

## PRE-LAB QUESTIONS

SUDHIKSHA 24BAD117

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

Experiment - 2  
Regression and Optimization (Linear and Logistic)

- PRE LAB QUESTIONS:

- 1) What is regression analysis, and how is it used to predict continuous values such as ocean temperature in linear regression?  
  
Ans. Regression analysis is a statistical technique used to model the relationship between a dependent variable and one or more independent variables. In linear regression, this relationship is represented using a straight-line equation, which helps predict continuous values such as ocean temperature.
- 2) How does logistic regression differ from linear regression in terms of:
  - a) type of output produced
  - b) logistic regression?

- a) Linear regression produces continuous numerical values. Logistic regression produces probability values between 0 and 1, which are later converted into class labels.
- b) Linear regression is used for regression problems. Logistic regression is used for classification problems. eg. disease detection, spam classification.

3) In a regression equation

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

what does each regression coefficient  $\beta_i$  indicate about the corresponding feature?

Ans. Each regression coefficient represents the change in the output variable for a one-unit change in corresponding feature, assuming all other features remain constant. The sign indicates the direction of influence, while the magnitude indicates the strength of impact.

4) Why is feature scaling (standardization or normalization) necessary before applying:

- a) gradient - descent based linear regression
- b) logistic regression?

a) Gradient - descent - based linear regression feature scaling ensures faster and more stable convergence of gradient descent by preventing features with large values from dominating the optimization.

b) Feature scaling improves numerical stability and speeds up convergence, especially when gradient based optimization techniques are used.



5) What is a cost function, and how is it used to evaluate the performance of regression model?

Ans. A cost function measures the difference between the predicted values and actual values. It provides a single numerical value that represents model error, and lower cost values indicate better model performance.

6) What is optimization in machine learning, and what does it mean to minimize a loss function?

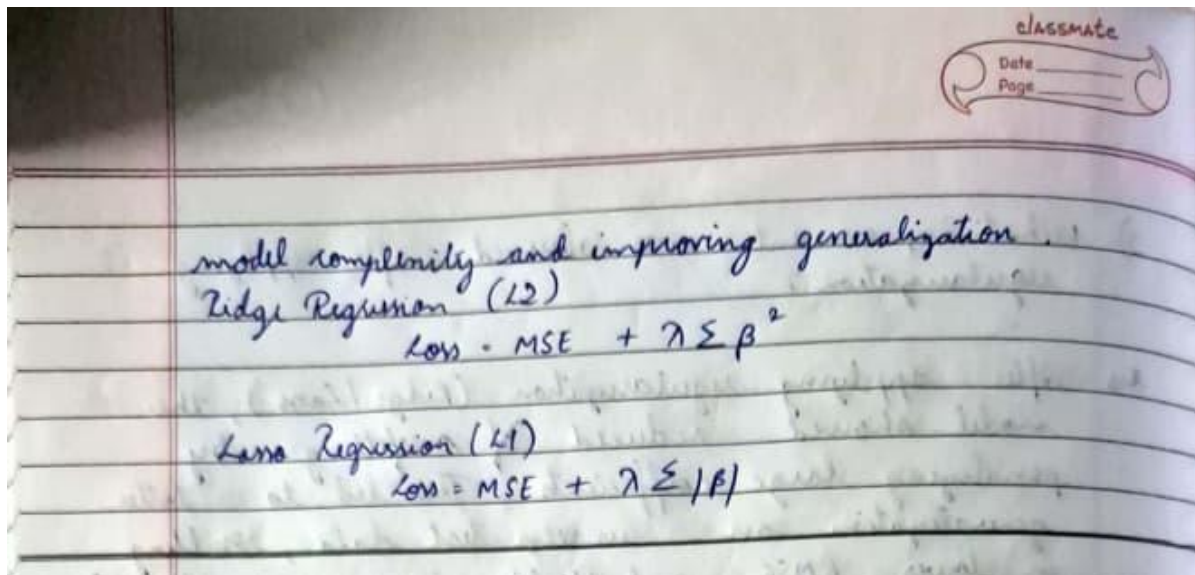
Ans. Optimization is the process of adjusting model parameters values that result in smaller improved performance. Minimizing a loss function means minimizing prediction error, leading to more accurate model.

7) Why is gradient descent preferred over solving equations directly when training regression models on large datasets?

Ans. Gradient descent is preferred because solving equations directly becomes computationally expensive for large datasets. Gradient descent is more efficient, scalable, and can handle high-dimensional data effectively.

8) What is overfitting in regression models, how does regularization help reduce it?

Ans. Overfitting occurs when a model learns noise in the training data and performs poorly on new data. Regularization techniques such as Ridge ( $L_2$ ) and Lasso ( $L_1$ ) add penalty terms to the loss function, reducing



## IN-LAB EXERCISE

### OBJECTIVE:

To implement and optimize **linear regression** and **logistic regression** models using real-world datasets and analyze their performance.

### SCENARIO 1:

Predict ocean **water temperature** using environmental and depth-related features.

### Dataset (Kaggle - Public):

<https://www.kaggle.com/datasets/sohier/calcofi>

### Target Variable:

- Water Temperature (T\_degC)

### Sample Input Features

- Depth (m)
- Salinity
- Oxygen
- Latitude
- Longitude

### IN-LAB TASKS

- Import necessary Python libraries (NumPy, Pandas, Matplotlib, Seaborn, Scikit-learn).
- Load the CalCOFI dataset into a Pandas DataFrame
- Select relevant numerical features and target variable.
- Handle missing values using mean/median imputation.
- Perform feature scaling using StandardScaler.
- Split the dataset into training and testing sets.
- Train a **Linear Regression** model using Scikit-learn.
- Predict water temperature for test data.
- Evaluate model performance using:
  - Mean Squared Error (MSE)
  - Root Mean Squared Error (RMSE)
  - $R^2$  Score
- Visualize:
  - Actual vs Predicted temperature
  - Residual errors

- Optimize model performance using:
  - Feature selection
  - Regularization (Ridge / Lasso)

### SCENARIO 2:

Classify whether LIC stock price will **increase (1)** or **decrease (0)** based on historical data.

#### Dataset (Kaggle – Public):

<https://www.kaggle.com/datasets/debashis74017/lic-stock-price-data>

#### Target Variable (Derived):

- Price Movement
  - 1 → Closing price > Opening price
  - 0 → Closing price ≤ Opening price

#### Input Features:

- Open
- High
- Low
- Volume

### IN-LAB TASKS (Logistic Regression)

- Import required Python libraries.
- Load LIC stock dataset into Pandas.
- Create a binary target variable (Price Movement).
- Handle missing values.
- Perform feature scaling.
- Split the dataset into training and testing sets.
- Train a **Logistic Regression** model.
- Predict stock movement for test data.
- Evaluate classification performance using:
  - Accuracy
  - Precision
  - Recall
  - F1-Score
  - Confusion Matrix
- Visualize:
  - ROC Curve
  - Feature importance
- Optimize model using:
  - Hyperparameter tuning (C, penalty)
  - Regularization

## CODE:

```
1 #SCENARIO 1: Predict Ocean Water Temperature
2 #Sudhiksha 24BAD117
3
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7
8 from sklearn.model_selection import train_test_split
9 from sklearn.preprocessing import StandardScaler
10 from sklearn.linear_model import LinearRegression, Ridge, Lasso
11 from sklearn.metrics import mean_squared_error, r2_score
12
13
14 df = pd.read_csv("bottle.csv", low_memory=False)
15 df.columns = df.columns.str.strip()
16 features = ['Depthm', 'Salnty', 'O2mL_L']
17 target = 'T_degC'
18
19 df = df[features + [target]]
20 df.fillna(df.mean(), inplace=True)
21 X = df[features]
22 y = df[target]
23
24
25 scaler = StandardScaler()
26 X_scaled = scaler.fit_transform(X)
27 X_train, X_test, y_train, y_test = train_test_split(
28     X_scaled, y, test_size=0.2, random_state=42
29 )
30
31
32 lr = LinearRegression()
33 lr.fit(X_train, y_train)
34
35 y_pred = lr.predict(X_test)
36 mse = mean_squared_error(y_test, y_pred)
37 rmse = np.sqrt(mse)
38 r2 = r2_score(y_test, y_pred)
39
40 print("Linear Regression Performance:")
41 print("MSE :", mse)
42 print("RMSE:", rmse)
43 print("R2 :", r2)
44
45 plt.figure(figsize=(6,5))
46 plt.scatter(y_test, y_pred, alpha=0.5)
47 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'r--')
48 plt.xlabel("Actual Temperature")
49 plt.ylabel("Predicted Temperature")
50 plt.title("Actual vs Predicted Water Temperature")
51 plt.show()
52
53 # Residual Plot
54 residuals = y_test - y_pred
55 plt.figure(figsize=(6,5))
56 plt.scatter(y_pred, residuals, alpha=0.5)
57 plt.axhline(0, color='r', linestyle='--')
58 plt.xlabel("Predicted Temperature")
59 plt.ylabel("Residual Error")
60 plt.title("Residual Plot")
61 plt.show()
62
63 # Ridge Regression
64 ridge = Ridge(alpha=1.0)
65 ridge.fit(X_train, y_train)
66 ridge_pred = ridge.predict(X_test)
67 print("Ridge Regression R2:", r2_score(y_test, ridge_pred))
68
```

```

1  SUDHIKSHA 248AD117
2  # SCENARIO 2 - MULTI-FILE VERSION
3
4  import pandas as pd
5  import numpy as np
6  import matplotlib.pyplot as plt
7
8  from sklearn.model_selection import train_test_split, GridSearchCV
9  from sklearn.preprocessing import StandardScaler
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.metrics import (
12     accuracy_score, precision_score, recall_score, f1_score,
13     confusion_matrix, roc_curve, roc_auc_score
14 )
15
16 files = [
17     "LICI - minute ata.csv",
18     "LICI - 3 minute data.csv",
19     "LICI - 5 minute data.csv",
20     "LICI - 10 minute data.csv",
21     "LICI - Daily data.csv"
22 ]
23
24
25 def find_col(df, possible):
26     for name in possible:
27         if name in df.columns:
28             return name
29     return None
30
31
32 for file in files:
33
34     print("\n" + "="*50)
35     print("Processing file:", file)
36
37     df = pd.read_csv(file)
38     df.columns = df.columns.str.strip().str.lower()
39
40     open_col = find_col(df, ['open'])
41     high_col = find_col(df, ['high'])
42     low_col = find_col(df, ['low'])
43     close_col = find_col(df, ['close', 'adj_close', 'adj_close'])
44     volume_col = find_col(df, ['volume'])
45
46     if None in [open_col, high_col, low_col, close_col, volume_col]:
47         print("Required columns missing → Skipping file")
48         continue
49
50     df['price_movement'] = np.where(df[close_col] > df[open_col], 1, 0)
51     df = df[[open_col, high_col, low_col, volume_col, 'price_movement']].dropna()
52
53     X = df[[open_col, high_col, low_col, volume_col]]
54     y = df['price_movement']
55
56     scaler = StandardScaler()
57     X_scaled = scaler.fit_transform(X)
58
59     X_train, X_test, y_train, y_test = train_test_split(
60         X_scaled, y, test_size=0.2, random_state=42, stratify=y
61     )
62
63     model = LogisticRegression(max_iter=1000)
64     model.fit(X_train, y_train)

```



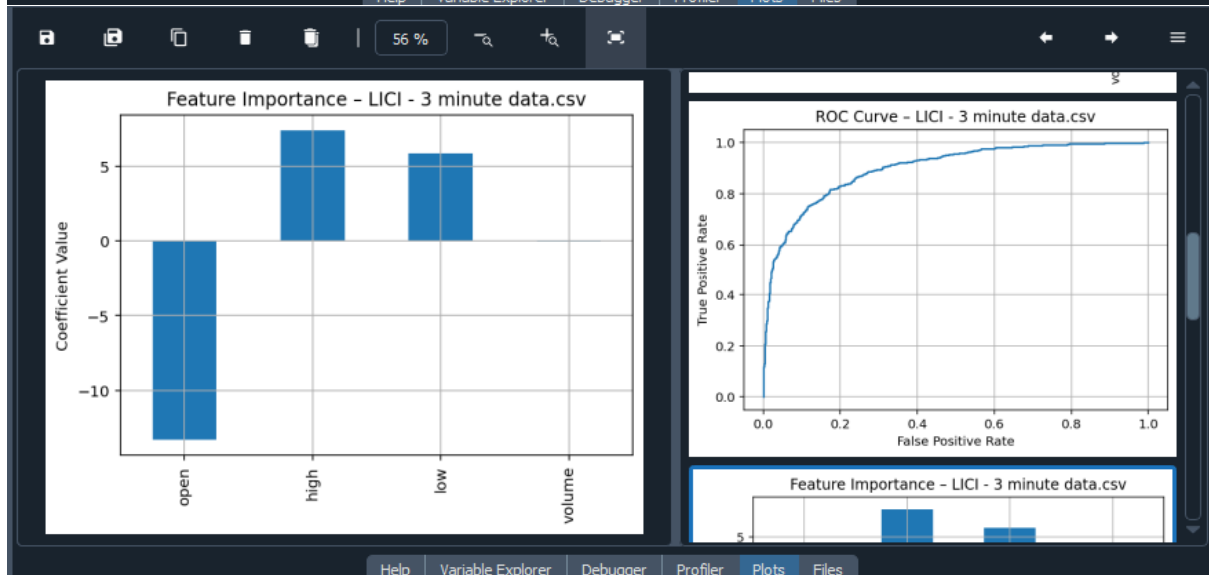
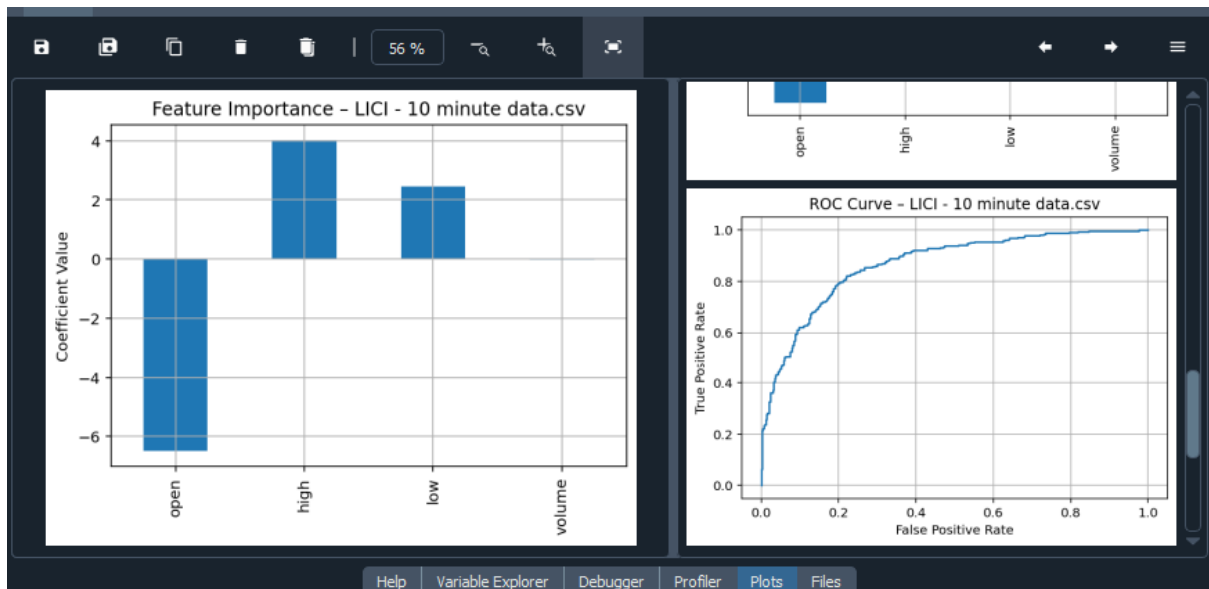
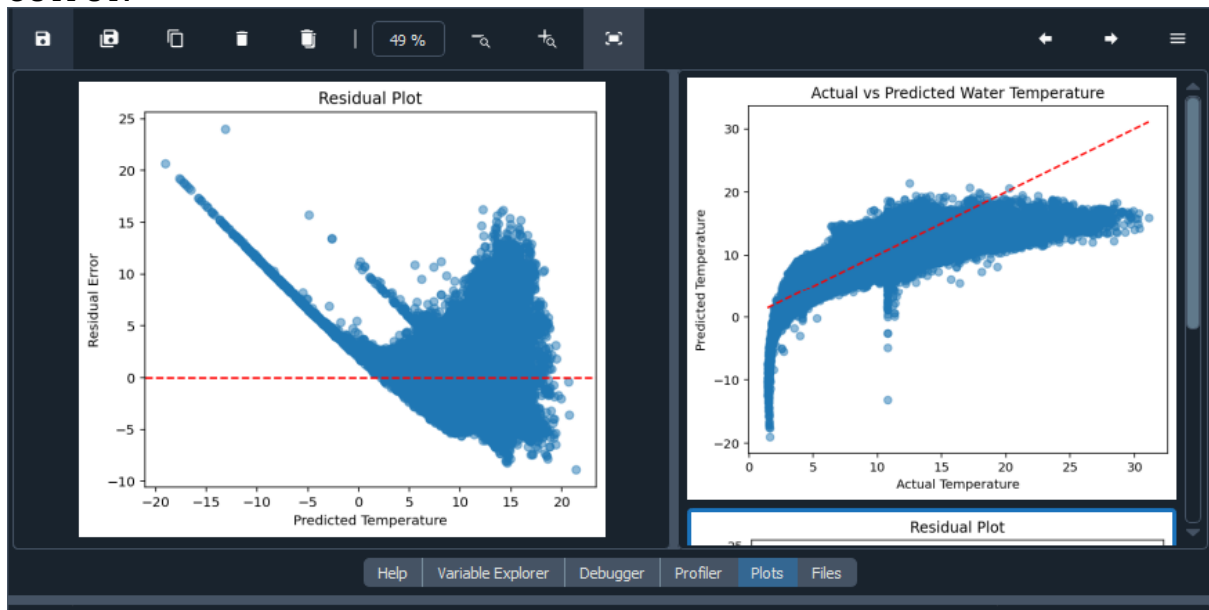
```

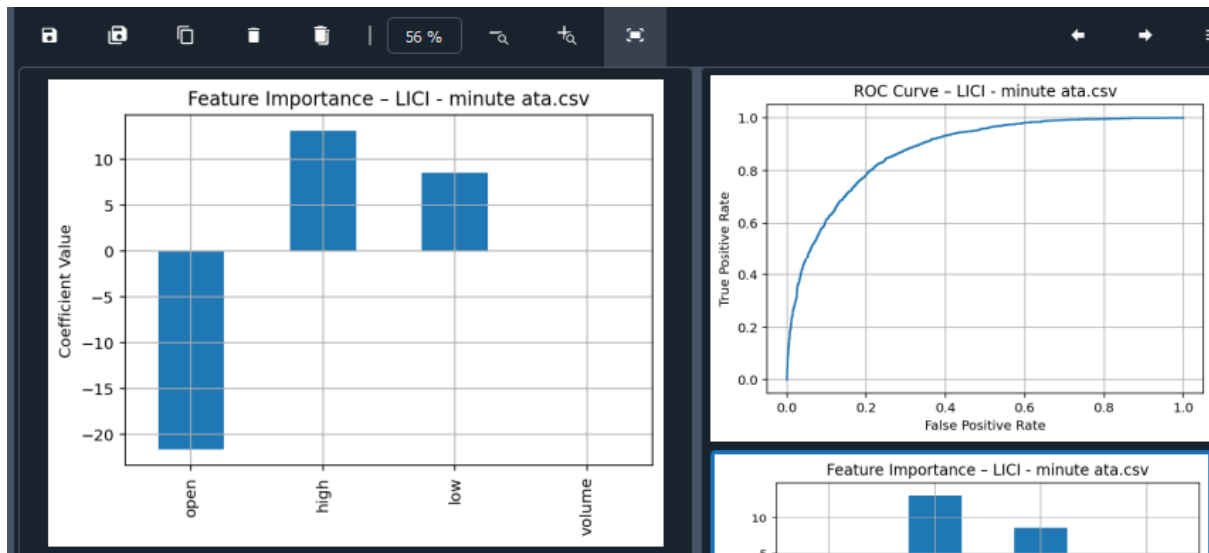
66     y_pred = model.predict(X_test)
67     y_prob = model.predict_proba(X_test)[: , 1]
68
69     acc = accuracy_score(y_test, y_pred)
70     prec = precision_score(y_test, y_pred)
71     rec = recall_score(y_test, y_pred)
72     f1 = f1_score(y_test, y_pred)
73     cm = confusion_matrix(y_test, y_pred)
74     auc = roc_auc_score(y_test, y_prob)
75
76     print("\n--- Performance Metrics ---")
77     print("Accuracy :", acc)
78     print("Precision:", prec)
79     print("Recall   :", rec)
80     print("F1-Score :", f1)
81     print("ROC-AUC  :", auc)
82     print("Confusion Matrix:\n", cm)
83
84
85     # ROC Curve
86     fpr, tpr, _ = roc_curve(y_test, y_prob)
87
88     plt.figure()
89     plt.plot(fpr, tpr)
90     plt.xlabel("False Positive Rate")
91     plt.ylabel("True Positive Rate")
92     plt.title(f"ROC Curve - {file}")
93     plt.grid()
94     plt.show()
95
96     # Feature Importance
97     importance = pd.Series(model.coef_[0], index=X.columns)
98
99     plt.figure()
100
101     # Feature Importance
102     importance = pd.Series(model.coef_[0], index=X.columns)
103
104     plt.figure()
105     importance.plot(kind="bar")
106     plt.title(f"Feature Importance - {file}")
107     plt.ylabel("Coefficient Value")
108     plt.grid()
109     plt.show()
110
111     # Hyperparameter tuning
112     param_grid = {
113         "C": [0.01, 0.1, 1, 10],
114         "penalty": ["l2"],
115         "solver": ["lbfgs"]
116     }
117
118     grid = GridSearchCV(
119         LogisticRegression(max_iter=1000),
120         param_grid,
121         cv=5,
122         scoring="accuracy"
123     )
124
125     grid.fit(X_train, y_train)
126
127     print("\nBest Hyperparameters:", grid.best_params_)
128     print("Tuned Accuracy:", accuracy_score(y_test, grid.best_estimator_.predict(X_test)))
129
130     print("\nAll files processed successfully.")

```

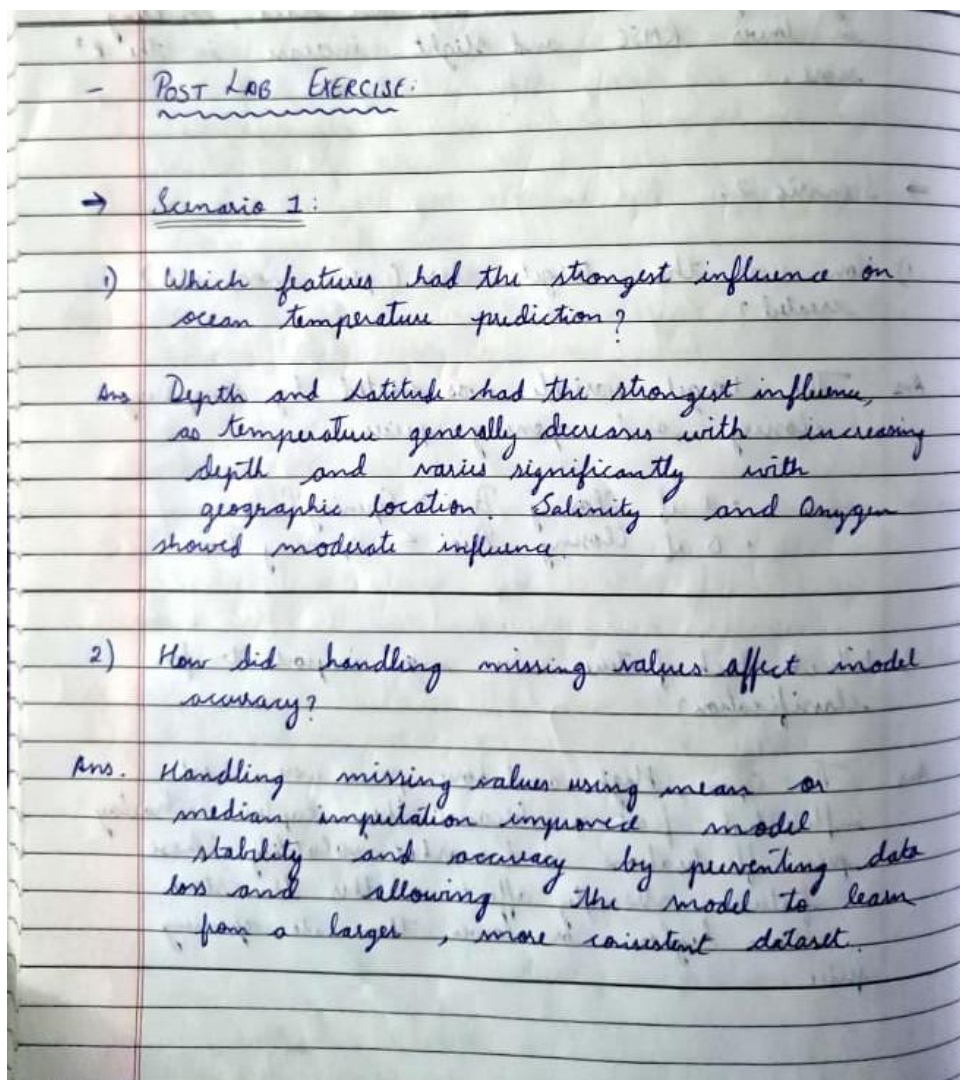


## OUTPUT:





### POST-LAB QUESTIONS (PROVIDE BRIEF ANSWERS TO THE FOLLOWING QUESTIONS)



3) What improvement was observed after applying regularization?

Ans After applying regularization (Ridge/Lasso), the model showed reduced  $\Delta$  overfitting by penalizing large coefficients. This led to better generalization on unseen test data, resulting in lower RMSE and slight increase in the  $R^2$  score.

⇒ Scenario 2:

1) How was the target class (price movement) created?

Ans The target variable was created by comparing closing and opening prices:

- 1 if Closing Price > Opening Price
- 0 if Closing Price  $\leq$  Opening Price

2) Which stock features were most important for classification?

Ans The Open, High and Low prices were most influential features because they capture intraday price fluctuations and market volatility. These features strongly affect whether the stock closes higher or lower than its opening price.



3) How did feature scaling impact convergence and accuracy?

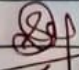
Ans. Feature scaling standardized all input variables to common scale, preventing features with larger numerical ranges from dominating the learning process. It improved speed of Logistic Regression algorithm, resulted in more consistent coefficient estimation.

4) What changes were observed after hyperparameter tuning?

Ans. Hyperparameter tuning, particularly adjusting the regularization strength and penalty type, improved the model's balance between bias and variance. The tuned model achieved higher accuracy and better precision.

#### ASSESSMENT

Description	Max Marks	Marks Awarded
Pre Lab Exercise	5	5
In Lab Exercise	10	10
Post Lab Exercise	5	5
ViVa	10	6
Total	30	26

Faculty Signature:   
6/26