPSDS13</b>
<b>Course Name</b>	<b>Data Mining for Business Intelligence</b>
<b>Periods per week (1 Period is 60 minutes)</b>	<b>4</b>
<b>Credits (Theory + Internals)</b>	<b>4</b>

<b>Unit</b>	<b>Contents</b>	<b>No. of Lectures</b>
<b>I</b>	Introduction: What is Data mining?, Why Data Mining?Major Issues in Data Mining Data Objects and Attribute Types, Basic Statistical Descriptions of Data, Data Visualization	<b>12</b>
<b>II</b>	Data Preprocessing: Data Cleaning, Data Integration, Data Reduction, Data transformation and discretization Data Warehousing and Online Analytical Processing: Data Warehouse Modeling, Data warehouse Design and Usage, Implementation	<b>12</b>
<b>III</b>	Data Cube Technology : Concepts, Methods,Multidimensional Data Analysis Mining frequent Patterns, Associations and correlations: Basic Concepts and Methods	<b>12</b>
<b>IV</b>	Advanced Pattern Mining, Classification: Basic Concepts, Advanced Methods Cluster Analysis : Basic Concepts and Methods, Advanced Cluster analysis	<b>12</b>
<b>V</b>	Outlier Detection, Data Mining Trends, Mining Complex data types, Data Mining Applications	<b>12</b>



## Books and References

<b>Sr. No.</b>	<b>Title</b>	<b>Author/s</b>	<b>Publisher</b>	<b>Edition</b>	<b>Year</b>
1	Data Mining: Concepts and Techniques	Jiawei Han, Micheline Kamber, Jian Pei	Morgan Kaufmann	Third	2012
2	Data Mining for Business Intelligence: Concepts, Techniques and Applications	Galit Shmueli, Nitin Patel, Peter Bruce	Wiley	Second	2010
3	Mining of Massive Datasets	Jure Leskovec , Anand Rajaraman, Jeffrey D. Ullman			2014

## Data Science - I

### Learning Objective:

To acquaint learners about the fact that Data is Science in today's world.

### Learning Outcome:

Students will be able to develop models using given data, and use that model to analyze data, predict data with accuracy check which is the key factor when analyzing data.

### Theory Component:

M. Sc (Data Science)	Semester – I – SIPSDS14
Course Name	Data Science - I
Periods per week (1 Period is 60 minutes)	4
Credits(Theory + Internals)	4

Unit	Contents	No. of Lectures
I	<b>Getting Started with R</b> : Installation, Getting started with the R interface <b>R Nuts and Bolts</b> : Entering Input, Evaluation, R Objects, Numbers, Attributes , Creating Vectors, Mixing Objects, Explicit Coercion, Matrices, Lists, Factors, Missing Values, Data Frames , <b>Names</b> <b>Getting Data In and Out of R</b> : Reading and Writing Data, Reading Data Files with read.table(), Reading in Larger Datasets with read.table, Calculating Memory Requirements for R Objects <b>Using the readr Package</b> <b>Using Textual and Binary Formats for Storing Data</b> : Using dput() and dump()	12
II	<b>Interfaces to the Outside World</b> : File Connections, Reading Lines of a Text File, Reading From a URL Connection <b>Subsetting R Objects</b> : Subsetting a Vector, Subsetting a Matrix, Subsetting Lists, Subsetting Nested Elements of a List, Extracting Multiple Elements of a List, Partial Matching, Removing NA Values <b>Vectorized Operations</b> , : Vectorized Matrix Operations <b>Dates and Times</b> : Dates in R, Times in R, Operations on Dates and Times	12
III	<b>Managing Data Frames with the dplyr package</b> : Data Frames, The dplyr Package, dplyr Grammar, Installing the dplyr package, select(), filter(), arrange(), rename(), mutate(), group_by(), %>% <b>Control Structures</b> : if-else, for Loops, Nested for loops, while Loops, repeat Loops, next, break <b>Functions</b> : Functions in R, Your First Function, Argument Matching, Lazy	12

	Evaluation, The ... Argument, Arguments Coming After the ... Argument <b>Scoping Rules of R</b> : A Diversion on Binding Values to Symbol, Scoping Rules, Lexical Scoping: Why Does It Matter?, Lexical vs. Dynamic Scoping, Application: Optimization, Plotting the Likelihood	
IV	<b>Coding Standards for R</b> : Loop Functions, Looping on the Command Line, lapply(), sapply(), split(), Splitting a Data Frame, tapply, apply(), Col/Row Sums and Means, Other Ways to Apply, mapply(), Vectorizing a Function <b>Debugging</b> : Something's Wrong!, Figuring Out What's Wrong, Debugging Tools in R, Using traceback(), Using debug(), Using recover()	12
V	<b>Profiling R Code</b> : Using system.time(), Timing Longer Expressions, The R Profiler, Using summaryRprof() <b>Simulation</b> : Generating Random Numbers, Setting the random number seed, Simulating a Linear Model, Random Sampling <b>Data Analysis Case Study</b> : Changes in Fine Particle Air Pollution in the U.S. : Synopsis, Loading and Processing the Raw Data, Results	12

### Books and References

Sr. No.	Title	Author/s	Publisher	Edition	Year
1	R Programming for Data Science	Roger D Peng		1 <sup>st</sup>	2015
2	Data Science from Scratch	Joel Grus	O'Reilly Media, Inc.	2 <sup>nd</sup>	2019
3	An Introduction to Statistical Learning	Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani	Springer Science & Business Media, 2013	illustrated	2013
4	Practical Statistics for Data Scientists	Peter Bruce, Andrew Bruce	O'Reilly Media, Inc.	3 <sup>rd</sup>	2018

**Practical Component: (SEMESTER I)**

<b>M. Sc (Data Science)</b>	<b>Semester – I – SIPSDSP11</b>
<b>Course Name</b>	<b>Statistical Methods and Linear Programming</b>
<b>Periods per week (1 Period is 60 minutes)</b>	<b>4</b>
<b>Credits</b>	<b>2</b>

**List of Practical: (Implement using R/Python programming language)**

<b>1</b>	Linear Regression
<b>2</b>	Regression and prediction.
<b>3</b>	Classification
<b>4</b>	Resampling
<b>5</b>	Subset Selection
<b>6</b>	Shrinkage
<b>7</b>	Reduction
<b>8</b>	Nonlinear Models
<b>9</b>	Tree-Based Methods
<b>10</b>	Linear programming problem.
<b>11</b>	Transportation problem.
<b>12</b>	Assignment problem.
<b>13</b>	PERT/CPM problem.

<b>M. Sc (Data Science)</b>	<b>Semester – I – SIPSDSP12</b>
<b>Course Name</b>	<b>Advanced Database Management Systems Practical</b>
<b>Periods per week (1 Period is 60 minutes)</b>	<b>4</b>
<b>Credits</b>	<b>2</b>

### List of Practicals:

1	a	<p>Create a global conceptual schema Emp ( Eno, Ename, Address, Email, Salary) and insert 10 records. Divide Emp into vertical fragments Emp1 ( Eno, Ename, Address) and Emp2 ( Eno, Email, Salary) on two different nodes. Fire the following queries:</p> <ol style="list-style-type: none"> <li>Find the salary of an Employee where employee number is known.</li> <li>Find the Email where the employee name is known.</li> <li>Find the employee name and Email where employee number is known.</li> <li>Find the employee name whose salary is &gt; 10000</li> </ol>
	b	<p>Create a global conceptual schema product_log(product_id, product_name, product_desc, cost, profit) and insert 10 records. Divide product_log into vertical fragments product_m4(product_id, product_name, product_desc) and product_m4(product_id, cost, profit) on two different nodes. Fire the following queries:</p> <ol style="list-style-type: none"> <li>Display cost and profit of each product</li> <li>Display product name where profit is less than Rs.20</li> <li>Display product name, details where cost is between 200 to 500</li> <li>Display product name beginning with 'LA' and profit is 10% of product cost</li> </ol>
2	a	<p>Create a global conceptual schema Emp (Eno, Ename, Address, Email, Salary) and insert 10 records. Divide Emp into horizontal fragments using the condition that Emp1 contains tuples with salary &lt; 10000 and Emp2 with 10000 &lt; salary &lt; 20000 on two different nodes. Fire the following queries:</p> <ol style="list-style-type: none"> <li>Find the salary of all employees</li> <li>Find the Email of all employees where salary=15000</li> <li>Find the employee name and Email where employee number is known</li> <li>Find the employee name and address where employee number is known</li> </ol>
	b	<p>Create a global conceptual schema cust_pdtls (cust_id, cust_name, cust_addr) and insert 10 records. Create two more schemas cust_bill(cust_id, cust_mobile, cust_billamt) and cust_totbill(cust_id, cust_totalamt) on two different nodes. Fire the following queries:</p> <ol style="list-style-type: none"> <li>List out the customer name operating more than 2 mobiles.</li> <li>Display the customer name where the total bill is greater than 2000.</li> <li>Display the total bill for all the customers.</li> </ol>

		iv. Display the customer name who is with us for the last 4 months.
3	a	<p>Create a global conceptual schema Emp(Eno;Ename;Address;Email;Salary) and insert 10 records. Store the replication of Emp into two different nodes and fire queries :</p> <ol style="list-style-type: none"> <li>i. Find the salary of all employees.</li> <li>ii. Find the email of all employees where salary = 15000.</li> <li>iii. Find the employee name and email where employee number is known.</li> <li>iv. Find the employee name and address where employee number is known.</li> </ol>
4	a	<p>Using Object Oriented databases create the following types:</p> <ol style="list-style-type: none"> <li>i. AddrType1 (Pincode: number, Street: char, City: char, State: char, No: number)</li> <li>ii. BranchType (address: AddrType1, phone1: integer, phone2: integer)</li> <li>iii. AuthorType (name: char, addr AddrType1)</li> <li>iv. PublisherType (name: char, addr: AddrType1, branches: BranchTableType)</li> <li>v. AuthorListType as varray, which is reference to AuthorType</li> </ol> <p>Next create the following tables:</p> <ol style="list-style-type: none"> <li>i. BranchTableType of BranchType</li> <li>ii. authors of AuthorType</li> <li>iii. books (title: varchar, year: date, published_by ref PublisherType, authors AuthorListType)</li> <li>iv. Publishers of PublisherType</li> </ol> <p><b>Insert 10 records into the above tables and fire the following queries:</b></p> <ol style="list-style-type: none"> <li>i. List all of the authors that have the same address as their publisher</li> <li>ii. List all the authors that have the same pin code as their publisher.</li> <li>iii. List all books that have 2 or more authors.</li> <li>iv. List the title of the book that has the most authors:</li> <li>v. List the name of the publisher that has the most branches.</li> <li>vi. Name of authors who have not published more than a book.</li> <li>vii. all the branches that belong to the publisher 'tata' to the publisher 'joshi'</li> <li>viii. List all the authors who have published more than one book.</li> <li>ix. List all books (title) where the same author appears more than once on the list of authors (assuming that an integrity constraint requiring that the name of an author is unique in a list of authors has not been specified).</li> </ol>
5	a	<p>Using Object Oriented databases, create the following types:</p> <ol style="list-style-type: none"> <li>i. state61(st_code: number, st_name: varchar2, st_district: varchar2, st_pincode: number)</li> <li>ii. contact_detail61(residence_no: number, office_no: number, email: varchar2, fax: number, mobile: number)</li> <li>iii. address61(road_no: varchar2, road_name: varchar2, landmark:varchar, state: state61, contact: contact_detail61)</li> </ol>

		<p>iv. staff61(staff_id: number, staff_name: varchar2, staff_address: address61, staff_deptno: number, staff_sal: number, staff_other: varchar2, dob: date) define method getAge() to calculate age using dob</p> <p>v. dept61(dept_id: number, location: varchar2, dept_name: varchar2, emp: staffTableType)</p> <p>Next create the following tables:</p> <p>i. staffTableType of staff61</p> <p>ii. dpt_refernce of dept61 with nested relation (emp)</p> <p><b>Insert records into the above tables and fire the following queries:</b></p> <p>i. Display staff ID and department name of all employees.</p> <p>ii. How many workers are in particular department.</p> <p>iii. Find department name for particular staff name</p> <p>iv. Display department-wise report</p> <p>v. Display age and birth date of particular employee</p>
6	a	<p>Create a table Employee with attributes employee_id, first_name, last_name, email, hire_date, job_id, salary, resume as clob and picture as blob to insert an employee's picture. Fire the following queries.</p> <p>i. Use of substr and instr function.</p> <p>ii. Use of OUTPUT.PUT_LINE.</p> <p>And also perform the following :</p> <p>i. For appending data into clob datatype.</p> <p>ii. Selecting CLOB Values by Using SQL</p> <p>iii. Removing LOBs</p>
	b	<p>Create a table Emp with the attributes Eno as employee number, Ename as employee name, Eaddress as employee address and photo as employee picture. Also create a table Company with attributes Eno, designation, age. Fire the following queries :</p> <p>i. Find the name and designation of all the employees.</p> <p>ii. Find the name and age of all the employees.</p> <p>iii. Find the name and photo of a particular employee.</p>
7	a	<p>Create a table tbl Emp_Appnt, which stores the account number,name, and valid time say, recruitment data retirement date. Insert 10 records and fire the following queries</p> <p>i. Find all the employees who join the company on 2/3/2001</p> <p>ii. Find all the employees who will retired on 2/3/2001</p>
	b	<p>Create a table tbl_shares, which stores the, name of company, number of shares, and price per share at transaction time. Insert 10 records and fire the following queries.</p> <p>i. Find all the names of a company whose share price is more than Rs.100 at 11:45 A.M.</p>

		ii. Find the name of company which has highest shares price at 5.00 P.M.
8	a	<p>Create a table employee which stores the employee number, employee name, email, address and salary. Create a table log_employee which stores employee number, old salary, updated salary and date.</p> <p>Create the following triggers :</p> <ol style="list-style-type: none"> <li>On insert of an employee record in the employee table, the corresponding va entered in the log_employee table.</li> <li>On update of any record in the employee table, the corresponding record mu the log_employee table.</li> </ol> <p>Insert 10 records and fire the following queries:</p> <ol style="list-style-type: none"> <li>Display the latest salary of all the employees.</li> <li>Display employee name that has got more than 2 user events.</li> <li>Display employee name that has got an increment of 5000 in one increment.</li> <li>Display employee name and salary of all the employees at second increment</li> <li>Display employee name, total salary and total increment.</li> </ol>
9	a	<p>Create table emp (eno, ename, hrs, pno, super_no) and project (pname, pno, thrs, head_no) where thrs is the total hours and is the derived attribute. Its value is the sum of all employees working on that project. eno and pno are primary keys, head_no is foreign key to emp relation. Insert 10 tuples and write triggers to do the following:</p> <ol style="list-style-type: none"> <li>Creating a trigger to insert new employee tuple and display the new total hours from project table.</li> <li>Creating a trigger to change the hrs of existing employee and display the new total hours from project table.</li> <li>Creating a trigger to change the project of an employee and display the new total hours from project table.</li> <li>Creating a trigger to delete the project of an employee.</li> </ol>
	b	<p>Create table stud1 (roll_no,name) and stud2 (roll_no,name) . Insert 10 tuples and write triggers to do the following:</p> <ol style="list-style-type: none"> <li>Create a trigger such that when a student record is inserted into the table stud1, the same record should be inserted into the table stud2.</li> </ol>
	c	<p>Create a table emp(dept_no,eno,ename,salary) and a table dept(dept_no,total_sal) where the employee table stores the list of employees belonging to which department and their respective salaries. The dept table shows the total salary given to all the employees belonging to the same department. Insert 10 tuples and write triggers to do the following:</p> <ol style="list-style-type: none"> <li>Create a trigger such that on insert of record in the emp table the salaries of employees belonging to the same department should get added in the dept table.</li> </ol>



		<ul style="list-style-type: none"> <li>ii. Create a trigger such that if a record is deleted from the emp table then the salary of the respective employee belonging to a specific department should get deducted from the dept table.</li> </ul>
10	a	<p>Create a table employee having dept_id as number datatype and employee_spec as XML datatype(XM_Type). The employee_spec is a schema with attributes emp_id, name, email, acc_no, managerEmail, dataOf Joning. Insert 10 tuples into employee table. Fire the following queries on XML database.</p> <ul style="list-style-type: none"> <li>i. Retrieve the names of employee.</li> <li>ii. Retrieve the acc_no of employees.</li> <li>iii. Retrieve the names, acc_no, email of employees.</li> <li>iv. Update the 3<sup>rd</sup> record from the table and display the name of an employee.</li> <li>v. Delete 4<sup>th</sup> record from the table.</li> </ul>
	b	<p>Create a table candidate having cand_id as varchar2 datatype and biodata as XML datatype ( XML type). The biodata is a schema with attributes <b>Name, address, skill – compskill – 1) language 2) networking, expr – 1) prog 2) prjmgr, objectives. Fire the following queries on XML database</b></p> <ul style="list-style-type: none"> <li>i. Display candidate name who is good in java and having experience more than 5 years</li> <li>ii. Display candidate having project manager level experience</li> <li>iii. Display name and skill of all candidates</li> <li>iv. Delete record for address = borivali</li> <li>v. Update experience of a particular candidate</li> </ul>

<b>M. Sc (Data Science)</b>	<b>Semester – I – SIPS DSP13</b>
<b>Course Name</b>	<b>Data Mining for Business Intelligence Practical</b>
<b>Periods per week (1 Period is 60 minutes)</b>	<b>4</b>
<b>Credits</b>	<b>2</b>

**List of Practicals :**

1	<p>The dataset ToyotaCorolla.xls contains data on used cars on sale during the late summer of 2004 in The Netherlands. It has 1436 records containing details on 38 attributes, including Price, Age, Kilometers, HP, and other specifications.</p> <p>a. Explore the data using the data visualization (matrix plot) capabilities of XLMiner. Which of the pairs among the variables seem to be correlated?</p> <p>b. We plan to analyze the data using various data mining techniques described in future chapters. Prepare the data for use as follows:</p> <p>i. The dataset has two categorical attributes, Fuel Type and Metallic.  (a) Describe how you would convert these to binary variables.  (b) Confirm this using XLMiner’s utility to transform categorical data into dummies.</p>
2	<p>The file ApplianceShipments.xls contains the series of quarterly shipments (in million \$) of U.S. household appliances between 1985 and 1989 (data courtesy of Ken Black).</p> <p>a. Create a well-formatted time plot of the data using Excel.  b. Does there appear to be a quarterly pattern? For a closer view of the patterns, zoom in to the range of 3500–5000 on the y axis.  c. Create four separate lines for Q1, Q2, Q3, and Q4, using Excel. In each, plot a line graph. In Excel, order the data by Q1, Q2, Q3, Q4 (alphabetical sorting will work), and plot them as separate series on the line graph. Zoom in to the range of 3500–5000 on the y axis. Does there appear to be a difference between quarters?  d. Using Excel, create a line graph of the series at a yearly aggregated level (i.e., the total shipments in each year). e. Re-create the above plots using an interactive visualization tool. Make sure to enter the quarter information in a format that is recognized by the software as a date. f. Compare the two processes of generating the line graphs in terms of the effort as well as the quality of the resulting plots. What are the advantages of each?</p>
3	<p>Sales of Toyota Corolla Cars. The file ToyotaCorolla.xls contains data on used cars (Toyota Corollas) on sale during late summer of 2004 in The Netherlands. It has 1436 records containing details on 38 attributes, including Price, Age, Kilometers, HP, and other specifications. The goal will be to predict the price of a used Toyota Corolla based on its specifications.</p> <p>a. Identify the categorical variables.</p>

b. Explain the relationship between a categorical variable and the series of binary dummy variables derived from it.

c. How many dummy binary variables are required to capture the information in a categorical variable with N categories?

d. Using XLMiner’s data utilities, convert the categorical variables in this dataset into dummy binaries, and explain in words, for one record, the values in the derived binary dummies.

e. Use Excel’s correlation command (Tools > DataAnalysis > Correlation menu) to produce a correlation matrix and XLMiner’s matrix plot to obtain a matrix of all scatterplots. Comment on the relationships among variables. 1The data are available at <http://lib.stat.cmu.edu/DASL/Stories/HealthyBreakfast.html>.

4 Predicting Housing Median Prices. The file BostonHousing.xls contains information on over 500 census tracts in Boston, where for each tract 14 variables are recorded. The last column (CAT.MEDV) was derived from MEDV, such that it obtains the value 1 if MEDV>30 and 0 otherwise. Consider the goal of predicting the median value (MEDV) of a tract, given the information in the first 13 columns.

Partition the data into training (60%) and validation (40%) sets.

a. Perform a k-NN prediction with all 13 predictors (ignore the CAT.MEDV column), trying values of k from 1 to 5. Make sure to normalize the data (click “normalize input data”). What is the best k chosen? What does it mean?

b. Predict the MEDV for a tract with the following information, using the best k:

CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD
0.2	0	7	0	0.538	6	62	47	4
TAX	PTRATIO	B	LSTAT					
307	21	360	10					

(Copy this table with the column names to a new worksheet and then in “Score new data” choose “from worksheet.”)

c. Why is the error of the training data zero?

d. Why is the validation data error overly optimistic compared to the error rate when applying this k-NN predictor to new data?

e. If the purpose is to predict MEDV for several thousands of new tracts, what would be the disadvantage of using k-NN prediction? List the operations that the algorithm goes through in order to produce each prediction.

5 Automobile Accidents. The file Accidents.xls contains information on 42,183 actual automobile accidents in 2001 in the United States that involved one of three levels of injury: NO INJURY, INJURY, or FATALITY. For each accident, additional information is recorded, such as day of week, weather conditions, and road type. A firm might be interested in developing a system for quickly classifying the severity of an accident based on initial reports and associated data in the system (some of which rely on GPS-assisted

	<p>reporting). Our goal here is to predict whether an accident just reported will involve an injury (<math>MAX\_SEV\_IR = 1</math> or <math>2</math>) or will not (<math>MAX\_SEV\_IR = 0</math>). For this purpose, create a dummy variable called INJURY that takes the value “yes” if <math>MAX\_SEV\_IR = 1</math> or <math>2</math>, and otherwise “no.”</p> <ol style="list-style-type: none"> <li>a. Using the information in this dataset, if an accident has just been reported and no further information is available, what should the prediction be? (<math>INJURY = Yes</math> or <math>No</math>) Why?</li> <li>b. Select the first 12 records in the dataset and look only at the response (<math>INJURY</math>) and the two predictors <math>WEATHER\_R</math> and <math>TRAF\_CON\_R</math>.       <ol style="list-style-type: none"> <li>i. Create a pivot table that examines <math>INJURY</math> as a function of the 2 predictors for these 12 records. Use all 3 variables in the pivot table as rows/columns, and use counts for the cells.</li> <li>ii. Compute the exact Bayes conditional probabilities of an injury (<math>INJURY = Yes</math>) given the six possible combinations of the predictors.</li> <li>iii. Classify the 12 accidents using these probabilities and a cutoff of 0.5.</li> <li>iv. Compute manually the naive Bayes conditional probability of an injury given <math>WEATHER\_R = 1</math> and <math>TRAF\_CON\_R = 1</math>.</li> <li>v. Run a naive Bayes classifier on the 12 records and 2 predictors using XLMiner. Check detailed report to obtain probabilities and classifications for all 12 records. Compare this to the exact Bayes classification. Are the resulting classifications equivalent? Is the ranking (= ordering) of observations equivalent?</li> </ol> </li> <li>c. Let us now return to the entire dataset. Partition the data into training/validation sets (use XLMiner’s ”automatic” option for partitioning percentages).       <ol style="list-style-type: none"> <li>i. Assuming that no information or initial reports about the accident itself are available at the time of prediction (only location characteristics, weather conditions, etc.), which predictors can we include in the analysis? (Use the Data_Codes sheet.)</li> <li>ii. Run a naive Bayes classifier on the complete training set with the relevant predictors (and <math>INJURY</math> as the response). Note that all predictors are categorical. Show the classification matrix.</li> <li>iii. What is the overall error for the validation set?</li> <li>iv. What is the percent improvement relative to the naive rule (using the validation set)?</li> <li>v. Examine the conditional probabilities output. Why do we get a probability of zero for <math>P(INJURY = No   SPD\_LIM = 5)</math>?</li> </ol> </li> </ol>
6	<p>Car Sales. Consider again the data on used cars (ToyotaCorolla.xls) with 1436 records and details on 38 attributes, including Price, Age, KM, HP, and other specifications. The goal is to predict the price of a used Toyota Corolla based on its specifications.</p> <ol style="list-style-type: none"> <li>a. Use XLMiner’s neural network routine to fit a model using the XLMiner default values for the neural net parameters, except normalizing the data. Record the RMS error for the training data and the validation data. Repeat the process, changing the number of epochs (and only this) to 300, 3000, and 10,000.       <ol style="list-style-type: none"> <li>i. What happens to the RMS error for the training data as the number of epochs increases?</li> <li>ii. What happens to the RMS error for the validation data?</li> <li>iii. Comment on the appropriate number of epochs for the model.</li> </ol> </li> <li>b. Conduct a similar experiment to assess the effect of changing the number of layers in the network as well as the gradient descent step size.</li> </ol>

7	<p>Online Statistics Courses. Consider the data in the file CourseTopics.xls. These data are for purchases of online statistics courses at statistics.com. Each row represents the courses attended by a single customer. The firm wishes to assess alternative sequencings and combinations of courses. Use association rules to analyze these data and interpret several of the resulting rules.</p>
8	<p>University Rankings. The dataset on American College and University Rankings (available from <a href="http://www.dataminingbook.com">www.dataminingbook.com</a>) contains information on 1302 American colleges and universities offering an undergraduate program. For each university there are 17 measurements, including continuous measurements (such as tuition and graduation rate) and categorical measurements (such as location by state and whether it is a private or public school). Note that many records are missing some measurements. Our first goal is to estimate these missing values from “similar” records. This will be done by clustering the complete records and then finding the closest cluster for each of the partial records. The missing values will be imputed from the information in that cluster.</p> <ol style="list-style-type: none"> <li>Remove all records with missing measurements from the dataset (by creating a new worksheet).</li> <li>For all the continuous measurements, run hierarchical clustering using complete linkage and Euclidean distance. Make sure to normalize the measurements. Examine the dendrogram: How many clusters seem reasonable for describing these data?</li> <li>Compare the summary statistics for each cluster and describe each cluster in this context (e.g., “Universities with high tuition, low acceptance rate. . .”). Hint: To obtain cluster statistics for hierarchical clustering, use Excel’s Pivot Table on the Predicted Clusters sheet.</li> <li>Use the categorical measurements that were not used in the analysis (State and Private/Public) to characterize the different clusters. Is there any relationship between the clusters and the categorical information?</li> <li>Can you think of other external information that explains the contents of some or all of these clusters?</li> <li>Consider Tufts University, which is missing some information. Compute the Euclidean distance of this record from each of the clusters that you found above (using only the measurements that you have). Which cluster is it closest to? Impute the missing values for Tufts by taking the average of the cluster on those measurements.</li> </ol>
9	<p>Forecasting Wal-Mart Stock: show plots, summary statistics, and output from fitting an AR(1) model to the series of Wal-Mart daily closing prices between February 2001 and February 2002. (Thanks to Chris Albright for suggesting the use of these data, which are publicly available, e.g., at <a href="http://finance.yahoo.com">http://finance.yahoo.com</a> and are in the file WalMartStock.xls.) Use all the information to answer the following questions.</p> <ol style="list-style-type: none"> <li>Create a time plot of the differenced series.</li> <li>Which of the following is/are relevant for testing whether this stock is a random walk? <ul style="list-style-type: none"> <li>• The autocorrelations of the close prices series</li> <li>• The AR(1) slope coefficient</li> <li>• The AR(1) constant coefficient</li> </ul> </li> <li>Does the AR model indicate that this is a random walk? Explain how you reached your conclusion.</li> </ol>

	<p>d. What are the implications of finding that a time series is a random walk? Choose the correct statement(s) below.</p> <ul style="list-style-type: none"> <li>• It is impossible to obtain useful forecasts of the series.</li> <li>• The series is random.</li> <li>• The changes in the series from one period to the other are random. FIGURE 16.19</li> </ul>
10	<p>Souvenir Sales: The file SouvenirSales.xls contains monthly sales for a souvenir shop at a beach resort town in Queensland, Australia, between 1995 and 2001. [Source: R. J. Hyndman, Time Series Data Library, <a href="http://www.robjhyndman.com/TSDL">http://www.robjhyndman.com/TSDL</a>; accessed on December 20, 2009.] Back in 2001, the store wanted to use the data to forecast sales for the next 12 months (year 2002). They hired an analyst to generate forecasts. The analyst first partitioned the data into training and validation sets, with the validation set containing the last 12 months of data (year 2001). She then fit a regression model to sales, using the training set.</p> <ol style="list-style-type: none"> <li>a. Create a well-formatted time plot of the data.</li> <li>b. Change the scale on the x axis, or on the y axis, or on both to log scale in order to achieve a linear relationship. Select the time plot that seems most linear.</li> <li>c. Comparing the two time plots, what can be said about the type of trend in the data?</li> <li>d. Why were the data partitioned? Partition the data into the training and validation set as explained above.</li> </ol>

<b>M. Sc (Data Science)</b>	<b>Semester – I – SIPSDSP14</b>
<b>Course Name</b>	<b>Data Science – I Practical</b>
<b>Periods per week (1 Period is 60 minutes)</b>	<b>4</b>
<b>Credits</b>	<b>2</b>

**List of Practical:**

**Using various online data Sets available in Kaggle like CaptaincyOne, ToyotaCorolla, airquality etc. perform the following (from Practical 3):**

1	<ul style="list-style-type: none"> <li>a. Reading data files using read.table(), read.csv(). Using readr package to read data files using read_table(), read_csv()</li> <li>b. Storing Data using dump() and dput()</li> <li>c. Reading data using connection interfaces that is using File connections, URL Connections, gzip connection and bzip Connection</li> </ul>
2	<ul style="list-style-type: none"> <li>a. Create a subset of the following types of data : Matrix, List, Data Frames</li> <li>b. Represent Date and Time in R and Perform operations on Dates and Times</li> </ul>
3.	<p>Write for loops to:</p> <ul style="list-style-type: none"> <li>a. Compute the mean of every column in mtcars.</li> <li>b. Determine the type of each column in nycflights13::flights.</li> <li>c. Compute the number of unique values in each column of iris.</li> <li>d. Generate 10 random numbers from distributions with means of -10, 0, 10, and 100.</li> </ul>
4	Manage Data Frames with the dplyr package, use the following functions select(), filter(), arrange(), rename(), mutate(), group_by()
5	Apply built-in and user defined functions on any data set and understand argument matching, lazy evaluation, the ... argument , arguments coming after the ... argument
6	Use the following functions on any data set : lapply(), split(), sapply(), apply(), tapply(), mapply()
7	Generate Random numbers using : rnorm, dnorm, pnorm, rpois and apply the functions summary() and plot() on the generated data
8	Use any data set to show the use of pipeable functions.

## **SEMESTER – II**



## Advanced Statistical Methods

### Learning Objectives:

The purpose of this course is to familiarize students with basics of Statistics, essential for prospective researchers and professionals.

### Learning Outcomes:

- Enable learners to know descriptive statistical concepts
- Enable study of probability concept required for Data Science learners
- Enable learners to know different types statistical testing methods used in daily life.

### Theory Component:

<b>M. Sc (Data Science)</b>	<b>Semester – II – SIPSDS21</b>
<b>Course Name</b>	<b>Advanced Statistical Methods</b>
<b>Periods per week (1 Period is 60 minutes)</b>	<b>4</b>
<b>Credits (Theory + Internals)</b>	<b>4</b>

<b>Unit</b>	<b>Contents</b>	<b>No. of Lectures</b>
<b>I</b>	<b>Standard distributions:</b> random variable; discrete, continuous, expectation and variance of a random variable, pmf, pdf, cdf, reliability, Introduction and properties without proof for following distributions; binomial, normal, chi-square, t, F. examples	<b>12</b>
<b>II</b>	<b>Hypothesis testing:</b> one sided, two sided hypothesis, critical region, p-value, tests based on t, Normal and F, confidence intervals. Analysis of variance : one-way, two-way analysis of variance	<b>12</b>
<b>III</b>	<b>Non-parametric tests:</b> need of non-parametric tests, sign test, Wilcoxon's signed rank test, run test, Kruskal-Walis tests. Post-hoc analysis of one-way analysis of variance : Duncan's test Chi-square test of association	<b>12</b>
<b>IV</b>	Time Series Analysis and Forecasting Economic time series Different components, illustration, additive and multiplicative models, determination of trend, seasonal and cyclical fluctuations. Time-series as discrete parameter stochastic process, auto covariance and autocorrelation functions and their properties. Exploratory time Series analysis, tests for trend and seasonality, exponential and moving average smoothing.	<b>12</b>
<b>V</b>	Detailed study of the stationary processes: (1) moving average (MA), (2) auto regressive (AR), (3) ARMA and (4) AR integrated MA (ARIMA) models. Box-Jenkins models, choice of AR and MA periods. Discussion (without proof) of estimation of mean, auto covariance and autocorrelation functions under large sample theory, estimation of ARIMA model	<b>12</b>

	parameters. Spectral analysis of weakly stationary process, periodogram and correlogram analyses, computations based on Fourier transform, non stationary process, introduction to forecasting	
--	--	--

### Books and References

Sr. No.	Title	Author/s	Publisher	Edition	Year
1	Probability, Statistics, Design of Experiments and Queuing theory, with applications of Computer Science	Trivedi, K.S.	Prentice Hall of India, New Delhi	2 <sup>nd</sup>	2009
2	Fundamentals of Mathematical Statistics	Gupta, S.C. and Kapoor, V.K.	S. Chand and Sons, New Delhi	11 <sup>th</sup>	2002
3	Applied Statistics, S	Gupta, S.C. and Kapoor, V.K.	. Chand and Son's, New Delhi	7 <sup>th</sup>	2002
4	Common statistical tests.	Kulkarni, M.B., Ghatpande, S.B. and Gore, S.D.	Satyajeet Prakashan, Pune	6 <sup>th</sup>	1999

## Machine Learning

### Learning Objective:

- To introduce several fundamental concepts and methods for machine learning.
- To familiarize the students with some basic learning algorithms and techniques and their applications, as well as general questions related to analyzing and handling large data sets.

### Learning Outcome:

The student will be able to:-

- Understand the implementation procedures for the machine learning algorithms.
- Design Java/Python programs for various Learning algorithms and apply appropriate data sets to the Machine Learning algorithms.
- Identify and apply Machine Learning algorithms to solve real world problems.

### Theory Component:

<b>M. Sc (Data Science)</b>	<b>Semester – II – SIPS22</b>
<b>Course Name</b>	<b>Machine Learning</b>
<b>Periods per week (1 Period is 60 minutes)</b>	<b>4</b>
<b>Credits (Theory + Internals)</b>	<b>4</b>

<b>Unit</b>	<b>Contents</b>	<b>No. of Lectures</b>
<b>I</b>	<b>Introduction to Machine Learning:</b> What is machine learning?, Types of learning, Applications of Machine Learning algorithms. <b>Supervised Learning:</b> Learning a Class from Examples, Vapnik-Chervonenkis (VC) Dimension, Probably Approximately Correct (PAC) Learning, Noise, Learning Multiple Classes, Regression, Model Selection and Generalization, Dimensions of a Supervised Machine Learning Algorithm	<b>12</b>
<b>II</b>	<b>Bayesian Decision Theory:</b> Classification, Losses and Risks, Discriminant Functions, Utility Theory, Association Rules <b>Parametric Methods:</b> Maximum Likelihood Estimation, Evaluating an Estimator: Bias and Variance, The Bayes' Estimator, Parametric Classification, Regression <b>Multivariate Methods:</b> Multivariate Data, Parameter Estimation, Estimation of Missing Values, Multivariate Normal Distribution, Multivariate Classification	<b>12</b>
<b>III</b>	<b>Dimensionality Reduction:</b> Subset Selection, Principal Components Analysis, Factor Analysis, Multidimensional Scaling, Linear Discriminant Analysis, Isomap	<b>12</b>

	<p><b>Clustering:</b> Mixture Densities, k-Means Clustering, Expectation-Maximization Algorithm, Hierarchical Clustering</p> <p><b>Non-Parametric Methods:</b> Nonparametric Density Estimation, Generalization to Multivariate Data, Nonparametric Classification, Condensed Nearest Neighbor</p>	
IV	<p><b>Decision Trees:</b> Univariate Trees, Pruning, Rule Extraction from Trees, Learning Rules from Data, Multivariate Trees</p> <p><b>Linear Discrimination:</b> Generalizing the Linear Model, Geometry of the Linear Discriminant, Pairwise Separation, Parametric Discrimination, Logistic Discrimination.</p> <p><b>Bayesian Estimation:</b> Estimating the Parameter of a Distribution, Bayesian Estimation of the Parameters of a Function, Gaussian Processes</p>	12
V	<p><b>Hidden Markov Models:</b> Discrete Markov Processes, Three Basic Problems of HMMs, Evaluation Problem, Finding the State Sequence, Learning Model Parameters</p> <p><b>Graphical Models:</b> Example of Graphical Models, d-Separation, Belief Propagation, Undirected Graphs, Learning the Structure of a Graphical Model</p> <p><b>Reinforcement Learning:</b> Elements of Reinforcement Learning ,Model-Based Learning ,Temporal Difference Learning, Generalization, Partially Observable States, Support Vector Machines</p>	12

### Books and References

Sr. No.	Title	Author/s	Publisher	Edition	Year
1	Introduction to Machine Learning	Ethem Alpaydm	The MIT Press Cambridge	Second Edition	2010
2	UNDERSTANDING MACHINE LEARNING : From Theory to Algorithms	Shai Shalev-Shwartz, Shai Ben-David	Cambridge University Press	First Edition	2014
3	A first course in Machine Learning	Simon Rogers and Girolami	CRC Press	Second Edition	2016
4	Machine Learning	Rudolph Russell			2018
5	Machine Learning: Algorithms and Applications	Mohssen Mohammed, Badruddin Khan, Eihab Bashier	CRC Press		2017
6	Machine Learning: An Applied Mathematics Introduction	Paul Wilmott			2019

## Linear Algebra

### Learning Objectives:

To offer the learner the relevant linear algebra concepts through Data science applications.

### Learning Outcomes:

- Appreciate the relevance of linear algebra in the field of computer science.
- Understand the concepts through program implementation
- Instill a computational thinking while learning linear algebra and linear programming.
- Linear Programming (LP), also known as linear optimization is a mathematical programming technique to obtain the best result or outcome.

### Theory Component:

<b>M. Sc (Data Science)</b>	<b>Semester – II – SIPSDS23</b>
<b>Course Name</b>	<b>Linear Algebra</b>
<b>Periods per week (1 Period is 60 minutes)</b>	<b>4</b>
<b>Credits (Theory + Internals)</b>	<b>4</b>

<b>Unit</b>	<b>Contents</b>	<b>No. of Lectures</b>
<b>I</b>	<b>Field:</b> Introduction to complex numbers, numbers in Python , Abstracting over fields, Playing with GF(2), Vector Space: Vectors are functions, Vector addition, Scalar-vector multiplication, Combining vector addition and scalar multiplication, Dictionary-based representations of vectors, Dot-product, Solving a triangular system of linear equations. Linear combination, Span, The geometry of sets of vectors, Vector spaces, Linear systems, homogeneous and otherwise	<b>12</b>
<b>II</b>	<b>Matrix:</b> Matrices as vectors, Transpose, Matrix-vector and vector-matrix multiplication in terms of linear combinations, Matrix-vector multiplication in terms of dot-products <b>Null space:</b> General description, Computing sparse matrix-vector product, Linear functions, Matrix-matrix multiplication, Inner product and outer product, From function inverse to matrix inverse <b>Basis:</b> Coordinate systems, Two greedy algorithms for finding a set of generators, Minimum Spanning Forest and GF(2), Linear dependence, Basis ,Unique representation, Change of basis, first look, Computational problems involving finding a basis	<b>12</b>
<b>III</b>	<b>Dimension:</b> Dimension and rank, Direct sum, Dimension and linear functions, The annihilator <b>Linear transformations :</b> properties, matrix of a linear transformation, change of basis, range and kernel, rank and nullity, Rank, Nullity theorem	<b>12</b>

	<b>Gaussian elimination:</b> Echelon form, Gaussian elimination over GF(2), Solving a matrix-vector equation using Gaussian elimination, Finding a basis for the null space, Factoring integers,	
<b>IV</b>	<b>Inner Product:</b> The inner product for vectors over the reals, Orthogonality, <b>Orthogonalization:</b> Projection orthogonal to multiple vectors, Projecting orthogonal to mutually orthogonal vectors, Building an orthogonal set of generators, Orthogonal complement	<b>12</b>
<b>V</b>	<b>Eigenvector:</b> Modeling discrete dynamic processes, Diagonalization of the Fibonacci matrix, Eigenvalues and eigenvectors, Coordinate representation in terms of eigenvectors, The Internet worm, Existence of eigenvalues, Markov chains, Modeling a web surfer: Page Rank. <b>Linear Algebra:</b> Applications, vectorized code, image recognition, dimensionality reduction.	<b>12</b>

### Books and References

<b>Sr. No.</b>	<b>Title</b>	<b>Author/s</b>	<b>Publisher</b>	<b>Edition</b>	<b>Year</b>
1	Coding the Matrix Linear Algebra through Applications to Computer Science	PHILIP N. KLEIN	Newtonian Press	1	2013
2	Linear Algebra and Its Applications	Gilbert Strang	Cengage Learning	4 <sup>th</sup>	2007
3	Linear Algebra and Its Applications	David C Lay	Pearson Education India	3 <sup>rd</sup>	2002
4	Linear Algebra and Probability for Computer Science Applications	Ernest Davis, A K Peters	A K Peters	1	2012
5	Operation research	SD Sharama	Kedarnath	2017	2012

## Research Methodology

### Learning Objective:

To develop the aptitude for research and the ability to explore research techniques to solve real world problems

### Learning Outcome:

- The learner will be able to critically analyze, synthesize and solve complex unstructured business and real world problems with scientific approach.
- The learner will develop analytical skills by applying scientific methods.

### Theory Component:

M. Sc (Data Science)	Semester – II – SIPSDS24
Course Name	Research Methodology
Periods per week (1 Period is 60 minutes)	4
Credits (Theory + Internals)	4

Unit	Contents	No. of Lectures
I	<b>Introduction to Research:</b> Objectives of research, Types of Research, Research approaches, Research methods versus methodology, Research Process. Formulation of the research problem: Selecting the problem, Technique involved in defining a problem.	12
II	<b>Research Design:</b> Meaning, Need and Features of a research design, Different research designs, Basic principles of Experimental Designs, Sampling Design: Implications and Steps in Sampling Design, Types of Sampling Designs.	12
III	<b>Data Collection Methods:</b> Primary data and Secondary data, Processing and Analysis of Data, Statistics in research, Sampling theory, Concept of Standard Error, Estimation, Sample size and its determination <b>Testing of hypotheses:</b> Procedure and flow diagram for hypothesis testing, Parametric Tests, Chi-Square Test, Analysis of Variance and Covariance, Non-parametric tests	12
IV	<b>Multivariate analysis techniques:</b> Classification, Variables, Factor Analysis, Path Analysis, Interpretation and Report Writing :Technique and Precaution in interpretation, Report Writing, <b>Use of tools / techniques for Research:</b> methods to search required information effectively, Reference Management Software like Zotero/Mendeley, Software for paper formatting like LaTeX/MS Office, <b>Referencing styles</b>	12

<b>V</b>	<b>Ethical Issues in Research, Plagiarism</b> and Self Plagairism, Avoiding plagiarism, Why cite?, Basics of citation <b>Fundamentals of Patents:</b> What is a patent?, Conditions for grant of patent, Inventions that are not Patentable, Process and Product Patent, Procedure of the process of registration and grant of patents, Transfer and Infringement of Patent Rights, Surrender of Patents, Challenges in Patents	<b>12</b>
----------	---	-----------

### Books and References

Sr. No.	Title	Author/s	Publisher	Edition	Year
1	Research Methodology – Methods and Techniques	C.R.Kothari, Gaurav Garg	New Age	4e	
2	Research Methodology – a step by step guide for beginners	Ranjit Kumar	Sage Publications	3e	2011
3	Research Methodology	Panneerselvam	PHI Learning	2e	2014
4	Business Research Methods	William G.Zikmund, B.J Babin, J.C. Carr, Atanu Adhikari, M.Griffin	Cengage	8e	2016
5	Business Research Methods	Alan Bryman and Emma Bell	Oxford University Press	3e	2011
6	Intellectual Property Rights	Neeraj Pandey, Khushdeep Dharni	PHI Learning		2014
7	The complete guide to referencing and avoiding plagiarism	Colin Neville	Open University Press	2e	2010
8	Cite Right	Charles Lipson	The University of Chicago Press		2006



**Practical Component: (SEMESTER II)**

<b>M. Sc (Data Science)</b>	<b>Semester – II – SIPSDSP21</b>
<b>Course Name</b>	<b>Advanced Statistical Methods Practical</b>
<b>Periods per week (1 Period is 60 minutes)</b>	<b>4</b>
<b>Credits</b>	<b>2</b>

**List of Practical:**

1. Problems based on binomial distribution
2. Problems based on normal distribution
3. Property plotting of binomial distribution
4. Property plotting of normal distribution
5. Plotting pdf, cdf, pmf, for discrete and continuous distribution
6. t test, normal test, F test
7. Analysis of Variance
8. Non parametric tests- I,II
9. Kruskal-Walis tests
10. Wilcoxon's signed rank test
11. Time Series Analysis and Forecasting.
12. Box- Jenkins methodology.
13. Problems based Periodogram and Correlogram

<b>M. Sc (Data Science)</b>	<b>Semester – II – SIPSDSP22</b>
<b>Course Name</b>	<b>Machine Learning Practical</b>
<b>Periods per week (1 Period is 60 minutes)</b>	<b>4</b>
<b>Credits</b>	<b>2</b>

**List of Practical:**

**Implement the following in Java/python using pre-defined data sets.**

<b>1</b>	Implement and demonstrate the FIND-S algorithm for finding the most specific hypothesis based on a given set of training data samples.
<b>2</b>	Implement and demonstrate the Candidate-Elimination algorithm
<b>3</b>	Write a program to demonstrate the working of the decision tree based ID3 algorithm.
<b>4</b>	Write a program to implement the naïve Bayesian classifier for a sample training data set stored as a .CSV file. Compute the accuracy of the classifier, considering few test data sets.
<b>5</b>	Assuming a set of documents that need to be classified, use the naïve Bayesian Classifier model to perform this task.
<b>6</b>	Write a program to construct a Bayesian network considering medical data.
<b>7</b>	Apply EM algorithm to cluster a set of data stored in a .CSV file. Use the same data set for clustering using k-Means algorithm.
<b>8</b>	Write a program to implement k-Nearest Neighbour algorithm.

<b>M. Sc (Data Science)</b>	<b>Semester – II – SIPSDSP23</b>
<b>Course Name</b>	<b>Linear Algebra Practical</b>
<b>Periods per week (1 Period is 60 minutes)</b>	<b>4</b>
<b>Credits</b>	<b>2</b>

**List of Practical:** Implement using R/Python Programming.

1. Write a program which demonstrates the following:
  - a. Addition of two complex numbers
  - b. Displaying the conjugate of a complex number
  - c. Plotting a set of complex numbers
  - d. Creating a new plot by rotating the given number by a degree 90, 180, 270 degrees and also by scaling by a number  $a=1/2$ ,  $a=1/3$ ,  $a=2$  etc.
2. Write a program to do the following:
  - a. Enter two distinct faces as vectors  $u$  and  $v$ .
  - b. Find a new face as a linear combination of  $u$  and  $v$  i.e.  $au+bv$  for  $a$  and  $b$  in  $\mathbb{R}$ .
  - c. Find the average face of the original faces.
3. Write a program to do the following:
  - a. Enter a vector  $u$  as a  $n$ -list
  - b. Enter another vector  $v$  as a  $n$ -list
  - c. Find the vector  $au+bv$  for different values of  $a$  and  $b$
  - d. Find the dot product of  $u$  and  $v$
4. Write a program to do the following:
  - a. Enter an  $r$  by  $c$  matrix  $M$  ( $r$  and  $c$  being positive integers)
  - b. Display  $M$  in matrix format
  - c. Display the rows and columns of the matrix  $M$
  - d. Find the scalar multiplication of  $M$  for a given scalar.
  - e. Find the transpose of the matrix  $M$ .
5. Write a program to do the following:
  - a. Find the vector –matrix multiplication of a  $r$  by  $c$  matrix  $M$  with an  $c$ -vector  $u$ .
  - b. Find the matrix-matrix product of  $M$  with a  $c$  by  $p$  matrix  $N$ .
6. Write a program to enter a matrix and check if it is invertible. If the inverse exists, find the inverse.
7. Write a program to convert a matrix into its row echelon form
8. Write a program to find Eigen values and vectors
9. Write a program to implement gaussian elimination method.
10. Write a program to implement concepts of orthogonalization.

<b>M. Sc (Data Science)</b>	<b>Semester – II – SIPSDSP24</b>
<b>Course Name</b>	<b>Research Methodology</b>
<b>Periods per week (1 Period is 60 minutes)</b>	<b>4</b>
<b>Credits</b>	<b>2</b>

**List of Practical:**

**(Using Google scholar/SPSS/Mendeley/End note etc)**

<b>1</b>	Defining a research problem
<b>2</b>	Literature Review using search tools like google scholar
<b>3</b>	Research design
<b>4</b>	Sampling Design
<b>5</b>	Usage of measurement and scaling techniques
<b>6</b>	Testing of Hypothesis
<b>7</b>	Implement data analysis techniques
<b>8</b>	Writing a research report