

## Charotar University of Science and Technology [CHARUSAT]

### Chandubhai S. Patel Institute of Technology [CSPIT]

#### Department of Computer Science & Engineering

Date:23/06/2025

#### Practical List

Academic Year	:	2025-26	Semester	:	5
Course Code	:	CSE303	Course Name	:	Machine Learning

Sr. No.	Aim	Hrs.	CO
<b>1.</b>	<b>Basics of Pandas</b> <ul style="list-style-type: none"> <li><b>1.1 Import CSV File</b> <ul style="list-style-type: none"> <li>• Use pd.read_csv.</li> <li>• Set pandas display options (float format, display width).</li> <li>• Rename 'month_year' to 'measuredate'.</li> <li>• Remove rows with missing 'avgtemp' values.</li> <li>• Display DataFrame info, head, shape, missing values, and summary statistics.</li> </ul> </li> <li><b>1.2 Import Excel File</b> <ul style="list-style-type: none"> <li>• Use pd.read_excel</li> <li>• Rename the 'Year' column to 'metro'.</li> <li>• Strip leading/trailing spaces from the 'metro' column.</li> <li>• Convert year columns (2000-2020) to numeric (float64) using pd.to_numeric with errors='coerce'.</li> <li>• Rename year columns to 'pcGDP' + year (e.g., 'pcGDP2000').</li> <li>• Remove rows where all year columns (2000-2020) are NaN.</li> <li>• Set 'metro' as the index.</li> <li>• Display DataFrame info, head, shape, dtypes, and summary statistics.</li> </ul> </li> <li><b>1.3 Import SQL File</b> <ul style="list-style-type: none"> <li>• Use pymssql and mysql.connector to connect to SQL Server and MySQL databases.</li> <li>• Execute the provided SQL query to fetch columns from the 'studentmath' table.</li> <li>• Use pd.read_sql to load query results into a DataFrame.</li> <li>• Close database connections.</li> <li>• Rearrange columns as specified in the document.</li> <li>• Set 'studentid' as the index.</li> <li>• Replace coded values in columns (e.g., 'famrel', 'freetime') with descriptive labels.</li> <li>• Convert specified columns to 'category' dtype.</li> </ul> </li> </ul>	<b>2</b>	<b>1</b>



	<ul style="list-style-type: none"><li>Calculate value counts or percentages for certain columns (e.g., 'famrel', 'freetime').</li><li>Display DataFrame info, head, dtypes, column counts, and unique value counts.</li></ul>		
2.	<b>Taking Measure of Data through Pandas</b>	2	1
	<ul style="list-style-type: none"><li>Load and inspect the nls97 and covidtotals datasets.</li><li>Set indices for nls97 (personid) and covidtotals (iso_code).</li><li>View column names and data types for both datasets.</li><li>Display the first few rows of both datasets.</li><li>Check the shape and unique index values of both datasets.</li><li>Select specific columns from nls97 using different methods ([]), loc, iloc).</li><li>Select multiple columns from nls97 based on a list of variables.</li><li>Select rows from nls97 using slicing, loc, and iloc.</li><li>Filter nls97 rows based on single and multiple conditions (e.g., nightlyhrssleep &lt;= 4 and childathome &gt;= 3).</li><li>Convert object columns in nls97 to categorical data type.</li><li>Generate frequency distributions for categorical variables in nls97 (e.g., maritalstatus, government responsibility columns).</li><li>Save frequency distributions of categorical columns to a text file.</li><li>Calculate descriptive statistics for continuous variables in covidtotals (e.g., total_cases, total_deaths).</li><li>Compute quantiles for continuous variables in covidtotals.</li><li>Create a histogram of total_cases in covidtotals.</li><li>Set up PandasAI with covidtotals for natural language queries.</li></ul>		
3.	<b>Outliers Multivariate</b>	2	1
	<ul style="list-style-type: none"><li>Load and inspect the covidtotals and nls97 datasets.</li><li>Set indices for covidtotals (iso_code) and nls97 (personid).</li><li>Generate descriptive statistics for covidtotals case data.</li><li>Calculate quantiles for covidtotals continuous variables.</li><li>Compute skewness and kurtosis for covidtotals variables.</li><li>Test covidtotals variables for normality using Shapiro-Wilk test.</li><li>Create Q-Q plots for total_cases and total_cases_pm in covidtotals.</li><li>Identify outliers in covidtotals using interquartile range (IQR) method.</li><li>Create a function to detect outliers across multiple covidtotals columns.</li><li>Examine outliers for total_deaths_pm with demographic variables.</li><li>Plot histograms of total_cases and log-transformed total_cases in covidtotals.</li><li>Generate a correlation matrix for covidtotals cumulative and demographic variables.</li><li>Categorize covidtotals cases and deaths into quantiles and create a crosstab.</li><li>Identify countries with unexpected case-death relationships in covidtotals.</li><li>Create scatterplots of total_cases vs. total_deaths and total_cases_pm vs. total_deaths_pm.</li><li>Filter nls97 for individuals with wage income but zero weeks worked in 2020.</li><li>Check nls97 for individuals ever enrolled in a four-year college.</li></ul>		



	<ul style="list-style-type: none"><li>Identify nls97 individuals with graduate enrollment but no bachelor's enrollment.</li><li>Find nls97 individuals with bachelor's or higher degrees but no four-year college enrollment.</li><li>Detect nls97 individuals with high wage income (3 standard deviations above mean).</li><li>Identify nls97 individuals with significant changes in weeks worked in 2021.</li><li>Create a crosstab of nls97 highest grade completed vs. highest degree for grades.</li></ul>		
<b>4.</b>	<b>Visualization through Matplotlib and Seaborn</b>	<b>2</b>	<b>1</b>
	<p><b>4.1 Histogram of Average Temperatures</b></p> <ul style="list-style-type: none"><li>Create a histogram of avgtemp from landtemps to visualize the distribution, including mean line, skewness (-0.3856), and kurtosis (2.794).</li></ul> <p><b>4.2 Q-Q Plot for Temperature Normality</b></p> <ul style="list-style-type: none"><li>Generate a Q-Q plot for landtemps.avgtemp to assess deviation from a normal distribution.</li></ul> <p><b>4.3 Stacked Histogram of COVID-19 Cases</b></p> <ul style="list-style-type: none"><li>Plot a stacked histogram of total_cases_pm from covidtotals for regions (Oceania/Aus, East Asia, Southern Africa, Western Europe) to compare distributions.</li></ul> <p><b>4.4 Boxplot of SAT Verbal Scores</b></p> <ul style="list-style-type: none"><li>Create a boxplot of nls97.satverbal to identify outliers, annotating median (500), quartiles (430, 570), and outlier thresholds (220, 780).</li></ul> <p><b>4.5 Grouped Boxplots of Weeks Worked</b></p> <ul style="list-style-type: none"><li>Use Seaborn to plot boxplots of nls97.weeksworked21 by highestdegree, ordered by degree level, highlighting no lower outliers for some groups.</li></ul> <p><b>4.6 Violin Plot of SAT Verbal Scores</b></p> <ul style="list-style-type: none"><li>Generate a violin plot of nls97.satverbal to show distribution shape and outliers, annotating quartiles, median, and frequency.</li></ul> <p><b>4.7 Scatter Plot of Latitude vs. Temperature</b></p> <ul style="list-style-type: none"><li>Plot landtemps.latabs vs. avgtemp, coloring points by elevation (low: blue, high: red) to explore bivariate relationships.</li></ul> <p><b>4.8 3D Scatter Plot of Temperature, Latitude, Elevation</b></p> <ul style="list-style-type: none"><li>Create a 3D scatter plot of landtemps variables (elevation, latabs, avgtemp) to visualize multivariate relationships.</li></ul> <p><b>4.9 Line Plot of COVID-19 Cases and Deaths</b></p> <ul style="list-style-type: none"><li>Plot daily new_cases and new_deaths from coviddaily (July 2023–March 2024) to examine trends over time.</li></ul> <p><b>4.10 Heat Map of Correlation Matrix</b></p> <ul style="list-style-type: none"><li>Generate a heat map of covidtotals correlation matrix, focusing on total_cases, total_deaths, total_cases_pm, total_deaths_pm, using Seaborn's coolwarm cmap.</li></ul> <p><b>4.11 Series Operations</b></p> <ul style="list-style-type: none"><li>Accessing Values with Bracket Slicing</li><li>Selecting Values with loc Accessor</li><li>Selecting Values with iloc Accessor</li></ul>		



	<ul style="list-style-type: none"><li>• Generating Summary Statistics with describe</li><li>• Computing Quantiles with quantile</li><li>• Computing Quantiles with quantile</li><li>• Complex Conditional Assignment with np.select</li><li>• Applying User-Defined Functions with apply</li><li>• String Cleaning with str Methods</li><li>• Date Parsing and Manipulation with to_datetime</li></ul>		
<b>5.</b>	<b>Linear Regression</b>  <b>Energy Efficiency Estimation for Smart Buildings</b> You are part of a smart infrastructure team working on developing an energy efficiency model for residential buildings. The dataset includes architectural and environmental features like wall area, roof area, glazing area, orientation, relative compactness, and overall height. Your goal is to predict heating load in kilowatts based on these attributes. You must analyze the role of each feature, check for linearity assumptions, and determine whether simple linear approaches suffice or polynomial transformations are necessary. Investigate both underfitting and overfitting scenarios by comparing training and testing errors. You are expected to critically evaluate and justify your modeling decisions using residual plots, correlation matrices, and error metrics. <b>Key Questions / Analysis / Interpretation to be Evaluated</b> 1. Conduct exploratory data analysis to determine if linear or nonlinear patterns exist. 2. Select relevant features by analyzing correlation and multicollinearity. 3. Build a predictive model for heating load, interpret its coefficients, and explain each feature's influence. 4. Calculate MAE, MSE, RMSE, and R <sup>2</sup> for both training and test sets. 5. Report on any signs of overfitting. 6. Assess bias and variance using learning curves and validate assumptions via residual plots. 7. Recommend whether polynomial transformations or feature interaction terms may improve performance. <b>Supplementary Problems</b> Predict electricity consumption using weather and occupancy data. <b>Key Skills Addressed</b> Regression modeling, error interpretation, residual analysis, transformation handling. <b>Applications</b> Linear regression, when mastered through practical implementations such as energy load prediction, serves as a foundational technique for numerous real-world tasks involving continuous value estimation. <b>Learning Outcome</b> Upon completing this practical: Students will master regression diagnostics, identify modeling pitfalls, and improve estimation accuracy using real-world features. <b>Dataset/Test Data</b> Source: UCI Energy Efficiency dataset	2	2

	<p>Link: <a href="https://archive.ics.uci.edu/dataset/242/energy+efficiency">https://archive.ics.uci.edu/dataset/242/energy+efficiency</a></p> <p><b>Dataset Information</b></p> <p>We perform energy analysis using 12 different building shapes simulated in Ecotect. The buildings differ with respect to the glazing area, the glazing area distribution, and the orientation, amongst other parameters. We simulate various settings as functions of the afore-mentioned characteristics to obtain 768 building shapes. The dataset valued responses. It can also be used as a multi-class classification problem if the response is rounded to the nearest integer.</p> <p><b>Tools/Technology</b></p> <p>Python, pandas, scikit-learn, matplotlib, seaborn.</p>		
<b>6.</b>	<p><b>Logistic regression</b></p> <p><b>Diagnosing Disease from Symptoms</b></p> <p>You are working with a health analytics company to develop a predictive system that flags high-risk patients for chronic illness based on demographic, biometric, and lifestyle attributes. Your dataset contains fields like age, BMI, blood pressure, glucose levels, and behavioral flags (e.g., smoking, alcohol). You are tasked with designing a binary classifier that not only predicts the outcome but also justifies its sensitivity to false positives and false negatives under different decision thresholds. Your solution should deal with class imbalance and emphasize the use of probability-based predictions instead of hard labels.</p> <p><b>Key Questions / Analysis / Interpretation to be Evaluated</b></p> <ol style="list-style-type: none"> <li>1. Which features have the highest impact? How is this validated?</li> <li>2. What does the confusion matrix reveal about the model's classification priorities?</li> <li>3. How does adjusting the decision threshold influence precision and recall?</li> <li>4. Why is ROC-AUC a more appropriate evaluation metric than accuracy in this scenario?</li> <li>5. Is the model biased toward a particular class? How can this bias be detected and addressed?</li> <li>6. How well does the model generalize across multiple validation folds?</li> </ol> <p><b>Supplementary Problems</b></p> <p>Fraud detection in financial transactions.</p> <p><b>Key Skills Addressed</b></p> <p>Binary classification, threshold tuning, evaluation metrics (ROC-AUC, F1), cost-sensitive modelling.</p> <p><b>Applications</b></p> <p>Logistic regression is a fundamental algorithm for binary and multi-class classification tasks, widely adopted in various sectors where decisions are made based on probability-driven outcomes. Its mathematical simplicity, interpretability, and real-time efficiency make it a core skill in the data science and applied machine learning toolkit.</p> <p><b>Learning Outcome</b></p> <p>Upon completing this practical:</p> <p>Students will build interpretable classifiers, handle imbalanced datasets, and make cost-sensitive decisions</p> <p><b>Dataset/Test Data</b></p> <p>SOURCE: PIMA Diabetes dataset / Synthetic healthcare dataset</p>	2	2



	<p>Link: <a href="https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database">https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database</a></p> <p><b>About Dataset</b></p> <p>Context This dataset is originally from the National Institute of Diabetes and Digestive and Kidney Diseases. The objective of the dataset is to diagnostically predict whether or not a patient has diabetes, based on certain diagnostic measurements included in the dataset. Several constraints were placed on the selection of these instances from a larger database. In particular, all patients here are females at least 21 years old of Pima Indian heritage.</p> <p>Content The datasets consists of several medical predictor variables and one target variable, Outcome. Predictor variables includes the number of pregnancies the patient has had, their BMI, insulin level, age, and so on.</p> <p><b>Tools/Technology</b></p> <p>Python, sklearn, imbalanced-learn, seaborn.</p>		
7.	<p><b>K-Nearest Neighbors (KNN)</b></p> <p><b>Recommending Products Based on Browsing Behaviour</b></p> <p>A personalized recommendation system is to be developed to predict which product category a user is likely to interact with next, based on their browsing behavior. The dataset contains features such as time spent per category, last purchased product, time since last login, search keywords, and average rating of viewed products. The task involves building a similarity-based model that identifies the top-k closest users or items and classifies the next likely interaction.</p> <p>The model must evaluate multiple distance metrics (e.g., Euclidean, Manhattan, Cosine) and analyze their influence on classification performance. Additionally, KNN should be extended to a regression context to predict user engagement time or next session duration. Scalability concerns, the effect of dimensionality, and the role of normalization in distance-based models must be critically analyzed.</p> <p><b>Key Questions / Analysis / Interpretation to be Evaluated</b></p> <ol style="list-style-type: none"><li>1. Which distance metric gave the best performance, and what justifies this choice?</li><li>2. How does the value of K influence model accuracy and variance?</li><li>3. At what K-values does the model exhibit overfitting or underfitting?</li><li>4. What insights do misclassified categories provide about potential class overlap?</li><li>5. How were the features normalized, and why is normalization critical in this context?</li><li>6. What are the scalability limitations of this approach when applied to large datasets?</li><li>7. How does the regression model perform (e.g., in predicting session time)? Compare MAE/MSE across K values.</li></ol> <p><b>Supplementary Problems</b></p> <p>News recommendation, personalized course suggestion.</p> <p><b>Key Skills Addressed</b></p> <p>Distance-based classification, multi-class confusion matrix, hyperparameter tuning.</p> <p><b>Applications</b></p> <p>K-Nearest Neighbors (KNN) is a non-parametric, instance-based learning algorithm widely used in both classification and regression tasks. Despite its simplicity, KNN is powerful in applications where pattern similarity, local structure, or case-based reasoning is important.</p>	2	2



	<p><b>Learning Outcome</b> Upon completing this practical: Students will develop and analyze similarity-based classification systems under practical constraints.</p> <p><b>Dataset/Test Data</b> Source: Retail user activity logs / Synthetic behavior data Link: <a href="https://archive.ics.uci.edu/dataset/352/online+retail">https://archive.ics.uci.edu/dataset/352/online+retail</a></p> <p><b>Dataset Information</b> This is a transactional data set which contains all the transactions occurring between 01/12/2010 and 09/12/2011 for a UK-based and registered non-store online retail. The company mainly sells unique all-occasion gifts. Many customers of the company are wholesalers.</p> <p><b>Tools/Technology</b> Python, scikit-learn, seaborn, NumPy.</p>		
<b>8.</b>	<b>Decision Tree &amp; Random Forest</b>	2	3
	<p><b>Customer Churn Prediction using Decision Tree &amp; Random Forest</b> In today's competitive telecom industry, retaining customers is crucial. Customer churn refers to the loss of clients who stop using a company's services. These practical aims to develop a predictive model using Decision Tree and Random Forest algorithms to identify customers likely to churn based on historical data such as service usage, billing information, and customer support interactions. By analyzing these patterns, the model will help the company proactively address issues and improve retention strategies. The key focus is on building accurate and interpretable models, identifying important features, and minimizing overfitting while ensuring better decision-making support for the business.</p> <p><b>Key Questions / Analysis / Interpretation to be Evaluated</b></p> <ol style="list-style-type: none"><li>1. What are the most important features influencing customer churn?</li><li>2. Can a decision tree model accurately classify churners vs non-churners?</li><li>3. Does the Random Forest model reduce overfitting compared to a single decision tree?</li><li>4. How does feature importance differ between the models?</li><li>5. How does model accuracy change with varying depths or number of trees?</li></ol> <p><b>Supplementary Problems</b> Handling missing or imbalanced data.</p> <p><b>Key Skills Addressed</b> Students will gain hands-on experience in handling data imbalance, building and tuning Decision Tree and Random Forest models, analyzing feature importance, and evaluating models to minimize overfitting</p> <p><b>Applications</b> Students will understand how machine learning can predict customer churn, aiding industries like telecom, banking, and e-commerce in reducing customer attrition through data-driven strategies.</p> <p><b>Learning Outcome</b> Upon completing this practical:</p>		



	<p>Students will learn to preprocess data, build interpretable predictive models, identify key factors driving churn, and validate their models effectively for real-world applications.</p> <p><b>Dataset/Test Data</b></p> <p>Source: Telco Customer Churn Dataset (Kaggle) Link: <a href="https://www.kaggle.com/datasets/blastchar/telco-customer-churn">https://www.kaggle.com/datasets/blastchar/telco-customer-churn</a> Features include: tenure, MonthlyCharges, Contract, CustomerServiceCalls, PaymentMethod, etc. Target: Churn (Yes/No).</p> <p><b>Tools/Technology</b></p> <p>Python, pandas, numpy for data handling sklearn for ML models matplotlib, seaborn for visualization Jupyter Notebook or Google Colab.</p>		
<b>9.</b>	<p><b>Support Vector Machine (SVM)</b></p> <p><b>Email Spam Detection using Support Vector Machine (SVM)</b></p> <p>With the rapid increase in email usage, spam messages have become a major nuisance, often carrying malicious links or irrelevant promotions. The objective of this practical is to build a Support Vector Machine (SVM) classifier to distinguish between spam and non-spam (ham) emails based on the textual content of emails. The model will analyze word frequencies, patterns, and metadata to identify whether an email is likely spam. By leveraging key SVM concepts such as margin maximization and kernel tricks, the solution aims to deliver high accuracy and generalization while visualizing how data is separated in high-dimensional space.</p> <p><b>Key Questions / Analysis / Interpretation to be Evaluated</b></p> <ol style="list-style-type: none"><li>1. Can SVM effectively separate spam and non-spam emails based on textual features?</li><li>2. How does changing the kernel (linear, polynomial, RBF) affect classification performance?</li><li>3. What is the optimal decision boundary, and how is it determined?</li><li>4. How does margin maximization improve generalization?</li><li>5. What insights can be gained by visualizing high-dimensional data?</li></ol> <p><b>Supplementary Problems</b></p> <p>Document classification.</p> <p><b>Key Skills Addressed</b></p> <p>Students will be able to create robust spam detection systems that improve email filtering, enhance cybersecurity, and moderate content, addressing challenges across various domains like communication and social media.</p> <p><b>Applications</b></p> <p>Students will understand how machine learning can predict customer churn, aiding industries like telecom, banking, and e-commerce in reducing customer attrition through data-driven strategies.</p> <p><b>Learning Outcome</b></p> <p>Upon completing this practical:</p>	2	3



	<p>Students will develop an understanding of SVM principles, implement text classification pipelines, and interpret the performance and decision boundaries of SVM for high-dimensional data.</p> <p><b>Dataset/Test Data</b> Dataset: SMS Spam Collection Dataset (UCI) Link: <a href="https://archive.ics.uci.edu/dataset/228/sms+spam+collection">https://archive.ics.uci.edu/dataset/228/sms+spam+collection</a> Features extracted from: email/SMS content, frequency of certain keywords, presence of special characters or links. Target: Spam or Ham</p> <p><b>Tools/Technology</b> Python, pandas, numpy for data handling sklearn for ML models matplotlib, seaborn for visualization Jupyter Notebook or Google Colab.</p>		
<b>10.</b>	<b>Model Evaluation &amp; Cross-Validation</b>	2	3
	<p><b>Heart Disease Prediction Using Model Evaluation &amp; Cross-Validation</b> Heart disease is one of the leading causes of death globally. Early and accurate prediction of heart disease based on patient health indicators can assist in timely diagnosis and treatment. This practical aims to build and evaluate a classification model (e.g., Logistic Regression, Decision Tree) using various performance metrics such as accuracy, precision, recall, F1-score, and ROC curve. To ensure robustness and generalizability of the model, k-fold cross-validation will be applied. The goal is to not only build a classifier but also interpret and compare the model's performance using different evaluation techniques.</p> <p><b>Key Questions / Analysis / Interpretation to be Evaluated</b></p> <ol style="list-style-type: none"><li>1. Is the model accurately predicting patients at risk of heart disease?</li><li>2. How do performance metrics vary when using different algorithms?</li><li>3. Which metric (precision, recall, F1) is more important in a medical context?</li><li>4. What does the ROC curve reveal about the model's discrimination ability?</li><li>5. How consistent is the model's performance across k-folds?</li></ol> <p><b>Supplementary Problems</b> Quality assurance in manufacturing.</p> <p><b>Key Skills Addressed</b> These practical enables students to evaluate models with metrics like ROC and F1-score, apply k-fold cross-validation for reliability, and compare algorithms to ensure robust predictions.</p> <p><b>Applications</b> Students will see how predictive models can assist in early detection of heart disease, supporting timely medical interventions and improving healthcare decision-making.</p> <p><b>Learning Outcome</b> Upon completing this practical: Students will learn to train and assess classification models, apply evaluation techniques for reliability, and interpret results, preparing them to solve real-world healthcare problems with machine learning.</p> <p><b>Dataset/Test Data</b> Dataset: UCI Heart Disease Dataset Link: <a href="https://archive.ics.uci.edu/dataset/45/heart+disease">https://archive.ics.uci.edu/dataset/45/heart+disease</a> Features include: age, sex,</p>		



	chest pain type, resting blood pressure, cholesterol, fasting blood sugar, etc. Target: Presence or absence of heart disease <b>Tools/Technology</b> Python pandas, numpy for data handling sklearn for modeling, evaluation, and cross-validation matplotlib, seaborn for plotting metrics and ROC curves Jupyter Notebook / Google Colab.		
<b>11.</b>	<b>K-Means and DBSCAN</b>  <b>Customer Segmentation Using K-Means and DBSCAN</b> Businesses often deal with large, diverse customer bases. Segmenting customers based on behavior allows companies to tailor marketing strategies, improve customer experience, and optimize resources. These practical aims to apply K-Means and DBSCAN clustering algorithms to a retail customer dataset to group customers with similar purchasing behavior. By analyzing features like annual income and spending score, we will explore clustering techniques, evaluate performance using inertia and silhouette scores, and visualize the clusters. This hands-on practical demonstrates how unsupervised learning can uncover hidden patterns in data without labeled outcomes. <b>Key Questions / Analysis / Interpretation to be Evaluated</b> 1. How many natural customer segments are there in the dataset? 2. What is the optimal number of clusters (for K-Means)? 3. How does DBSCAN perform compared to K-Means for density-based clustering? 4. What do inertia and silhouette scores reveal about clustering performance? 5. How can clusters be interpreted for business decision-making? <b>Supplementary Problems</b> Image compression and segmentation. <b>Key Skills Addressed</b> By completing this practical, students will apply centroid-based and density-based clustering techniques, evaluate cluster quality using inertia and silhouette scores, and visualize data to derive meaningful insights for business. <b>Applications</b> Students will learn to segment customers effectively, enabling businesses in retail and e-commerce to design targeted marketing strategies, improve customer satisfaction, and detect anomalies. <b>Learning Outcome</b> Upon completing this practical: Students will understand clustering algorithms, evaluate their performance, and interpret clusters to support business decisions with unsupervised learning techniques. <b>Dataset/Test Data</b> Dataset: Mall Customer Segmentation Dataset (Kaggle) <a href="https://www.kaggle.com/datasets/vjchoudhary7/customer-segmentation-tutorial-in-python">https://www.kaggle.com/datasets/vjchoudhary7/customer-segmentation-tutorial-in-python</a> Features: CustomerID, Gender, Age, Annual Income, Spending Score Optional: Normalize Annual Income and Spending Score for DBSCAN <b>Tools/Technology</b>	2	4



	Python pandas, numpy for data manipulation sklearn for K-Means, DBSCAN, silhouette score matplotlib, seaborn for visualization scipy, PCA from sklearn.decomposition for dimensionality reduction Jupyter Notebook / Google Colab.		
<b>12.</b>	<b>Principal Component Analysis (PCA)</b>	2	4
	<p><b>Face Recognition Feature Reduction using PCA</b></p> <p>Face recognition systems process high-dimensional image data, which increases computational cost and may lead to overfitting. This practical focuses on using Principal Component Analysis (PCA) to reduce the dimensionality of facial image datasets while retaining essential features. The goal is to apply PCA to compress the dataset, visualize it in 2D using principal components, and understand the impact of eigen decomposition and explained variance. This helps in improving model efficiency without significant loss of accuracy. Through this, learners explore how PCA simplifies data while preserving structure and key information.</p> <p><b>Key Questions / Analysis / Interpretation to be Evaluated</b></p> <ol style="list-style-type: none"><li>1. How does PCA reduce the dimensionality of image data?</li><li>2. What percentage of variance is retained by top components?</li><li>3. How are eigenvectors and eigenvalues used in PCA?</li><li>4. How do lower-dimensional visualizations help in understanding data clusters?</li><li>5. What are the trade-offs between compression and information loss?</li></ol> <p><b>Supplementary Problems</b></p> <p>Handwriting recognition.</p> <p><b>Key Skills Addressed</b></p> <p>Through this practical, students will standardize high-dimensional data, apply PCA for feature reduction, and analyze explained variance to balance dimensionality and performance.</p> <p><b>Applications</b></p> <p>Students will explore how dimensionality reduction optimizes face recognition systems and other high-dimensional tasks, reducing computation costs while maintaining essential data for accurate predictions.</p> <p><b>Learning Outcome</b></p> <p>Upon completing this practical:</p> <p>Students will learn to implement PCA, visualize data in reduced dimensions, and evaluate the trade-offs between data compression and accuracy for real-world applications.</p> <p><b>Dataset/Test Data</b></p> <p>Dataset: Olivetti Faces Dataset (sklearn) <a href="https://scikit-learn.org/stable/datasets/real_world.html#olivetti-faces-dataset">https://scikit-learn.org/stable/datasets/real_world.html#olivetti-faces-dataset</a> Features: Pixel values of grayscale facial images Labels (optional): Person identities (used for classification post PCA)</p> <p><b>Tools/Technology</b></p> <p>Python numpy, pandas for data handling sklearn.decomposition.PCA for dimensionality reduction matplotlib, seaborn for 2D plotting and visualizations sklearn.datasets for loading image datasets.</p>		



<b>13.</b>	<b>Artificial Neural Networks (ANN)</b>  <b>Handwritten Digit Recognition using Artificial Neural Networks (ANN)</b> Handwritten digit recognition is a classic problem in the field of pattern recognition and deep learning. These practical aims to build an Artificial Neural Network (ANN) using a Multi-Layer Perceptron (MLP) to classify images of handwritten digits from the MNIST dataset. The model will learn from pixel-level image data through forward propagation, adjusting weights using backpropagation, and applying activation functions such as ReLU and Softmax. The objective is to train the ANN to achieve high accuracy in predicting digits (0–9), understand the working of neural layers, and interpret training behavior using loss and accuracy curves. <b>Key Questions / Analysis / Interpretation to be Evaluated</b> 1. How does the MLP architecture (hidden layers, neurons) affect performance? 2. What role do activation functions play in learning non-linear patterns? 3. How does the network learn through forward and backward passes? 4. What is the effect of learning rate and number of epochs? 5. How well does the trained ANN generalize to unseen digit images? <b>Supplementary Problems</b> Automated postal address reading. <b>Key Skills Addressed</b> This practical will involve training an ANN using forward and backward propagation, understanding activation functions, and tuning hyperparameters to optimize model performance <b>Applications</b> Students will explore how dimensionality reduction optimizes face recognition systems and other high-dimensional tasks, reducing computation costs while maintaining essential data for accurate predictions. <b>Learning Outcome</b> Upon completing this practical: Students will gain the ability to implement and train ANNs, analyze training behavior, and achieve high accuracy in recognizing handwritten digits using modern deep learning frameworks. <b>Dataset/Test Data</b> Dataset: MNIST Handwritten Digit Dataset (60,000 training, 10,000 test images) Features: 28x28 grayscale images of digits (0–9) Target: Digit label (0–9) <b>Tools/Technology</b> Python TensorFlow or PyTorch for ANN implementation Keras (if using TensorFlow) for high-level APIs matplotlib, seaborn for plotting training metrics Jupyter Notebook / Google Colab for hands-on environment.	<b>2</b>	<b>5</b>
<b>14.</b>	<b>Convolutional Neural Network (CNN)</b>  To implement a Convolutional Neural Network (CNN) for the binary classification of images of dogs and cats. The goal is to build a model that can accurately distinguish between images of dogs and cats. In this case study, we will use the Kaggle Dogs vs. Cats dataset, which consists of 25,000 labeled images of dogs and cats. The dataset can be downloaded from the following link: <a href="https://www.kaggle.com/c/dogs-vs-cats/data">https://www.kaggle.com/c/dogs-vs-cats/data</a> Dataset:	<b>2</b>	<b>5</b>



- Training set: 25,000 images (12,500 dogs and 12,500 cats)
- Test set: To be split from the training set or use the provided test set for evaluation.

**Tasks to be Performed:**

1. Dataset Preparation:
    - Download and extract the dataset.
    - Split the dataset into training and validation sets.
    - Preprocess the images (resizing, normalization, and augmentation).
  2. Model Building:
    - Define the architecture of the CNN.
    - Choose appropriate layers (convolutional layers, pooling layers, dense layers, etc.).
    - Implement the model using a deep learning framework (e.g., TensorFlow/Keras, PyTorch).
  3. Model Training:
    - Compile the model with appropriate loss function and optimizer.
    - Train the model on the training set and validate it on the validation set.
    - Monitor the training process using metrics like accuracy and loss.
  4. Model Evaluation:
    - Evaluate the model's performance on the test set.
    - Analyze the results using confusion matrix, precision, recall, and F1-score.
  5. Model Optimization:
    - Implement techniques to improve model accuracy (e.g., hyperparameter tuning, regularization, dropout).
    - Retrain and evaluate the optimized model.
6. Model Deployment (Optional):
  - Save the trained model.
  - Implement a simple application to use the model for real-time image classification.

**Key Questions / Analysis / Interpretation to be Evaluated**

1. Understanding the Dataset:
  - How is the dataset structured, and what preprocessing steps were necessary?
  - What challenges did you encounter during data preprocessing?
2. Model Architecture:
  - What architecture did you choose for the CNN, and why?
  - How did you decide on the number of layers and their types?
3. Training Process:
  - What loss function and optimizer were used, and why?
  - How did you split the dataset for training and validation?
  - What metrics were used to monitor the training process?
4. Model Performance:
  - What was the accuracy of the model on the validation and test sets?
  - What do the confusion matrix and other evaluation metrics indicate about the model's performance?
5. Optimization Techniques:
  - What optimization techniques did you implement, and how did they affect the model's performance?
  - What were the best hyperparameters found during tuning?
6. Deployment and Application:

	<ul style="list-style-type: none"> <li>• How can the trained model be deployed for practical use?</li> <li>• What are the potential applications of the model in real-world scenarios?</li> </ul> <p>7. Reflection:</p> <ul style="list-style-type: none"> <li>• What were the main challenges you faced during the implementation of the CNN?</li> <li>• How would you improve the project if given more time or resources?</li> </ul> <p><b>Supplementary Problems</b></p> <p>3 class problem</p> <p><b>Key Skills Addressed</b></p> <p>Training and evaluating deep learning models Visualizing and interpreting model performance</p> <p>Performing hyperparameter tuning (epochs, learning rate, batch size)</p> <p>Applications</p> <p>Image Classification</p> <p><b>Learning Outcome</b></p> <p>Implement an CNN using modern frameworks like Keras or PyTorch. Apply forward and backward propagation effectively during training. Analyze training curves and optimize model parameters. Evaluate and interpret classification results on real-world datasets.</p> <p><b>Dataset/Test Data</b></p> <p>The dataset can be downloaded from the following link:  <a href="https://www.kaggle.com/c/dogs-vs-cats/data">https://www.kaggle.com/c/dogs-vs-cats/data</a></p> <p>Dataset:</p> <ul style="list-style-type: none"> <li>• Training set: 25,000 images (12,500 dogs and 12,500 cats)</li> <li>• Test set: To be split from the training set or use the provided test set for evaluation.</li> </ul> <p><b>Tools/Technology</b></p> <p>Keras, matplotlib</p>		
<b>14.</b>	<p><b>Q-Learning algorithm</b></p> <p>To implement the Q-Learning algorithm for solving a reinforcement learning problem and evaluate its performance.</p> <p><b>Tasks to be Performed</b></p> <ol style="list-style-type: none"> <li>1. Setup Environment: <ul style="list-style-type: none"> <li>• Install the OpenAI Gym toolkit.</li> <li>• Load the Frozen Lake environment.</li> <li>• Understand the environment's state and action space.</li> </ul> </li> <li>2. Implement Q-Learning Algorithm: <ul style="list-style-type: none"> <li>• Initialize the Q-table with zeros.</li> <li>• Define the hyperparameters: learning rate (alpha), discount factor (gamma), and exploration rate (epsilon).</li> <li>• Implement the Q-Learning update rule.</li> <li>• Implement an epsilon-greedy policy for action selection.</li> </ul> </li> <li>3. Training the Agent: <ul style="list-style-type: none"> <li>• Run the Q-Learning algorithm for a fixed number of episodes.</li> <li>• Update the Q-values based on the experiences gained by the agent.</li> <li>• Monitor the agent's performance over time.</li> </ul> </li> <li>4. Evaluation: <ul style="list-style-type: none"> <li>• Test the trained Q-Learning agent in the environment.</li> </ul> </li> </ol>	2	5



	<ul style="list-style-type: none"><li>• Measure the agent's performance using metrics such as average reward per episode.</li></ul> <p>5. Analysis:</p> <ul style="list-style-type: none"><li>• Analyze how different hyperparameters affect the learning process.</li><li>• Compare the performance of the Q-Learning agent with other baseline methods if available.</li></ul> <p><b>Key Questions / Analysis / Interpretation to be Evaluated</b></p> <ol style="list-style-type: none"><li>1. How does the Q-Learning algorithm update the Q-values?</li><li>2. What is the role of the learning rate (<math>\alpha</math>) and how does it affect the learning process?</li><li>3. How does the discount factor (<math>\gamma</math>) influence the Q-Learning algorithm?</li><li>4. What is the epsilon-greedy policy and why is it used in Q-Learning?</li><li>5. How does the exploration rate (<math>\epsilon</math>) impact the agent's learning and performance?</li><li>6. What challenges did you encounter while implementing the Q-Learning algorithm and how did you address them?</li><li>7. How does the performance of the Q-Learning agent change with different values of <math>\alpha</math>, <math>\gamma</math>, and <math>\epsilon</math>?</li><li>8. Can the Q-Learning algorithm be applied to other environments? If yes, what changes would be required in the implementation?</li><li>9. Discuss the importance of choosing the right hyper parameters for the Q-Learning algorithm.</li><li>10. How does the Q-Learning algorithm compare to other reinforcement learning algorithms?</li><li>11. Analyse the reward value for the different actions.</li></ol> <p><b>Supplementary Problems</b></p> <p>Frozen lake problem</p> <p><b>Key Skills Addressed</b></p> <p>Implementation of Q-Learning algorithm Hyperparameter tuning (<math>\alpha</math>, <math>\gamma</math>, <math>\epsilon</math>) Policy design using epsilon-greedy strategy</p> <p><b>Applications</b></p> <p>Spam detection News topic classification Social media content moderation Customer feedback analysis</p> <p><b>Learning Outcome</b></p> <p>Students will understand and implement a tabular Q-Learning algorithm Students will analyze how learning parameters affect agent performance Students will develop the ability to design, train, and evaluate RL agents</p> <p><b>Dataset/Test Data</b></p> <p>Environment: Frozen Lake v1 from OpenAI Gym Predefined 4x4 or 8x8 grid map with slippery surface, holes, and a goal No external dataset; environment provides states, actions, and rewards dynamically.</p>	
--	--	--