

SCHOOL OF COMPUTER SCIENCE AND ARTIFICIAL INTELLIGENCE		DEPARTMENT OF COMPUTER SCIENCE ENGINEERING	
Program Name: B. Tech		Assignment Type: Lab	Academic Year: 2025-26
Course Coordinator Name		Dr.Vairachilai Shenbagavel	
Instructor(s) Name		Srinivas Komakula	
Course Code	23CA201SE402	Course Title	Explainable AI (P)
Year/Sem	III/V	Regulation	R24
Date and Day of Assignment	28-07-2025	Time(s)	09:00AM -05:00PM
Duration	2 Hours	Applicable to Batches	23CSBTB35
Assignment Number: 01			

Q. No.	Question	Expected Time to complete
1	Book Heaven – Online Bookstore	

Context:

Book Heaven runs digital ads on Google and tracks their impact on weekly book orders.

Google Ads (₹1000s) (x)	Books Sold (y)
1	100
2	130
3	160
1	110
2	140

Objective:

Analyze the effect of Google ad spending on weekly book sales for BookHeaven by performing Linear Regression and interpreting SHAP values.

Requirements:

1. Perform Linear Regression Analysis

- Use the given dataset where:
 - **Independent Variable (x):** Google Ads (in ₹1000s)
 - **Dependent Variable (y):** Books Sold

2. Calculate the Baseline Value

- Compute the **mean of all book sales (y values)**.

3. Calculate SHAP Values

- For each record, calculate the difference between the **predicted value** and the **baseline**.
- This difference is the **SHAP value**, attributed to the amount spent on Google Ads.

4. Compute Final Prediction

- Use the **linear regression model** to calculate predicted book sales for each ad spend value.
- Confirm that:

$$\text{Final Prediction} = \text{Baseline} + \text{SHAP Value}$$

$$\text{Final Prediction} = \text{Baseline} + \text{SHAP Value}$$

5. Interpret the Results

- Explain how the ad spending influenced each predicted sales value.
- Compare the predicted value to the actual value for each row.
- Identify **under prediction** or **over prediction**, and provide reasoning.

Deliverables:

A notebook or document containing:

- Linear regression implementation with coefficients
- Baseline (mean of y)
- Table of SHAP values and predictions
- Explanation of how each input influenced the prediction
- Comparison of predicted vs actual values, with over/under prediction notes
- Summary analysis covering:
 - Accuracy of the model
 - Trend analysis
 - SHAP interpretation insights

Q. No.	Question	Expected Time to complete
2	BookHeaven – Bookstore Sales Prediction using Multiple Linear Regression and SHAP Analysis	

Objective: Understand how customer footfall and promotional activities impact daily book sales using Multiple Linear Regression and interpret the outcomes using SHAP value analysis.

Given Dataset:

Footfall (x_1)	Promotions (1/0) (x_2)	Sales (y)
100	1	1500
80	0	1000
120	1	1700
90	0	1100
60	1	900

Tasks:

1. **Perform Multiple Linear Regression Analysis**
 - Use Footfall and Promotions as independent variables
 - Use Sales as the dependent variable
2. **Calculate the Baseline Value**
 - Compute the mean of all sales values
3. **Calculate SHAP Values**
 - Calculate SHAP value
 - Distribute SHAP contributions between Footfall and Promotions based on model coefficients
4. **Compute Final Prediction for Each Record**
 - Use the regression equation
 - Verify: Prediction = Baseline + SHAP (Footfall) + SHAP (Promotions)
5. **Interpret the Results**
 - For each entry, explain how Footfall and Promotions influenced the predicted sales
 - Compare predicted vs actual values
 - Indicate any overprediction or underprediction and suggest potential reasons

Q. No.	Question	Expected Time to complete
3	Regression with Diabetes Dataset	

Objective:

Understand how patient features influence disease progression using Multiple Linear Regression and SHAP value analysis.

Tasks

1. *Perform Multiple Linear Regression Analysis*

- Use all available features from the Diabetes dataset as independent variables.
 - Fit a Multiple Linear Regression model to predict disease progression.
2. *Calculate the Baseline Value*
- Compute the **mean** of the target variable (disease progression scores) from the training data.
 - This will serve as the **baseline prediction**.
3. *Calculate SHAP Values*
- Apply SHAP to compute **feature contributions** to each prediction.
 - Use model coefficients to proportionally attribute the difference from the baseline to each feature.
4. *Compute Final Prediction for Each Record*
- For every test record, verify that:
Prediction = Baseline + SHAP(Feature₁) + SHAP(Feature₂) + ... + SHAP(Feature_n)
5. *Interpret the Results*
- For each patient record:
 - Explain how each feature contributed to the predicted disease progression.
 - Compare the **predicted value** vs the **actual observed value**.
 - Comment on whether the model **overpredicted or underpredicted** and **why**, based on SHAP values.

Q. No.	Question	Expected Time to complete
4	Regression with Student Performance Dataset	

Objective:

Investigate how student background and behavior influence final exam scores using Multiple Linear Regression and SHAP value analysis.

Tasks

1. *Perform Multiple Linear Regression Analysis*

- Use all relevant student attributes (e.g., study time, parental education, absences, etc.) as independent variables.
- Fit a regression model to predict the **final exam score**.

2. *Calculate the Baseline Value*

- Compute the **mean of the final exam scores** from the training set.
- This serves as the **baseline prediction** (expected value).

3. *Calculate SHAP Values*

- Use SHAP to compute the contribution of each student attribute to the final exam score prediction.
- Distribute the prediction deviation from the baseline among the features.

4. *Compute Final Prediction for Each Record*

- For each student record, confirm:
Predicted Score = Baseline + SHAP(Feature₁) + SHAP(Feature₂) + ... + SHAP(Feature_n)

5. *Interpret the Results*

- For every prediction:
 - Explain how different features (e.g., study time, failures, health) impacted the exam score.
 - Compare predicted score to actual score.
 - Comment on overprediction or underprediction and possible reasons behind it.