

**St. Xavier's College (Autonomous), Kolkata**  
**Department of Statistics**  
**MSc in Data Science**  
**Syllabus**  
**Semester 1 & 2**

Semester	<b>ONE</b>
Paper Number	<b>1</b>
Paper Code	<b>MDSC 4111</b>
Paper Title	<b>Introduction to Data Science and Data Structures</b>
No. of Credits	<b>6</b>
Course description	Composite Paper Module 1: 2 classes/week. Using Excel, R, Tableau. Module 2: 2classes/week. Using Python No. of classes assigned Theory: 4 classes per week Practical: 4 classes per week
Course Objective	At the end of the course, the students should be able to <ol style="list-style-type: none"> <li>1. Compute basic descriptive measures</li> <li>2. Perform exploratory data analysis with descriptive statistics on given data sets.</li> <li>3. Gain experience in creating a data visualisation for an application domain of their choice.</li> <li>4. Critically evaluate and interpret a data visualisation.</li> <li>5. Analyze, evaluate, and implement data structures &amp; algorithms using python.</li> </ol>
Syllabus	<b>Module1: Introduction to Data Science</b> <b>Unit 1:</b> Motivation for data Science: Story telling with data with illustration from different fields. Data abstraction and data wrangling. Storage of data. [4]

	<p>Data: Types of data, scales of measurement. [3]</p> <p>Univariate Data: Descriptive measures related to univariate metric data. [4]</p> <p>Bivariate Data: Descriptive measures related to bivariate metric data: Correlations, linear and polynomial regressions. Descriptive measures related to bivariate categorical data: Measures of associations in a contingency table. [7]</p> <p><b>Unit 2:</b>  Exploratory data analysis: Philosophy of EDA, Basic tools of EDA (plots, graphs and summary statistics). [4]</p> <p>Data Visualization: Basic principles, ideas and tools for data visualization. Visualization of qualitative, quantitative, temporal, spatial and panel data. [4]</p> <p><b>Module2: Data Structures</b></p> <p>Introduction to Data Structures, Arrays, Linked Lists, Stacks, Queues, Binary Trees, Threaded Binary Trees, Binary Search Trees, AVL trees, Sets, Tuples, Dictionaries, Trie  Searching and Sorting algorithms  Basic ideas on complexity analysis, Big-Oh, Big-Omega, Big-Theta notations. [26]</p>
List of Practical	<ol style="list-style-type: none"> <li>1. Data visualization techniques for non-frequency data</li> <li>2. Data visualization techniques for univariate metric data</li> <li>3. Data visualization techniques for bivariate metric data</li> <li>4. Descriptive measures for univariate metric data</li> <li>5. Descriptive measures for bivariate metric data</li> <li>6. Descriptive measures for bivariate categorical data</li> <li>7. Introduction to Python.</li> <li>8. Practical based on Data Structures using Python.</li> </ol>
Reading/Reference Lists	<ol style="list-style-type: none"> <li>1. Jeffrey M. Stanton (2013). Introduction to Data Science.</li> </ol>

	<ol style="list-style-type: none"> <li>2. The Visual Display of Quantitative Information (2nd Edition). E. Tufte. Graphics</li> <li>3. Trevor Hastie, Robert Tibshirani and Jerome Friedman. Elements of Statistical Learning, Second Edition. ISBN 0387952845. 2009. (free online)</li> <li>4. Yule G.U. and Kendall M.G (1994) : An Introduction to the theory of Statistics. 14<sup>th</sup> Edn. Universal Book stall, Delhi.</li> <li>5. Hogg, R.V., Tanis, E.A. and Rao J.M. (2009): Probability and Statistical Inference, Seventh Ed, Pearson Education, New Delhi.</li> <li>6. Data Structures and Algorithms in Python, Michael T Goodrich, Roberto Tamassia, Michael H Goldwasser, John Wiley</li> <li>7. A First Course on Data Structures in Python, Donald R. Sheehy.</li> </ol>
--	--

Semester	<b>ONE</b>
Paper Number	<b>2</b>
Paper Code	<b>MDSC 4112</b>
Paper Title	<b>Probability</b>
No. of Credits	<b>6</b>
Course description	Composite Paper One Module. Applications Using R No. of classes assigned Theory: 4 classes per week Practical: 3 classes per week
Course Objective	At the end of the course, the students are expected to have <ul style="list-style-type: none"> <li>• Knowledge of basic ideas of Probability.</li> <li>• Knowledge of different types of random variables and their probability distributions.</li> <li>• Knowledge of different discrete and continuous standard theoretical distribution and their uses in modelling data through R.</li> <li>• Construction of Mixed distributions and their uses modelling data through R.</li> <li>• Basic knowledge of prior and posterior distributions.</li> </ul>

Syllabus	<p>Introduction to Probability: random experiments, sample space, events and algebra of events. Definitions of Probability – classical, statistical and axiomatic. [5]</p> <p>Conditional Probability: Theorem of compound probability, theorem of total probability, Bayes theorem and its applications, independent events. [4]</p> <p>Random variables and their probability distributions: PMF, PDF and CDF, statement of properties of CDF, Empirical distribution functions and their properties, illustrations and properties of random variables. Moments. Joint, marginal and conditional probability distributions, Joint PMF, PDF and CDF, statement of properties of Joint CDF, independence of variables. Markov's and Chebyshev's inequalities. Mixed random variables. Construction of probability distributions of mixed random variables. [12]</p> <p>Standard Univariate Discrete Theoretical Distributions: Binomial, Poisson, geometric, negative binomial, hypergeometric, uniform (Genesis, Statement of properties and applications). [8]</p> <p>Standard Univariate Continuous Theoretical Distributions: Rectangular, normal, exponential, Cauchy, beta, gamma, lognormal, logistic, double exponential and Pareto (Genesis, Statement of properties and applications). [10]</p> <p>Bivariate Normal Distribution (Genesis, Statement of properties and applications). [4]</p> <p>Truncated Distributions. [5]</p>
List of Practical	<ol style="list-style-type: none"> <li>1. Introduction to R.</li> <li>2. Generation of random samples from given distributions.</li> <li>3. Fitting of probability models by introducing the concepts of Method of moments and Likelihood methods</li> <li>4. Modelling data through fitting of Standard Univariate Discrete Distributions.</li> <li>5. Modelling data through fitting of Standard Univariate Continuous Distributions.</li> </ol>

	6. Modelling data through fitting of truncated distributions. 7. Modelling data through fitting of Mixed Distributions. 8. Simulation based studies: Practical based on Markov's/ Chebyshev's inequalities.
Reading/Reference Lists	1. Ronald E. Walpole; Raymond H. Myers; Sharon L. Myers; Keying E. Ye <i>Probability and Statistics for Engineers and Scientists</i> , by Pearson, Ninth Edition (2013). 2. Sheldon Ross <i>A First Course in Probability</i> , Pearson, Ninth Edition (2018). 3. Prabhanjan N. Tattar, Suresh Ramaiah, B. G. Manjunath, <i>A Course in Statistics with R</i> ; Wiley, (2018).

Semester	<b>ONE</b>
Paper Number	<b>3</b>
Paper Code	<b>MDSC 4113</b>
Paper Title	<b>Linear Algebra and Elements of Statistical Inference</b>
No. of Credits	<b>6</b>
Course description	Composite Paper Module 1: 2 classes/week Module 2: 2classes/week No. of classes assigned Theory: 4 classes per week Practical: 4 classes per week
Course Objective	After completion of the course a student is expected to have an idea of <ul style="list-style-type: none"> <li>• Matrix algebra and determinants.</li> <li>• Vector spaces, subspaces, their dimensions and basis.</li> <li>• Rank of a matrix and systems of linear equations.</li> <li>• Characteristic roots and vectors along with the understanding of classification of quadratic forms.</li> <li>• Applications on least squares and dimension reduction.</li> </ul>

	<ul style="list-style-type: none"> <li>• Understand the concept of an iid sample.</li> <li>• Conceptualise drawing samples from theoretical distributions.</li> <li>• Conceptualise level, size, power of a test and the errors associated with a testing problem.</li> <li>• Applying the results of sampling distributions to build test statistics and critical regions.</li> <li>• Construct confidence intervals for parameters and develop their relation with testing of hypotheses.</li> </ul>
Syllabus	<p><b>Module1: Linear Algebra</b></p> <p><b>Unit 1:</b> Vectors – Concept of a vector, length of a vector, Angle between two vectors, Orthogonal and orthonormal vectors, Linear dependence and independence of vectors, Vector spaces, Spanning set of a vector space, Basis of a vector space, Dimension of a vector space, Projection of a vector on a vector space, Orthogonal Basis, Orthocomplement of a vector space, Gram-Schmidt orthogonalization procedure, Row space &amp; column space of a matrix.</p> <p>Matrices (as a vector of vectors), Square matrices, Matrix operations (Addition, subtraction, multiplication by a scalar and by a matrix, Kronecker Product), Null matrix, Identity matrix, symmetric and skew symmetric matrices, orthogonal matrices, Rank of a matrix, singular and non-singular matrices, A few important results on the rank of a matrix, Inverse of a matrix, idempotent matrices, Elementary Transformations on a matrix, Reduction of a matrix to echelon, and diagonal forms by elementary transformations, QR factorization, Trace of a matrix, Partitioning of matrices and simple properties. [11]</p> <p><b>Unit2:</b> Direct and iterative methods for solving a linear system of equations: Gaussian elimination, LU factorization, Cholesky method, Householder's matrices, Jacobi's method, Gauss-Seidel method. [5]</p> <p>Characteristic roots and Characteristic vector, Properties of characteristic roots, Spectral Decomposition, Singular value decomposition. [5]</p> <p>Quadratic forms: Classification &amp; canonical reduction. [3]</p>

	<p><b>Module 2 : Elements of Statistical Inference</b></p> <p><b>Unit 3: Sampling Distributions</b>  <i>Basic Concepts:</i> Concept of an iid sample, Statistic and its standard error. Drawing of random samples from theoretical distributions. Illustrations with R. (3)  <i>Techniques of Sampling Distributions:</i> Distribution Function, Moment Generating Function and Transformation of variables technique to obtain sampling distribution of statistics. (2)  <i>Basic Sampling Distributions from Univariate Normal Distributions :</i> Chi-square, t and F. Degrees of freedom. Sampling Distributions of sample mean, sample variance, their independence, linear combinations of normal variables. Definitions of Non-central chi-square, t and F distributions. Illustrations through simulations. Illustrations through simulations. (5)   <i>Basic Sampling Distributions from Bivariate Normal Distributions :</i> Sampling distributions of sample correlation coefficient and linear regression coefficients. Illustrations through simulations. (2)</p> <p><b>Unit 4: Elements of Inference Problems</b>  <i>Problems and Paradigms of Inference:</i> Estimation and Testing of Hypotheses Problems. Parametric and Nonparametric Inference. Classical and Bayesian Inference. (3)  <i>Estimation:</i> Basic Criteria of a good estimator – Sufficiency, Unbiasedness, Minimum Variance, Consistency and Efficiency. OPEF. (4)  <i>Interval Estimation:</i> Methods of finding confidence intervals. Shortest confidence intervals. Confidence belts. (3)   <i>Testing of Hypothesis:</i> Null and Alternative Hypothesis. Simple and Composite Hypothesis. Type-1 and Type-2 Errors. Level and Power of a Test. Power Function. p-value of a test. p-hacking. Application to tests of significance. [4]</p>
List of Practical	<ol style="list-style-type: none"> <li>1. Orthogonalisation of vectors.</li> <li>2. Solution of linear system of equations using direct and iterative methods.</li> <li>3. Characteristic roots and vectors, spectral decomposition and singular value decomposition</li> <li>4. Classification of Quadratic Forms</li> <li>5. Drawing random samples from Discrete Distributions</li> </ol>

	<ol style="list-style-type: none"> <li>6. Drawing random samples from Continuous Distributions</li> <li>7. Studying the basic sampling distributions arising from a normal distribution through simulations</li> <li>8. Studying consistency and efficiency of estimators through simulations</li> <li>9. Obtaining empirical level and power of tests through simulations</li> <li>10. Applying standard tests to real life data</li> <li>11. Obtaining confidence intervals of parameters based on real life data</li> </ol>
Reading/Reference Lists	<ol style="list-style-type: none"> <li>1. Stephen Boyd and Lieven Vandenberghe, Introduction to Applied Linear Algebra: Vectors, Matrices, and Least Squares (Cambridge University Press, 3rd edition)</li> <li>2. Lloyd N. Trefethen and David Bau, III: Numerical linear algebra, SIAM (1997)</li> <li>3. Matrix Algebra: Theory, Computations and Applications in Statistics by J.E. Gentle, Springer, 2007</li> <li>4. Fundamentals of Matrix Computations by D.S. Watkins, 2nd ed., Wiley, New York, 2002.</li> <li>5. Gilbert Strang; Linear Algebra and its Applications; Academic Press; Second Edition.</li> <li>6. Goon A.M., Gupta M.K., Das Gupta.B.: Fundamentals of Statistics, Vol. 1, World Press, 2010.</li> <li>7. Goon, A.M. Gupta, M.K. and Dasgupta, B. : An outline of Statistical Theory, Vol. 1, World Press, 2010.</li> <li>8. Ismay, C. and Kim, A.Y., Statistical Inference via Data Science, A Modern Dive into R and the Tidyverse, CRC Press Talor and Francis group, 2020.</li> <li>9. Moulin, P. and Venugopal, V.V., Statistical Inference for Engineers and Data Scientists, Cambridge University Press.</li> <li>10. Caffo, B., Statistical Inference for Data Science, Leanpub, 2016.</li> </ol>

Semester	<b>ONE</b>
Paper Number	<b>4</b>
Paper Code	<b>MDSC 4114</b>
Paper Title	<b>Data Base Management System and Data warehousing</b>
No. of Credits	<b>6</b>
Course description	Composite Paper One Module. Practical using SQL. No. of classes assigned Theory: 4 classes per week



	Practical: 3 classes per week
Course Objective	<p>At the end of the course, the students are expected to have</p> <ul style="list-style-type: none"> <li>• The students will learn the fundamental concepts of database management system</li> <li>• They will learn how to design databases and how to improve the schemas.</li> <li>• The students will also learn about transaction processing, concurrency control and recovery techniques</li> <li>• They would also be familiarized with query processing and optimization techniques</li> <li>• Finally they will be introduced to data warehousing with a discussion on fundamental concepts.</li> </ul>
Syllabus	<p>Fundamental concepts of DBMS: Data Abstraction; Data Models (Entity Relationship Model, Relational Model) Database Languages; Structured Query Language. Database design: Constraints; Referential Integrity; Functional Dependencies; Normalization. File Organization; Indexing and Hashing.</p> <p>Transaction Management, Concurrency Control &amp; Recovery: ACID properties; rollback; recoverable schedules; lock-based protocols; deadlock handling; failure; recovery techniques.</p> <p>Query processing and optimization.</p> <p>Data warehousing: Basic Concepts, OLTP, Data Warehouse Architecture; Schemas; Data Marts; Data Warehouse Design; On-line Analytical Processing (OLAP).</p>
Practical	Practical based on the theory component.
Reading/Reference Lists	<ol style="list-style-type: none"> <li>1. R. Elmasri, S.B. Navathe, Fundamentals of Database Systems 6th Edition, Pearson Education, 2010.</li> <li>2. A. Silberschatz, H.F. Korth, S. Sudarshan, Database System Concepts 6th Edition, McGraw Hill, 2010.</li> <li>3. “Data Mining: Concepts and techniques”, J Han and M Kamber, Third Edition, Elsevier.</li> <li>4. C.J Date, An Introduction to Database systems; Pearson/Addison Wesley; 2003</li> </ol>

Semester	<b>TWO</b>
Paper Number	<b>5</b>
Paper Code	<b>MDSC 4211</b>
Paper Title	<b>Statistical Inference</b>
No. of Credits	<b>6</b>
Course description	Composite Paper Module 1: 2 classes/week Module 2: 2classes/week No. of classes assigned Theory: 4 classes per week Practical: 3 classes per week
Course Objective	At the end of the course, the students should be able to <ol style="list-style-type: none"> <li>1. Apply likelihood and moment methods to obtain estimates of parameters.</li> <li>2. Construct tests of hypotheses based on the likelihood function.</li> <li>3. Apply the chi-square to test goodness of fit and homogeneity on real life data.</li> <li>4. Apply resampling techniques to obtain estimates, standard errors of estimates and confidence intervals of parameters.</li> <li>5. Identify the parameters of a Gauss Markov model.</li> <li>6. Differentiate between ANOVA and regression models.</li> <li>7. Apply the theorems of least squares to carry out tests in ANOVA and regression models and identify worth of a concomitant variable in an ANOCOVA Model.</li> <li>8. Apply nonparametric tests to data where the parent distribution is unknown in structure.</li> <li>9. Differentiate between classical and Bayesian school of thoughts.</li> <li>10. Select suitable priors of parameters and obtain the posteriors.</li> </ol>
Syllabus	<b>Module-1</b> <b>Unit 1 : Parametric Methods (14 L)</b> <i>Parametric Methods:</i> Method of Moments, Maximum Likelihood Estimators. Likelihood Ratio, Rao's Score and Wald Tests. Statements of their large sample properties. Pearsonian Chi-square and its uses. (14) <b>Unit2: Bayesian Inferential Methods and Resampling Techniques(12 L)</b>

	<p><i>Priors and Posteriors:</i> Prior and posterior distributions, conjugate priors. Objective priors. Mixtures of conjugate priors. (6)</p> <p>Concept of Jackknife and Bootstrap. Resampling methods in estimation. Bootstrap Confidence Intervals. Cross-validation studies. Illustrations with R. (6)</p> <p><b>Module-2</b></p> <p><b>Unit 1 : Nonparametric Methods (10 L)</b></p> <p>Basic tests of location and scale. Tests of Goodness of fit, Homogeneity and Associations. (10)</p> <p><b>Unit 2 : Linear Models (16 L)</b></p> <p><i>The Gauss-Markov Model:</i> Least Square Estimators. Normal Equations and their solutions. Best Linear Unbiased Estimators. The Gauss Markov Theorem. Error and Error Variance. (3)</p> <p><i>Linear Models:</i> ANOVA, Regression and ANOCOVA Models and some related testing problems. (10)</p> <p><i>Simultaneous confidence intervals:</i> Bonferroni, Scheffe, Tukey, HSU and Duncan's Methods. Comparisons. (3)</p>
List of Practical	<ol style="list-style-type: none"> <li>1. Parametric methods of estimation.</li> <li>2. Likelihood based tests of hypotheses.</li> <li>3. Application of the chi-square tests to real life data.</li> <li>4. Point and interval estimation using resampling techniques.</li> <li>5. ANOVA, Regression and ANOCOVA Models.</li> <li>6. Nonparametric tests of hypotheses.</li> </ol>
Reading/Reference Lists	<ol style="list-style-type: none"> <li>1) Goon A.M., Gupta M.K., Das Gupta.B.: Fundamentals of Statistics, Vol. 1, World Press, 2010.</li> <li>2) Christensen R., Johnson W., Branscum A., Bayesian Ideas and Data Analysis: An Introduction for scientists and statisticians, Chapman &amp; Hall, 2010.</li> <li>3) Faraway, J., Linear Models with R, CRC Press, Second Edition. 2014.</li> <li>4) Faraway, J., Extending the Linear Model with R, CRC Press, Second Edition. 2016.</li> <li>5) Trevor Hastie, Robert Tibshirani and Jerome Friedman. Elements of Statistical Learning, Second Edition.</li> <li>6) Ismay, C. and Kim, A.Y., Statistical Inference via Data Science, A Modern Dive into R and the Tidyverse, CRC Press Talor and Francis group, 2020.</li> <li>7) Moulin, P. and Venugopal, V.V., Statistical Inference for Engineers and Data Scientists, Cambridge University Press.</li> </ol>

	8) Caffo, B., Statistical Inference for Data Science, Leanpub, 2016. 9) Nonparametric Statistical Inference, Gibbons and Chakraborty, CRC Press, First Edition.
--	--

Semester	<b>TWO</b>
Paper Number	<b>6</b>
Paper Code	<b>MDSC 4212</b>
Paper Title	<b>Multivariate Analysis</b>
No. of Credits	<b>6</b>
Course description	Composite Paper Module 1: Unit 1 ( 2 classes/week) Module 2: Unit 2 ( 2classes/week) No. of classes assigned Theory: 4 classes per week Practical: 3 classes per week
Course objective	At the end of the course, the students should be able to understand <ul style="list-style-type: none"> <li>• Multivariate Probability Distributions.</li> <li>• Sampling distributions of some statistics drawn from Multivariate Normal distribution.</li> <li>• Basic concepts and applications of Copula.</li> <li>• Multivariate Data Visualisation.</li> <li>• Application of multivariate techniques.</li> </ul>
Syllabus	<b>Module-I: Multivariate Probability Distributions</b>  <b>Multivariate Data Visualisation:</b> Mosaic Plots, Scatterplot Matrix, Bivariate qq-plots, Spider Web plots, DD Plots, Parallel coordinate plots, Trellis Displays. [6]

	<p><b>Multivariate Probability Distributions:</b> Random Vector, Mean vector &amp; Dispersion matrix, Probability mass/density functions, Marginal &amp; Conditional distributions, Multiple and partial correlation coefficient, Multinomial Distribution, Dirichlet Distribution, Multivariate Normal distribution and its properties. [12]</p> <p><b>Sampling from Multivariate Normal Distribution:</b> Sampling distribution for mean vector and variance-covariance matrix, Wishart Distribution, Hotelling <math>T^2</math> and Mahalanobis <math>D^2</math>. [4]</p> <p><b>Copula:</b> Definition and basic properties, Multivariate Distribution using Copula Functions. [4]</p> <p><b>Module-2: Multivariate Techniques</b></p> <p>Decomposition of data matrices by factors, Principal Component Analysis, Independent Component Analysis, Factor Analysis, Correspondence Analysis, Canonical Correlation Analysis, Discriminant Analysis, Cluster Analysis, Multidimensional Scaling. [26]</p>
List of Practical	<ol style="list-style-type: none"> <li>1. Analysis of multivariate frequency type data.</li> <li>2. Multivariate Probability distributions.</li> <li>3. Sampling from Multivariate Normal Distributions.</li> <li>4. Copula</li> <li>5. Multivariate Data Visualisation</li> <li>6. Multivariate Techniques.</li> </ol>
Reading/Reference Lists	<ol style="list-style-type: none"> <li>1. Johnson/Wichern; Applied Multivariate Statistical Analysis Sixth Edition, Pearson 2015</li> <li>2. Hardle Wolfgang, Simar Leopold: Applied Multivariate Statistical Analysis, Second Editon, Springer.</li> <li>3. Roger B. Nelsen: An introduction to Copulas, Second Editon, Springer.</li> </ol>

<b>Semester</b>	<b>TWO</b>
Paper Number	<b>7</b>

Paper Code	<b>MDSC 4213</b>
Paper Title	<b>Big Data Analytics</b>
No. of Credits	6
Course Description	Composite Paper One Module Number of classes: Theory – 4 per week Practical – 3 per week
Course Objective	<p>After completion of the course a student is expected to have</p> <ul style="list-style-type: none"> <li>○ Understanding of the challenges of computation related to big data.</li> <li>○ Gaining wholesome knowledge about various computational platforms available for big data analytics.</li> <li>○ Understanding the advantages and disadvantages of the big data analytics platforms, including the software frameworks.</li> <li>○ Gaining wholesome knowledge about parallel computation in various big data analytics platforms.</li> <li>○ Gaining hands-on experience in parallel computing with R and Python.</li> <li>○ Gaining hands-on experience in cloud computing.</li> </ul>
Syllabus	<p><b>Introduction:</b> Examples of big data in natural sciences, engineering, social media, industry, etc. Importance of analysing big data. Limitations of the traditional computational platforms in analysis of big data.</p> <p><b>Scaling of big data analytics platforms:</b> Horizontal and vertical scaling, Peer-to-peer networks, Hadoop, Spark, Berkeley Data Analysis Stack (BDAS), High Performance Computing (HPC) clusters, multicore processors, Graphics Processing Unit (GPU), Field Programmable Gate Arrays (FPGA).</p> <p><b>Distributed computing:</b> Importance of distributed computing for big data, Basic ideas of the communication systems for parallel computing in peer-to-peer networks (Message Passing Interface (MPI)), Hadoop (HDFS, YARN, Map Reduce), Spark, BDAS, (Tachyon+Mesos– improvement over Spark due to more aggressive memory exploitation). Communication systems for vertical scaling – MPI for HPC and multicore processors; CUDA for GPUs, Hardware Descriptive Language (HDL) for FPGA. The concept of cloud computing.</p> <p><b>Comparisons of different big data platforms:</b> communication mechanisms based on scalability, data I/O performance,</p>

	<p>fault tolerance, real-time processing, data size supported, iterative task support.</p> <p><b>Pseudocodes:</b> Illustrative examples of simple pseudocodes of the K-means algorithm in MapReduce, MPI and GPU based platforms.</p>
List of Practical	<ol style="list-style-type: none"> <li>1. Writing and implementing parallel codes in R and / Python.</li> <li>2. Implementing the parallel codes in cloud.</li> <li>3. Implementing Hadoop/Spark's machine learning algorithms in cloud with big data.</li> </ol>
Reference List	<ol style="list-style-type: none"> <li>1. Sourav Mazumder, Robin Singh Bhadoria and Ganesh Chandra Deka (2017), "Distributed Computing in Big Data Analytics", Springer.</li> <li>2. Martin Van Steen and Andrew S Tanenbaums: <i>Distributed Systems</i> 3rd Edition (2017)</li> <li>3. Singh, D. and G. K. Reddy (2014). A Survey on Platforms for Big Data Analytics, <i>Journal of Big Data</i> 1:8, 1–20</li> </ol>

Semester	<b>TWO</b>
Paper Number	<b>8</b>
Paper Code	<b>MDSC 4214</b>
Paper Title	<b>Predictive Analytics</b>
No. of Credits	<b>6</b>
Course Description	<p>Composite Paper</p> <p>One Module</p> <p>No. of classes assigned Theory: 4 classes per week</p> <p>Practical: 3 classes per week</p>
Course Objective	<p>At the end of the course, the students should be able to</p> <ul style="list-style-type: none"> <li>○ Develop the concept of regression and classification</li> <li>○ Understand the different model selection methods</li> </ul>

	<ul style="list-style-type: none"> <li>○ Apply different dimension reduction techniques to real life data</li> <li>○ Analyse current data and make future predictions</li> </ul>
Syllabus	<p><b><i>UNIT 1: Introduction</i></b></p> <p>Diagnostic versus prognostic models. Regression versus classification problems. Bias-variance trade off. [4]</p> <p><b><i>UNIT 2: Linear Regression</i></b></p> <p>Least square method, simple linear regression, multiple linear regression with quantitative &amp; qualitative predictors. Dummy variables, regression diagnostics (Outlier detection, leverage, Influential point, Cook's distance, Model selection via AIC and BIC, adjusted R-Square). K-nearest neighbour regression . [10]</p> <p><b><i>UNIT 3: Classification</i></b></p> <p>Logistic regression, multiple logistic regression, multi-category logistic regression (model, parameter estimation and prediction). Multiclass discriminant analysis. Decision trees (CART and CHAID) [15]</p> <p><b><i>UNIT 4: Model selection and Regularization</i></b></p>



	<p>Subset selection method (forward and backward stepwise selection), Shrinkage methods: Penalized likelihood and Bayesian linear regression; Ridge regression, the LASSO and Elastic NET. Applications of dimension reduction techniques. [15]</p> <p><b><i>UNIT 5: Generalized linear model</i></b></p> <p>Components of GLM, link functions (logit, probit, log link). Fitting of GLM (parameter estimation and prediction). Contingency tables, odds ratio and log linear models. Generalized linear mixed models (Inference for model parameters and predictions). [8]</p>
List of Practical	<ol style="list-style-type: none"> <li>1. Fitting of linear regression and regression diagnostics</li> <li>2. Logistic regression</li> <li>3. Multiclass discriminant analysis</li> <li>4. CART and CHAID</li> <li>5. Subset selection methods</li> <li>6. Ridge Regression</li> <li>7. LASSO</li> <li>8. Elastic Net</li> <li>9. Fitting of Generalized linear models</li> </ol> <p>Fitting of Generalized linear mixed effects model</p>
Reading/Reference Lists	<ol style="list-style-type: none"> <li>1. James, Witten, Hastie and Tibshirani: <i>An Introduction to Statistical Learning</i>. Second edition, Springer.</li> <li>2. Hastie, Tibshirani, Friedman: <i>The Elements of Statistical Learning, Data Mining, Inference and Prediction</i>. Second Edition, Springer Series in Statistics.</li> <li>3. McCullagh, P &amp; Nelder, J.A.(1995), <i>Generalized Linear Models</i>. Chapman and Hall.</li> <li>4. Agresti, A. (2007): <i>An Introduction to Categorical data analysis</i>. Wiley</li> </ol>

