| 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** | 23CSBTB38 |

**Assignment Number:** 01

| Q. No. | Question | Expected Time to complete |
|---|---|---|
| 1 | Edu Spark – Educational YouTube Channel | |

**Context:**

Edu Spark uploads new videos to increase total weekly views.

| Videos Uploaded (x) | Weekly Views (y) |
|---|---|
| 1 | 500 |
| 2 | 750 |
| 3 | 950 |
| 1 | 550 |
| 2 | 800 |

**Objective:**

Analyze the effect of video uploads on weekly views for EduSpark by performing Linear Regression and interpreting SHAP values.

Requirements:

1. **Perform Linear Regression Analysis**
   o Use the given dataset where:
      ▪ **Independent Variable (x):** Videos Uploaded
      ▪ **Dependent Variable (y):** Weekly Views
2. **Calculate the Baseline Value**
   o Compute the **mean of all weekly views (y values)**.
3. **Calculate SHAP Values**
   o For each record, calculate the difference between the **predicted value** and the **baseline**.
   o This difference is the **SHAP value**, attributed to the number of videos uploaded.
4. **Compute Final Prediction**
   o Use the **linear regression model** to calculate predicted weekly views for each video count.
   o Confirm that:
      $$\text{Final Prediction} = \text{Baseline} + \text{SHAP Value}$$
5. **Interpret the Results**
   o Explain how the number of videos uploaded influenced each predicted view count.
   o Compare the predicted value to the actual value for each row.
   o 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:
    o Accuracy of the model
    o Trend analysis
    o SHAP interpretation insights

| Q. No. | Question | Expected Time to complete |
|---|---|---|
| 2 | **FreshBasket – Grocery App Usage Retention Prediction using Multiple Linear Regression and SHAP Analysis** | |

**Objective:**
Evaluate how the number of push notifications and average delivery time influence user retention using Multiple Linear Regression, and explain the results through SHAP value interpretation.

**Given Dataset:**

| Notifications $(x_1)$ | Avg Delivery Time (min) $(x_2)$ | Retention (%) $(y)$ |
|---|---|---|
| 5 | 30 | 75 |
| 7 | 25 | 85 |
| 4 | 35 | 70 |
| 6 | 20 | 90 |
| 3 | 40 | 65 |

**Tasks:**

1. **Perform Multiple Linear Regression Analysis**
    o Use Notifications and Average Delivery Time as independent variables
    o Use Retention (%) as the dependent variable
2. **Calculate the Baseline Value**
    o Compute the mean of all retention values
3. **Calculate SHAP Values**
    o Calculate SHAP Value
    o Distribute SHAP contributions between Notifications and Delivery Time based on model coefficients
4. **Compute Final Prediction for Each Record**
    o Use the regression equation
    o Verify: Prediction = Baseline + SHAP (Notifications) + SHAP (Delivery Time)
5. **Interpret the Results**
    o For each entry, explain how notifications and delivery time influenced the prediction
    o Compare predicted vs actual retention
    o State whether the model overpredicted or underpredicted and suggest why

| 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(Feature$_1$) + SHAP(Feature$_2$) + ... + 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(Feature$_1$) + SHAP(Feature$_2$) + ... + 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.