

# Contents

<b>1</b>	<b>Multivariate Linear Regression</b>	<b>3</b>
1.1	Introduction . . . . .	3
1.1.1	Why it matters . . . . .	3
1.1.2	Scope . . . . .	3
1.2	Deep Explanation . . . . .	4
1.2.1	Matrix Notation . . . . .	4
1.2.2	Estimating the Coefficients (Ordinary Least Squares - OLS) . . .	4
1.2.3	Assumptions of OLS Linear Regression . . . . .	5
1.2.4	Interpretation of Coefficients . . . . .	5
1.2.5	Model Evaluation . . . . .	5
1.3	Example: Python Implementation . . . . .	5
1.4	Summary / Key Takeaways . . . . .	6



# Chapter 1

## Multivariate Linear Regression

### 1.1 Introduction

Multivariate linear regression is a statistical technique used to model the linear relationship between a dependent variable and multiple independent variables. Unlike simple linear regression, which involves only one independent variable, multivariate linear regression allows for the analysis of how changes in several predictors simultaneously affect the outcome. It's a cornerstone in statistical modeling and machine learning, providing a powerful tool for prediction, inference, and understanding complex relationships in data.

The core idea is to find a linear equation that best describes how the dependent variable can be predicted from a set of independent variables. This “best fit” is typically determined by minimizing the sum of squared differences between the observed and predicted values, a method known as Ordinary Least Squares (OLS).

#### 1.1.1 Why it matters

- **Prediction:** It allows for forecasting the value of a dependent variable based on multiple influencing factors.
- **Inference:** It helps in understanding which independent variables have a significant impact on the dependent variable and the magnitude and direction of those effects.
- **Control for Confounding:** By including multiple variables, it can account for the effects of other factors, leading to a more accurate understanding of individual predictor relationships.

#### 1.1.2 Scope

Multivariate linear regression is widely applied across various fields, including economics (predicting GDP based on multiple indicators), finance (stock price prediction), social sciences (factors influencing educational attainment), engineering (predicting material properties), and healthcare (predicting disease risk).

## 1.2 Deep Explanation

At its heart, multivariate linear regression seeks to establish a linear equation of the form:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \epsilon \quad (1.1)$$

Where:

- $Y$ : The dependent (response) variable.
- $X_1, X_2, \dots, X_p$ : The  $p$  independent (predictor) variables.
- $\beta_0$ : The intercept, representing the expected value of  $Y$  when all  $X_i = 0$ .
- $\beta_1, \beta_2, \dots, \beta_p$ : The coefficients for each independent variable, representing the change in  $Y$  for a one-unit change in the corresponding  $X_i$ , holding all other variables constant.
- $\epsilon$ : The error term, representing the portion of  $Y$  that cannot be explained by the linear relationship with the predictors.

### 1.2.1 Matrix Notation

For computational efficiency, especially with large datasets, the equation is often expressed in matrix form:

$$Y = X\beta + \epsilon \quad (1.2)$$

Where:

- $Y$ : A vector of observed dependent variable values ( $n \times 1$ ).
- $X$ : The design matrix ( $n \times (p + 1)$ ), containing a column of ones for the intercept and columns for each independent variable.
- $\beta$ : A vector of coefficients ( $(p + 1) \times 1$ ), including  $\beta_0, \beta_1, \dots, \beta_p$ .
- $\epsilon$ : A vector of error terms ( $n \times 1$ ).

### 1.2.2 Estimating the Coefficients (Ordinary Least Squares - OLS)

The goal of OLS is to find  $\hat{\beta}$  that minimizes the sum of squared residuals:

$$SSR = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = (Y - X\hat{\beta})^T (Y - X\hat{\beta}) \quad (1.3)$$

By setting the derivative to zero, we get the normal equation:

$$(X^T X)\hat{\beta} = X^T Y \quad (1.4)$$

Solving for  $\hat{\beta}$ :

$$\hat{\beta} = (X^T X)^{-1} X^T Y \quad (1.5)$$

provided that  $(X^T X)$  is invertible.

### 1.2.3 Assumptions of OLS Linear Regression

- **Linearity:** The relationship between  $Y$  and  $X$  is linear.
- **Independence of Errors:** The error terms are independent.
- **Homoscedasticity:** The variance of the errors is constant.
- **Normality of Errors:** Errors are normally distributed.
- **No Multicollinearity:** Independent variables are not highly correlated.
- **No Endogeneity:** Predictors are uncorrelated with the error term.

### 1.2.4 Interpretation of Coefficients

Each  $\beta_j$  represents the expected change in  $Y$  for a one-unit increase in  $X_j$ , holding all other variables constant.

### 1.2.5 Model Evaluation

- **R-squared ( $R^2$ ):** Proportion of variance in  $Y$  explained by  $X$ :

$$R^2 = 1 - \frac{SSR}{SST} = 1 - \frac{\sum(Y_i - \hat{Y}_i)^2}{\sum(Y_i - \bar{Y})^2} \quad (1.6)$$

- **Adjusted R-squared:** Adjusts for the number of predictors.
- **F-statistic:** Tests overall model significance.
- **t-tests:** Evaluate the significance of individual coefficients.
- **Residual Analysis:** Used to verify OLS assumptions.

## 1.3 Example: Python Implementation

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.linear_model import LinearRegression
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import mean_squared_error, r2_score
6
7 # Generate synthetic data
8 np.random.seed(42)
9 n_samples = 100
10 hours_studied = np.random.rand(n_samples) * 10
11 attendance = np.random.rand(n_samples) * 5
12 previous_score = np.random.rand(n_samples) * 100
13
14 # True coefficients
15 beta_0 = 30
16 beta_hours = 5
17 beta_attendance = -3

```

```

18 beta_prev = 0.5
19
20 errors = np.random.randn(n_samples) * 5
21 scores = (beta_0 +
22           beta_hours * hours_studied +
23           beta_attendance * attendance +
24           beta_prev * previous_score +
25           errors)
26
27 data = pd.DataFrame({
28     'Hours_Studied': hours_studied,
29     'Attendance': attendance,
30     'Previous_Score': previous_score,
31     'Score': scores
32 })
33
34 X = data[['Hours_Studied', 'Attendance', 'Previous_Score']]
35 y = data['Score']
36
37 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
38     =0.2, random_state=42)
39 model = LinearRegression()
40 model.fit(X_train, y_train)
41 y_pred = model.predict(X_test)
42
43 print("Model Coefficients:")
44 for feature, coef in zip(X.columns, model.coef_):
45     print(f"{feature}: {coef:.2f}")
46
47 print(f"Intercept: {model.intercept_:.2f}")
48
49 print(f"\nMSE: {mean_squared_error(y_test, y_pred):.2f}")
50 print(f"R-squared: {r2_score(y_test, y_pred):.2f}")

```

Listing 1.1: Python Implementation for Multivariate Linear Regression

```

1 Model Coefficients:
2 Hours_Studied: 4.88
3 Attendance: -3.07
4 Previous_Score: 0.50
5 Intercept: 30.60
6
7 MSE: 26.68
8 R-squared: 0.88

```

Listing 1.2: Sample Output

## 1.4 Summary / Key Takeaways

- **Definition:** Models linear relationships with multiple predictors.
- **Equation:**  $Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \epsilon$
- **Matrix Form:**  $Y = X\beta + \epsilon$
- **Estimator:**  $\hat{\beta} = (X^T X)^{-1} X^T Y$
- **Evaluation:** Use  $R^2$ , Adjusted  $R^2$ , F-statistic, t-tests.