

EXPERIMENT 2

Aim: Implement Multi Regression, Lasso, and Ridge Regression on real-world datasets

1. Dataset Source: Custom dataset : [Link](#)

2. Dataset Description

The dataset used in this experiment is obtained from Fitbit activity records. It contains timestamp-based physical activity data.

Independent Variables (Features)

- Steps
- Distance
- Active Minutes

Dependent Variable (Target)

- Calories burned

Each record corresponds to a specific timestamp and contains numerical values. The dataset is continuous in nature and suitable for regression analysis, as calories burned depend on physical activity metrics.

3. Mathematical Formulation of the Algorithm

Multiple Linear Regression

The prediction equation is:

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

Where:

x_1 = Steps

x_2 = Distance

x_3 = Active Minutes

The objective is to minimize Mean Squared Error:

$$MSE = (1/n) \sum (y - \hat{y})^2$$

Ridge Regression (L2 Regularization)

Ridge modifies the cost function as:

$$\text{Cost} = \text{MSE} + \lambda \sum \beta^2$$

It reduces the magnitude of coefficients to prevent overfitting.

Lasso Regression (L1 Regularization)

Lasso modifies the cost function as:

$$\text{Cost} = \text{MSE} + \lambda \sum |\beta|$$

It shrinks coefficients and may reduce some to zero.

4. Algorithm Limitations

- Assumes a linear relationship between features and target.
- Sensitive to outliers.
- Not suitable for complex non-linear patterns.
- Regularization may shrink important coefficients if not properly chosen.

5. Methodology / Workflow

1. Loaded Fitbit CSV files.
2. Selected timestamp and numeric columns.
3. Merged datasets using timestamp.
4. Removed missing values.
5. Selected features (Steps, Distance, Active Minutes) and target (Calories).
6. Split data into training and testing sets (80–20).
7. Trained Linear, Ridge, and Lasso regression models.
8. Evaluated models using MSE and R² score.

Workflow:

Data Loading → Data Cleaning → Merging → Feature Selection → Train-Test Split → Model Training → Evaluation

6. Performance Analysis

Model performance was evaluated using:

- Mean Squared Error (MSE)
- R² Score

All three models produced similar results due to strong linear correlation between activity metrics and calories burned. Ridge and Lasso slightly reduced coefficient values compared to Linear Regression.

7. Hyperparameter Tuning

In this experiment, a fixed alpha value was used for Ridge and Lasso regression. No multiple hyperparameter values were tested. The models were trained using the default regularization settings to compare their behavior with standard Linear Regression.

Code:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression, Ridge, Lasso
5 from sklearn.metrics import mean_squared_error, r2_score
6
7 def load_and_clean(file_path, new_column_name):
8     df = pd.read_csv(file_path)
9
10
11    if "timestamp" not in df.columns:
12        raise Exception(f"'timestamp' column not found in {file_path}")
13
14    numeric_cols = df.select_dtypes(include=['int64', 'float64']).columns
15
16    if len(numeric_cols) == 0:
17        raise Exception(f"No numeric column found in {file_path}")
18
19    value_column = numeric_cols[0]
20
21    df = df[["timestamp", value_column]]
22    df = df.rename(columns={value_column: new_column_name})
23
24    df["timestamp"] = pd.to_datetime(df["timestamp"])
25
26    return df
27
28 steps = load_and_clean("steps.csv", "steps")
29 distance = load_and_clean("distance.csv", "distance")
30 active = load_and_clean("active_minutes.csv", "active_minutes")
31 calories = load_and_clean("calories.csv", "calories")
32
33
34 data = steps.merge(distance, on="timestamp") \
35         .merge(active, on="timestamp") \
36         .merge(calories, on="timestamp")
37
38 data = data.dropna()
```

```
40     print("Merged Data Sample:")
41     print(data.head())
42
43
44     X = data[["steps", "distance", "active_minutes"]]
45     y = data["calories"]
46
47
48     X_train, X_test, y_train, y_test = train_test_split(
49         X, y, test_size=0.2, random_state=42
50     )
51
52
53     linear_model = LinearRegression()
54     ridge_model = Ridge(alpha=1.0)
55     lasso_model = Lasso(alpha=0.1)
56
57     linear_model.fit(X_train, y_train)
58     ridge_model.fit(X_train, y_train)
59     lasso_model.fit(X_train, y_train)
60
61
62     y_pred_linear = linear_model.predict(X_test)
63     y_pred_ridge = ridge_model.predict(X_test)
64     y_pred_lasso = lasso_model.predict(X_test)
65
66
67     def evaluate(name, y_test, y_pred):
68         print(f"\n{name}")
69         print("MSE:", mean_squared_error(y_test, y_pred))
70         print("R2 Score:", r2_score(y_test, y_pred))
71
72     evaluate("Linear Regression", y_test, y_pred_linear)
73     evaluate("Ridge Regression", y_test, y_pred_ridge)
74     evaluate("Lasso Regression", y_test, y_pred_lasso)
```

```

print("\nCoefficient Comparison")
print("Features:", X.columns.tolist())
print("Linear:", linear_model.coef_)
print("Ridge :", ridge_model.coef_)
print("Lasso :", lasso_model.coef_)

plt.figure(figsize=(8,6))
plt.scatter(y_test, y_pred_linear, label="Linear")
plt.scatter(y_test, y_pred_ridge, label="Ridge")
plt.scatter(y_test, y_pred_lasso, label="Lasso")

plt.xlabel("Actual Calories")
plt.ylabel("Predicted Calories")
plt.title("Regression Model Comparison")
plt.legend()
plt.show()

```

Output:

