

## DATA MINING LAB

<b>Course Code</b>	23CS3551	<b>Year</b>	III	<b>Semester</b>	I
<b>Course Category</b>	PCC	<b>Branch</b>	CSE	<b>Course Type</b>	Practical
<b>Credits</b>	1.5	<b>L-T-P</b>	0-0-3	<b>Prerequisites</b>	Data Base Management Systems, Python Programming
<b>Continuous Internal Evaluation :</b>	30	<b>Semester End Evaluation:</b>	70	<b>Total Marks:</b>	100

Course Outcomes		
Upon successful completion of the course, the student will be able to:		
CO1	Apply various preprocessing techniques on different datasets for a given problem.	L3
CO2	Implement various experiments in Jupyter Notebook Environment and Colab.	L3
CO3	Develop an effective report based on various learning methods implemented.	L3
CO4	Apply technical knowledge for a given scenario and express with an effective oral communication	L3
CO5	Analyze the outputs and visualizations generated for different datasets.	L4

Syllabus		
Exp No.	CONTENTS	Mapped CO
1	Explore different Tools: Jupyter Notebook, PyTorch, TensorFlow, Google Colab, Kaggle. Explore the different datasets: Kaggle, UCI Machine Learning Repository.	CO1, CO2, CO3, CO4, CO5
2	<b>Apply</b> essential data preprocessing techniques to clean and prepare a given dataset (handling missing values, normalization, encoding categorical variables), and <b>compute</b> descriptive statistics (mean, median, mode, standard deviation) to summarize the data distribution and identify potential outliers or biases.	CO1, CO2, CO3, CO4, CO5
3	<b>Apply</b> the <b>K-Nearest Neighbors (KNN)</b> algorithm for both classification and regression problems. <b>Determine</b> the optimal number of neighbors (K) using <b>cross-validation</b> and <b>evaluate</b> using accuracy and <b>error metrics</b> .	CO1, CO2, CO3, CO4, CO5
4	Implement the Decision Tree algorithm for both a classification problem. Perform parameter tuning (e.g., max depth, min samples split) and evaluate using accuracy, precision, recall, F1-score (for classification)	CO1, CO2, CO3, CO4, CO5
5	<b>Apply</b> the <b>Naïve Bayes classification algorithm</b> on textual or categorical datasets and <b>analyze</b> its robustness using performance metrics like <b>confusion matrix, precision, recall, and accuracy</b> .	CO1, CO2, CO3, CO4, CO5
6	<b>Implement</b> a <b>Simple Perceptron</b> and a <b>Multi-Layer Perceptron (MLP)</b> using the MNIST dataset, and <b>evaluate</b> their classification performance using <b>accuracy, precision, recall, and F1-score</b> .	CO1, CO2, CO3, CO4, CO5

7	<b>Apply the Support Vector Machine (SVM) algorithm for classification tasks on multiple datasets and assess performance using confusion matrices, precision, recall, F1 score, and ROC-AUC curves. Also, visualize the margin and support vectors where applicable.</b>	<b>CO1,CO2,CO3, CO4, CO5</b>
8	<p>a. <b>Implement a Simple Linear Regression model to predict continuous output (e.g., house prices or CO<sub>2</sub> emissions).</b></p> <p>b. <b>Apply Logistic Regression on classification datasets (eg: Breast Cancer, Titanic Survival, or Spam Detection)</b></p>	<b>CO1,CO2,CO3, CO4, CO5</b>
9	<p>a. <b>Implement the K-Means clustering algorithm on synthetic and real-world datasets (e.g., customer segmentation). Evaluate using Silhouette Score, Davies-Bouldin Index, and visual clustering results.</b></p> <p>b. <b>Implement Hierarchical Agglomerative Clustering and compare it with K-Means using dendrograms, linkage criteria, and cluster validation indices.</b></p>	<b>CO1,CO2,CO3, CO4, CO5</b>
10	<b>Implement the Expectation-Maximization (EM) algorithm for clustering on Gaussian Mixture Models. Assess clustering performance using log-likelihood, BIC/AIC, and visual interpretation of clusters on 2D datasets.</b>	<b>CO1,CO2,CO3, CO4, CO5</b>
11	<b>Capstone Project: Design and implement an end-to-end Machine Learning pipeline involving problem identification, dataset selection, preprocessing, algorithm selection (classification, regression, or clustering), model building, parameter tuning, and performance evaluation using appropriate metrics.</b>	<b>CO1,CO2,CO3, CO4, CO5</b>

### Learning Resources

#### Text Books

1. Data Mining concepts and Techniques, 3<sup>rd</sup> edition, Jiawei Han, Michel Kamber, Elsevier, 2011.
2. Machine Learning with Python for Everyone, Mark E.Fenner, First Edition, 2020,Pearson.
3. Machine Learning: A Probabilistic Perspective, Kevin P. Murphy, 2012, MIT Press

#### Reference Books

1. "Machine Learning:An Algorithmic Perspective", Second Edition,Stephen Marsland, CRC Press
2. "Machine Learning in Action",Peter Harrington, DreamTech
3. "Introduction to Data Mining", Pang-Ning Tan, Michel Stenbach, Vipin Kumar, 7<sup>th</sup> Edition, 2019.

#### E-Resources & other digital material

1. <https://www.coursera.org/learn/machine-learning>
2. <https://github.com/atinesh-s/Coursera-Machine-Learning-Stanford>