

Department of Artificial Intelligence, SVNIT,SURAT
B.Tech-III ,SEM-V
Subject- Machine Learning(AI301)

LAB ASSIGNMENT-4

1.	<p>Objectives: Handling Class Imbalance and Hypothesis Testing</p> <p>Dataset: Use the Credit Card Fraud Detection dataset (available on Kaggle).</p> <ul style="list-style-type: none">• 284,807 transactions (rows).• “Class” column = 0 (non-fraud) or 1 (fraud).• Fraud cases are ~0.17% → strong class imbalance. <p>Part A : Perform EDA on dataset</p> <p>Part B : Exploring Class Imbalance</p> <ol style="list-style-type: none">1. Plot class distribution (pie chart / bar plot).2. Discuss why class imbalance is problematic for machine learning models. <p>Part C: Under-sampling and Over-sampling</p> <ol style="list-style-type: none">1. Apply Random Under-Sampling (RUS): balance classes by downsampling majority class.2. Apply Random Over-Sampling (ROS) or SMOTE (Synthetic Minority Over-sampling Technique). <p>Part D: Discussion & Insights</p> <ol style="list-style-type: none">1. Which method (under/over-sampling) worked best and why?2. Which metrics (precision, recall, F1) are most useful for fraud detection?3. How does hypothesis testing support feature selection in imbalanced datasets?
2.	<p>Objectives:</p> <ul style="list-style-type: none">• Understand the concept of simple and multiple linear regression.• Perform exploratory data analysis (EDA) before modeling.• Train and evaluate linear regression models.• Interpret coefficients and residuals. <p>Dataset : Boston housing</p> <p>Part A: Exploratory Data Analysis</p> <ol style="list-style-type: none">1. Load the dataset (e.g., Boston Housing or any regression dataset).2. Plot the distribution of the target variable. Comment on skewness.3. Compute and visualize the correlation matrix (heatmap). Which features are most correlated with the target? <p>Part B: Simple Linear Regression</p>

1. Select **one predictor feature** (e.g., "average number of rooms" in housing dataset).
2. Formulate the hypothesis equation:
$$y = \beta_0 + \beta_1 x + \epsilon$$
3. Fit a simple linear regression model. Report:
 - Intercept (β_0)
 - Slope (β_1)
4. Plot scatter plot of feature vs target with the regression line.
5. Interpret the slope: What does a unit increase in x mean for y?

Part C: Multiple Linear Regression

1. Select at least **3 predictor variables**.
2. Fit a multiple linear regression model. Write down the fitted equation.
3. Report model coefficients and interpret any one of them.
4. Evaluate the model using:
 - R^2 score
 - Mean Squared Error (MSE)
 - Root Mean Squared Error (RMSE)
5. Compare performance with the simple regression model. Which one is better? Why?