# Notes on Linear Regression: A Quick Reference Guide

## 1 Introduction

Linear regression is the statistical technique and machine learning algorithm that models the relationship between a dependent variable and one or more independent or explanatory variables [5]. The method attempts to fit a linear equation to the observed data, thereby facilitating an understanding of the association between these variables [1]. In practice, it is important to determine the presence of a linear relationship through exploratory data analysis before applying linear regression. Data visualization methods, such as a scatterplot can serve as a preliminary tool in evaluating the potential for a linear model to evaluate the data relationship adequately.

### 1.1 Linear Regression Model

The linear regression model assumes that any dependent variable $(Y)$ can be expressed as a linear function of an independent variable $(X)$, plus an error term $(\varepsilon)$ that accounts for discrepancies. The model is defined by the equation:

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

where $\beta_0$ is the intercept, $\beta_1$ the slope, and $\varepsilon_i$ the error term for each observation $i$. The linear regression model assumes a linear relationship between the variables, whereby the dependent variable can be predicted within the scope of the independent variable's influence [1].

### 1.2 Error Term

The error term $\varepsilon$ plays a crucial role in linear regression. It captures the discrepancy between the actual observed values $(Y)$ and the predicted values $(\hat{Y})$. These predicted values are calculated using the model's equation: $\hat{Y} = \beta_0 + \beta_1 X$. The error term $\varepsilon$ is derived from the difference between the observed values $(Y)$ and those predicted by the linear model $(\hat{Y})$.

$$\varepsilon_i = Y_i - \hat{Y}_i$$

Here, $\varepsilon_i$ represents the error term for the $i^{th}$ observation, $Y_i$ is the actual observed value, and $\hat{Y}_i$ is the predicted value for the $i^{th}$ observation, calculated from the linear regression model. The predicted value $\hat{Y}_i$ is calculated using the equation:

$$\hat{Y}_i = \beta_0 + \beta_1 X_i$$

Thus, inserting this into the equation for $\varepsilon_i$, we get:

$$\varepsilon_i = Y_i - (\beta_0 + \beta_1 X_i)$$

This difference captures the randomness and the unexplained variability in the dependent variable, Y, that the independent variable, X, does not account for [3]. It acknowledges the complexity of real-world data where other unknown factors affect $Y$. Including the error term helps in estimating the accuracy of the model's predictions and provides insights into the potential improvement of the model by considering additional or different independent variables. Essentially, the error term should be seen as reflecting the total influence of factors not included in the model. This interpretation hinges on three critical conditions: first, the influence of these omitted factors does not correlate with the influence of the included variables; second, their influence is consistent across different observations; and third, on average, this influence is neutral, meaning it does not systematically add to or subtract from the dependent variable. This refined understanding highlights the importance of acknowledging the impact of unobserved variables on the outcomes being studied [4].

## 2 Importance in Machine Learning

In machine learning, linear regression is a fundamental algorithm valued for its ability to manage high-dimensional data and predict continuous variables effectively. Linear regression is utilized extensively for predictive tasks where the objective is to forecast outcomes based on observed trends. In financial analysis, it is applied to project market changes and asset prices, providing a quantitative framework for investment strategies [11]. The algorithm's capacity to model continuous data makes it equally relevant in the field of econometrics for predicting economic indicators [8].

Linear regression can be employed within healthcare analytics, from treatment assesment to facilitating resource management and operational efficiency. For example, it can be utilized to model and predict patient length of stay post-surgery, serving as a key metric for hospital resource planning [13]. This predictive application allows for more informed decision-making in patient care logistics and resource allocation, thereby enhancing the operational efficacy of healthcare institutions.

Despite its utility, linear regression is characterized by an inherent assumption of linearity, which may not be suitable for all datasets, particularly those with complex, nonlinear interactions. However, the transparency and interpretability of linear models are particularly advantageous in scenarios where the explanation of the influence of predictor variables is as critical as predictive accuracy.

The enduring relevance of linear regression in machine learning is multifaceted. It serves as a transparent, interpretable model for prediction; it acts as a foundational algorithm for more sophisticated methods; and it remains a vital tool across numerous sectors for both predictive and explanatory purposes. Subsequent sections will expound upon the mathematical formalism of linear regression, examine its various forms, and illustrate its applications through an empirical case study.

## 3 Theoretical Background

The theoretical foundation of linear regression is grounded in the principles of statistical inference, where the relationship between variables is expressed linearly. The method's conceptual framework encompasses the estimation of model parameters using least squares or maximum likelihood estimation techniques. These estimators, derived from foundational statistical theories, are considered the best linear unbiased estimators (BLUE) under certain conditions[5]. Specifically, this holds true when the estimator is linear, meaning it is a linear function of a random variable, such as the dependent variable Y in the regression model. It means that among all the unbiased linear estimators (estimators whose expected value equals the true value of the parameter being estimated), the BLUE has the smallest variance. It is therefore, the most accurate or "best" estimator that doesn't systematically overestimate or underestimate the true value of the parameter it aims to estimate, while also having the least amount of error compared to any other unbiased linear estimator. This concept is central to ensuring that the linear regression model is as precise and reliable as possible in estimating relationships between variables.

The methodology of linear re hinges on key assumptions: linearity, homoscedasticity, error independence, and normally distributed errors, which are critical for ensuring the model's estimates are unbiased and consistent.

This framework underscores the necessity of accurate data representation and the theoretical challenges of measurement without error.

# 4 Simple Linear Regression

Simple linear regression is the foundational method of regression analysis, delineating a direct linear relationship between an independent variable (X) and a dependent variable (Y). It is mathematically articulated as $Y = \beta_0 + \beta_1 X + \varepsilon$, where $\beta_0$ represents the y-intercept, $\beta_1$ denotes the slope of the regression line, and $\varepsilon$ signifies the error term. This model is pivotal for introducing regression analysis, shedding light on the magnitude and direction of the relationship between variables. Primarily utilized as a predictive mechanism, it also facilitates the exploration of potential causality.

## 4.1 Equation: Formulation and Interpretation

The simple linear regression model is defined by the equation:

$$Y = \beta_0 + \beta_1 X + \varepsilon, \tag{1}$$

**Components**

- $\beta_0$ is the intercept of the regression line, indicating the expected value of $Y$ when $X$ is zero.

- $\beta_1$ is the slope of the regression line, representing the average change in $Y$ for a one-unit increase in $X$.

- $\varepsilon$ is the error term, accounting for the difference between the observed values and those predicted by the model.

**Key Points:**

- **Estimation of $\beta_0$ and $\beta_1$:** The primary objective is to find the values of $\beta_0$ and $\beta_1$ that best fit the observed data.

- **Role of $\beta_1$:** It quantifies the effect of the independent variable $X$ on the dependent variable $Y$.

- **Interpretation of $\beta_0$:** It provides a baseline value of $Y$ in the absence of any influence from $X$.

This mathematical formulation allows for a clear understanding of how each component contributes to the linear regression model and lays the foundation for further analysis and interpretation.

## 4.2 Type of Data: When to Use Simple Linear Regression

Simple Linear Regression (SLR) is best applied to data meeting the following criteria:

- **Single Predictor:** SLR analyzes the effect of a single independent variable on a dependent variable as shown in table 1.

- **Quantitative Variables:** Both the predictor and response variables must be quantitative and continuous.

- **Linear Relationship:** A linear relationship between the variables, ideally confirmed through a scatterplot as in illustrated in figure 1.

**Applications:** SLR is useful for prediction and understanding the influence of one variable on another, applicable in fields like economics for forecasting consumer spending based on income, or in biology for modeling growth rates.

| X | Y |
|---|---|
| 1 | 3.496714 |
| 2 | 4.861736 |
| 3 | 7.647689 |
| 4 | 10.523030 |
| 5 | 10.765847 |
| 6 | 12.765863 |
| 7 | 16.579213 |
| 8 | 17.767435 |
| 9 | 18.530526 |
| 10 | 21.542560 |

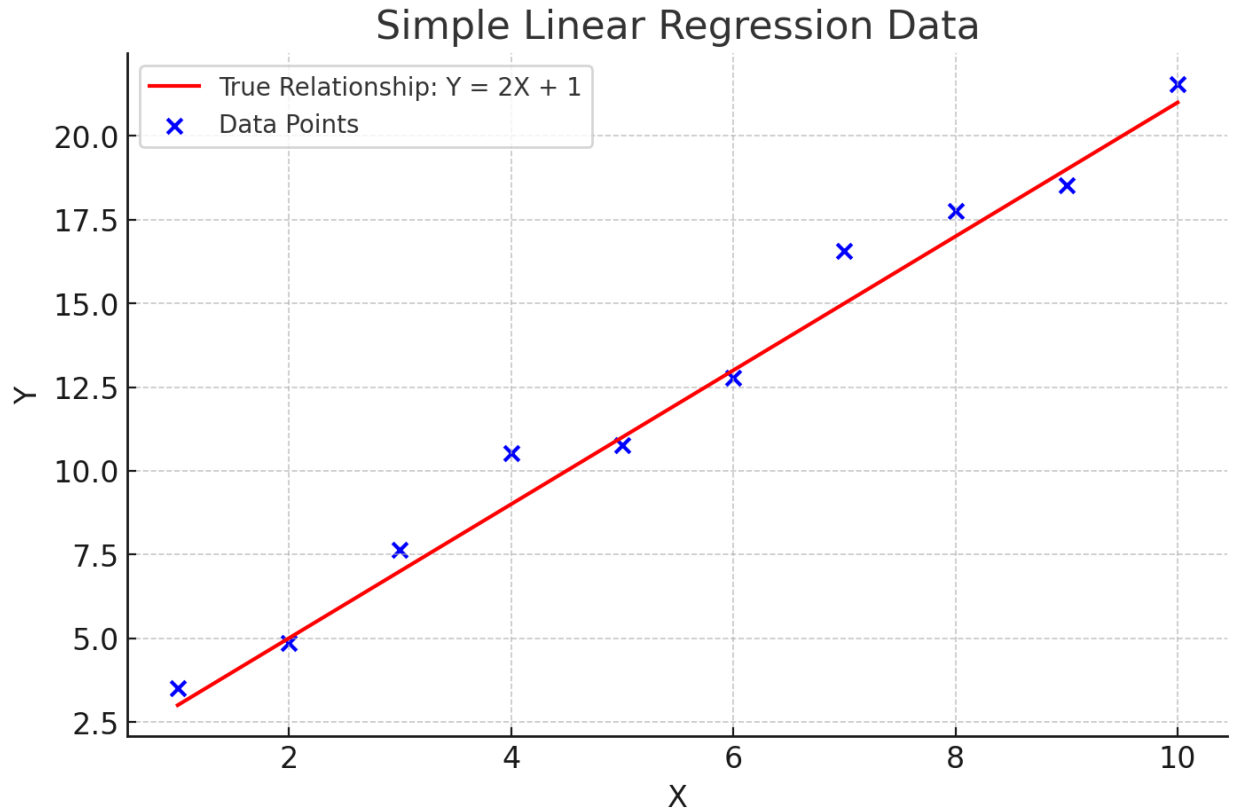Table 1: Data points exhibiting linear behavior



Figure 1: A graph illustrating a dataset with the 10 data points (fig 1 exhibiting linear behavior, alongside the true linear relationship $Y = 2X + 1$ for comparison. The blue dots represent the data points, which include some variability around the true linear relationship shown by the red line. This graph effectively demonstrates how simple linear regression can be used to estimate the underlying relationship between the independent variable $X$ and the dependent variable $Y$.

# 5 Multiple Linear Regression

Multiple Linear Regression (MLR) extends Simple Linear Regression to predict an outcome based on two or more independent variables. It models the linear relationship between a single response variable and multiple predictors, offering a comprehensive analysis of the effects of various factors.

## 5.1 Equation: Formulation and Interpretation

The general form of the Multiple Linear Regression model is represented as:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_n X_n + \varepsilon, \tag{2}$$

**Components:**

- $Y$: The dependent variable we aim to predict.
- $X_1, X_2, \ldots, X_n$: The independent variables or predictors.
- $\beta_0$: The y-intercept of the regression line.
- $\beta_1, \beta_2, \ldots, \beta_n$: The coefficients of the independent variables.
- $\varepsilon$: The error term.

**Key Points:**

- **Coefficient Estimation and Variable Influence:** The primary objective in MLR is to estimate the coefficients $(\beta_0, \beta_1, \ldots, \beta_n)$, where each $\beta_i$ (for $i > 0$) represents the expected change in $Y$ for a one-unit increase in $X_i$, assuming all other predictors remain constant. This process not only aims to best predict the dependent variable based on the multiple independent variables but also assesses the relative influence of each predictor on the outcome.

- **Intercept Interpretation ($\beta_0$):** Provides the expected value of $Y$ when all independent variables $(X_1, X_2, \ldots, X_n)$ are zero. It acts as a baseline from which the influence of the predictors can be measured.

- **Understanding the Error Term ($\varepsilon$):** Reflects the variation in $Y$ not explained by the predictors, capturing the residuals of the model.

This formulation seeks to minimize the sum of squared residuals, ensuring the best possible prediction of $Y$ based on the $X$ variables.

## 5.2 Type of Data: When to Use Multiple Linear Regression

Multiple Linear Regression is particularly useful in scenarios where a single outcome is influenced by multiple factors. Key considerations for using MLR include:

- **Predictive Analysis:** When the goal is to predict the value of a dependent variable based on several independent variables.
- **Effect of Variables:** To understand the effect of several independent variables on a dependent variable simultaneously.
- **Multifactor Influences:** In situations where the outcome is known to be influenced by multiple factors, MLR can quantify the strength and direction of these influences.
- **Real-World Applications:** MLR is widely applied in economics, finance, health sciences, and social sciences to model complex relationships between variables.

Choosing MLR requires careful consideration of the data's structure, the relationships among variables, and the specific research questions being addressed. It's essential to ensure that the assumptions of MLR are

met before proceeding with model fitting and interpretation. By carefully selecting variables and testing assumptions, researchers and practitioners can utilize MLR to uncover meaningful insights and make informed decisions.

# 6  Multivariate Linear Regression

Multivariate Linear Regression extends the concept of linear regression by modeling the relationships between multiple dependent variables and a set of predictor variables. Unlike simple linear regression, which predicts a single outcome variable from one or more predictors, a multivariate linear regression model assesses how a set of predictor variables (denoted as $Xs$) influence multiple outcome variables (denoted as $Ys$) [6]. This approach enables the simultaneous analysis of complex relationships where each predictor variable may affect multiple outcomes.

## 6.1  Equation: Vector Formulation

The general equation for a multivariate linear regression model, which encapsulates the relationships between multiple dependent variables and independent variables, is given by:

$$\begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_m \end{pmatrix} = \begin{pmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1n} \\ 1 & X_{21} & X_{22} & \cdots & X_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{m1} & X_{m2} & \cdots & X_{mn} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_m \end{pmatrix}$$

Here, $Y_1, Y_2, \ldots, Y_m$ represent the dependent variables for each of the $m$ outcomes. The matrix $X$ contains the values of $n$ predictor variables for each outcome, with each row representing an observation and each column a predictor variable. $\beta_0, \beta_1, \ldots, \beta_n$ are the regression coefficients that quantify the impact of each predictor on the outcome variables, and $\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_m$ are the error terms for each model, capturing the discrepancy between the observed and predicted values.

## 6.2  Type of Data: Scenarios for Multivariate Applications

This model is particularly useful in scenarios where the impact of predictors on multiple outcomes needs to be simultaneously assessed. Examples include socioeconomic studies where multiple aspects of well-being are predicted from demographic factors, or in clinical trials where the effects of treatments are measured across multiple health outcomes.

# 7  Multivariate Linear Regression

Multivariate Linear Regression (MvLR) advances the linear regression framework by modeling the influence of a set of predictor variables on multiple dependent variables simultaneously. While Simple Linear Regression (SLR) examines the relationship between a single predictor and a single outcome, and Multiple Linear Regression (MLR) explores the effect of multiple predictors on a single outcome, a multivariate linear regression model assesses how a set of predictor variables (denoted as $Xs$) influence multiple outcome variables (denoted as $Ys$) [6]. This distinction makes MvLR particularly suited for investigations into how various predictors impact several dependent variables, offering a comprehensive view of their interrelations.

## 7.1  Equation: Formulation and Interpretation

The Multivariate Linear Regression (MLR) model is defined through a matrix equation that represents the linear relationships between multiple dependent variables and their corresponding predictors. The model is structured as follows:

$$
\begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_m \end{pmatrix} = \begin{pmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1n} \\ 1 & X_{21} & X_{22} & \cdots & X_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{m1} & X_{m2} & \cdots & X_{mn} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_m \end{pmatrix}.
$$

**Components:**

- The left-hand side column vector represents the dependent variables ($Y_1$ to $Y_m$) for each observation.

- The matrix in the middle is the design matrix $\mathbf{X}$, including a leading column of ones for the intercept term $\beta_0$, followed by columns for each of the $n$ independent variables ($X_{11}$ to $X_{mn}$).

- The column vector following the design matrix consists of the regression coefficients ($\beta_0$ to $\beta_n$), quantifying the impact of each independent variable on the dependent variables.

- The final column vector represents the error terms ($\varepsilon_1$ to $\varepsilon_m$), capturing the deviation of the observed dependent variables from their predicted values based on the model.

This matrix equation elegantly summarizes the multivariate linear regression model, allowing for the estimation of multiple dependent variables from a single set of independent variables. Each coefficient in the $\beta$ vector directly corresponds to the influence of an independent variable across all observations, providing insights into how each predictor affects the outcome variables. The inclusion of the error vector $\boldsymbol{\varepsilon}$ acknowledges the presence of unexplained variance, ensuring the model's realism and applicability to real-world data analysis.

## 7.2 Type of Data: When to Use Multivariate Linear Regression

Multivariate Linear Regression is best applied to datasets requiring the analysis of multiple dependent variables influenced by a common set of predictor variables. This approach is ideal for comprehensive studies where understanding the simultaneous impact of predictors on several outcomes is crucial.

**Key Considerations**

- **Interconnected Outcomes:** MLR is most effective in scenarios where multiple outcomes are believed to be influenced by a common set of predictor variables, allowing for an integrated analysis of their effects.

- **Holistic Insights:** It provides insights into how changes in predictor variables can simultaneously affect multiple outcome variables, ideal for complex systems where outcomes are interdependent.

- **Versatile Applications:** MLR is suitable for diverse fields requiring an analysis of how predictors influence multiple aspects simultaneously, such as in socioeconomic studies, healthcare research, and business analytics.

MLR's utility shines in multidimensional analyses, where the interactions and combined effects of predictors on a spectrum of outcomes are of interest. It provides a nuanced understanding of how variables interplay across different dimensions of the dataset.

# 8 Polynomial Regression

Polynomial Regression extends the linear regression by modeling the dependent and independent variables as a polynomial [9]. This approach enables the model to capture nonlinear relationships within the data, providing a more flexible framework compared to simple linear regression. Non-linearity in data refers to the phenomenon where the relationship between independent variables and a dependent variable cannot be accurately described using a straight line. Instead, these relationships exhibit curves or more complex patterns, indicating that the effect of the independent variables on the dependent variable changes at different rates across the range of data. Despite modeling a nonlinear relationship between variables, Polynomial

Regression is considered a special case of linear regression because it is linear in the coefficients that are estimated from the data.

## 8.1 Equation: Formulation and Interpretation

The general form of the Polynomial Regression model can be expressed as:

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \cdots + \beta_n X^n + \varepsilon,$$

where:

- $Y$ represents the dependent variable.

- $X$ is the independent variable.

- $\beta_0, \beta_1, \ldots, \beta_n$ are the coefficients of the model, indicating the influence of each polynomial term of $X$ on $Y$.

- $X^n$ represents the $n$th polynomial term of $X$, allowing the model to capture nonlinear relationships.

- $\varepsilon$ is the error term, accounting for the variability in $Y$ not explained by the polynomial terms of $X$.

This formulation enables the Polynomial Regression model to fit more complex curves than a simple linear regression model, providing a better approximation of the underlying relationship when the data exhibits nonlinearity.

## 8.2 Type of Data: When to Use Polynomial Regression

Polynomial regression offers a flexible approach to modeling relationships between variables, making it particularly effective in situations where traditional linear regression models fall short. It is adept at capturing the nuances of nonlinear relationships, where the effect of an independent variable on the dependent variable changes at different levels of the independent variable. Key considerations for employing polynomial regression include:

- **Nonlinear Relationships:** When data visualization or preliminary analysis suggests that the relationship between variables is curved rather than straight, polynomial regression can model these complexities accurately.

- **Variable Interactions:** It allows for the exploration of how the interaction between two or more independent variables affects the dependent variable.

- **Flexibility in Model Fit:** Polynomial regression can be adjusted to fit the specific curvature of the data by selecting the appropriate degree for the polynomial, providing a tailored approach to data analysis.

**Applicable Case Studies:**

Polynomial regression's versatility makes it applicable across diverse fields for analyzing complex relationships. For instance, in environmental science, it is used to model the nonlinear growth rates of pollutants. In economics, polynomial regression can elucidate nonlinear trends in economic indicators, providing insights that linear models might miss. Similarly, in the biomedical field, it aids in understanding dose-response curves, where the effect of a treatment on an outcome variable may increase at a diminishing rate. These examples demonstrates the method's capabilities in capturing and analyzing the subtleties of nonlinear data relationships.

# 9 Estimation Techniques

Estimation methods such as Ordinary Least Squares (OLS), Maximum Likelihood Estimation (MLE), and Gradient Descent stand out for their widespread application and foundational role linear regression. These

techniques are capable of deriving meaningful insights from data by determining the values of parameters that best represent the underlying relationships within datasets. They enable researchers and analysts to make predictions, infer causal relationships, and understand patterns by fitting mathematical models to observed data. While Ordinary Least Squares (OLS) focuses on minimizing the discrepancies between observed and predicted values in linear regression models, Maximum Likelihood Estimation (MLE) offers a versatile framework for estimating model parameters by maximizing the probability of observing the given data under assumed conditions. Gradient Descent, on the other hand, provides an iterative optimization strategy for finding the minimum of the cost function, applicable in both OLS and MLE, especially useful in scenarios with large datasets or complex models.

## 9.1 Ordinary Least Squares (OLS)

The OLS method is commonly employed to estimate the parameters of a simple linear regression model, minimizing the sum of squared residuals to find the best-fitting line. Attributed to Carl Friedrich Gauss, OLS is a foundational approach in regression analysis, recognized for its effectiveness in statistical estimation under well-defined conditions [5]. This method adheres to the least-squares principle, which involves estimating the parameters of a regression model by minimizing the sum of squared residuals—the squared differences between observed and predicted values.

Central to the OLS methodology is the regression equation $Y = \beta_0 + \beta_1 X + \varepsilon$. The aim is to adjust $\beta_0$ and $\beta_1$ to minimize the discrepancy between the actual $Y$ values and those estimated by the model, $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$, across the dataset. The OLS method focuses on reducing the squared differences ($\hat{\varepsilon}^2$) between observed and predicted values. This approach ensures that the model minimizes errors uniformly across all observations, mitigating any disproportionate impact from outliers or anomalous observations on the model's accuracy.

Squaring the residuals ensures that larger deviations between observed and predicted values are given greater emphasis, promoting a more precise representation of the data's central tendency.

By adopting the least-squares criterion, OLS identifies a regression model that minimizes the sum of squared residuals, thereby deriving parameter estimates that provide the optimal fit to the data. This approach ensures that the estimators are unbiased and possess the minimum variance among all linear estimators for the given dataset.

## 9.2 Mathematical Calculations in Ordinary Least Squares

The Ordinary Least Squares (OLS) method, a cornerstone of linear regression analysis, operates on the principle of minimizing the sum of the squared differences between the observed dependent variable values and those predicted by the linear model. The mathematical foundation of OLS can be defined through a series of steps, as follows:

### 9.2.1 Objective Function

The objective of OLS is to minimize the sum of squared residuals (SSR), where a residual ($\hat{\varepsilon}_i$) is the difference between an observed value ($Y_i$) and its estimated value ($\hat{Y}_i$). Mathematically, the SSR is represented as:

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i,$$

$$\hat{Y}_i = \beta_0 + \beta_1 X_i$$

$$SSR = \sum_{i=1}^{n} (Y_i - (\beta_0 + \beta_1 X_i))^2.$$

$$SSR = \sum_{i=1}^{n} (\hat{\varepsilon}_i)^2 = \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2$$

For a comprehensive understanding of the statistical principles underlying linear regression and the OLS method, the reader is encouraged to consult the work by Soch et al. (2023) [12], which provides an extensive overview of statistical proofs and methodologies.

## 9.3 Estimation of Coefficients in OLS

**Step 1: Partial Derivatives.** The partial derivatives of SSR with respect to $\beta_0$ and $\beta_1$ are taken and set to zero:

$$\frac{\partial SSR}{\partial \beta_0} = -2 \sum_{i=1}^{n} (Y_i - \beta_0 - \beta_1 X_i) = 0,$$

$$\frac{\partial SSR}{\partial \beta_1} = -2 \sum_{i=1}^{n} X_i (Y_i - \beta_0 - \beta_1 X_i) = 0.$$

**Step 2: Solving the Equations.** Setting these derivatives to zero and solving the resulting system of equations yield the OLS estimates:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n} (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^{n} (X_i - \bar{X})^2},$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X},$$

where $\bar{X}$ and $\bar{Y}$ are the sample means of $X$ and $Y$, respectively.

**Interpretation:** $\hat{\beta}_1$ represents the change in $Y$ for a one-unit change in $X$, reflecting the covariance between $X$ and $Y$ relative to the variance of $X$. $\hat{\beta}_0$ adjusts the regression line to pass through the point of averages, ensuring an accurate model fit.

This process outlines the derivation of the OLS estimates, demonstrating the mathematical rigor underlying linear regression analysis.

### 9.3.1 Application of OLS

Upon determining $\hat{\beta}_0$ and $\hat{\beta}_1$, one can predict the value of $Y$ for any given $X$ using the model:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

This model facilitates the exploration of the relationship between $X$ and $Y$, allowing for predictions, inferences, and insights into the underlying data structure.

### 9.3.2 Summary

The OLS method is instrumental in identifying the linear relationship that best fits a given set of data, by minimizing the discrepancies between observed and predicted values. Through the calculation of the least squares estimates, OLS offers a powerful approach for regression analysis, providing a foundation for statistical inference and prediction in various research domains.

## 9.4 Maximum Likelihood Estimation (MLE)

Maximum Likelihood Estimation (MLE) is a statistical method used for estimating the parameters of a statistical model. Unlike Ordinary Least Squares (OLS) which minimizes the sum of squared residuals, MLE seeks to find the parameter values that maximize the likelihood function, given a set of observed data. The likelihood function measures the probability of observing the given sample data under specific model parameters [2].

### 9.4.1 The Likelihood Function

For a linear regression model $Y = \beta_0 + \beta_1 X + \varepsilon$, where $\varepsilon$ follows a normal distribution with mean 0 and variance $\sigma^2$, the likelihood function $L(\beta_0, \beta_1, \sigma^2)$ represents the joint probability of observing the specific set of $Y$ values given the predictors $X$, the model parameters $\beta_0$ and $\beta_1$, and the variance of the error term $\sigma^2$. Mathematically, it can be expressed as:

$$L(\beta_0, \beta_1, \sigma^2) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(Y_i - \beta_0 - \beta_1 X_i)^2}{2\sigma^2}\right),$$

where $n$ is the number of observations.

This expression involves several key components:

- $\prod_{i=1}^{n}$: Signifies the product of terms for all $n$ observations, combining the individual probabilities of observing each $Y_i$ given $X_i$, $\beta_0$, $\beta_1$, and $\sigma^2$.

- $\frac{1}{\sqrt{2\pi\sigma^2}}$: Acts as the normalization factor of the normal distribution's probability density function, ensuring the total area under the curve equals 1.

- $\exp\left(-\frac{(Y_i - \beta_0 - \beta_1 X_i)^2}{2\sigma^2}\right)$: The exponential function represents the likelihood of observing each $Y_i$, factoring in the squared difference between observed and predicted $Y_i$ values, normalized by the variance $\sigma^2$.

### 9.4.2 Significance

In the context of linear regression models with normally distributed errors, the likelihood function quantifies the probability of observing the given set of $Y$ values for specified model parameters ($\beta_0$, $\beta_1$) and error variance ($\sigma^2$). Maximizing this function with respect to the parameters enables the identification of estimates most likely to have generated the observed data, foundational to MLE in model fitting.

### 9.4.3 Estimation Process

The MLE process begins by transforming the likelihood function into a log-likelihood function. This step simplifies the subsequent mathematical operations due to the properties of logarithms.

## 9.5 Log-Likelihood Function

The log-likelihood function for the linear regression model is given by:

$$\ell(\beta_0, \beta_1, \sigma^2) = \log L(\beta_0, \beta_1, \sigma^2)$$

$$= -\frac{n}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma^2}\sum_{i=1}^{n}(Y_i - \beta_0 - \beta_1 X_i)^2.$$

This equation transforms the likelihood of observing our data under the assumed linear model into a form that is more tractable for optimization by taking the natural logarithm, denoted by log. The log-likelihood function consists of two main components:

- The first component $-\frac{n}{2}\log(2\pi\sigma^2)$ corresponds to the constant part of the normal distribution's probability density function for the residuals, where $n$ is the number of observations, and $\sigma^2$ is the variance of the residuals. This term ensures that the likelihood function integrates to 1 over all possible values of the response variable $Y$, adhering to the properties of a probability distribution.

- The second component $-\frac{1}{2\sigma^2}\sum_{i=1}^{n}(Y_i - \beta_0 - \beta_1 X_i)^2$ is the sum of the squared differences between the observed values $Y_i$ and the values predicted by the linear regression model $\beta_0 + \beta_1 X_i$. These differences are squared to address both positive and negative deviations equally, and they are scaled by the variance of the errors, $\sigma^2$, reflecting the degree of dispersion around the regression line. This term effectively quantifies the overall 'goodness of fit' of the model to the observed data.

The optimization of the log-likelihood function involves finding the values of $\beta_0$, $\beta_1$, and $\sigma^2$ that maximize this function, which is equivalent to minimizing the sum of squared differences between observed and predicted.

### 9.5.1 Optimization

The optimization step in MLE seeks to determine the parameter estimates that maximize the log-likelihood function $\ell(\beta_0, \beta_1, \sigma^2)$. This involves calculus-based optimization techniques:

- **Partial Derivatives:** The process begins by taking the partial derivatives of the log-likelihood function with respect to each parameter—$\beta_0$, $\beta_1$, and $\sigma^2$. These derivatives reflect the slope of the log-likelihood function with respect to the parameters and are used to locate its maximum point.

$$\frac{\partial \ell}{\partial \beta_0} = -\frac{1}{\sigma^2}\sum_{i=1}^{n}(Y_i - \beta_0 - \beta_1 X_i),$$

$$\frac{\partial \ell}{\partial \beta_1} = -\frac{1}{\sigma^2}\sum_{i=1}^{n}X_i(Y_i - \beta_0 - \beta_1 X_i),$$

$$\frac{\partial \ell}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}\sum_{i=1}^{n}(Y_i - \beta_0 - \beta_1 X_i)^2.$$

**Setting to Zero:** We then set each partial derivative to zero to obtain the system of likelihood equations:

$$\sum_{i=1}^{n}(Y_i - \beta_0 - \beta_1 X_i) = 0,$$

$$\sum_{i=1}^{n}X_i(Y_i - \beta_0 - \beta_1 X_i) = 0,$$

$$-\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}\sum_{i=1}^{n}(Y_i - \beta_0 - \beta_1 X_i)^2 = 0.$$

**Solving the Equations:** We solve these equations simultaneously to find the values of $\beta_0$, $\beta_1$, and $\sigma^2$ that maximize the log-likelihood function.

These equations are solved to obtain the MLE of $\beta_0$, $\beta_1$, and $\sigma^2$. However, they often do not have a closed-form solution and are typically solved using numerical optimization techniques. Computational tools such as R, Python, or specialized statistical software are employed to perform this optimization efficiently.

**Note on Computational Tools:** The complexity of these equations usually necessitates the use of computational methods. These tools implement algorithms that can handle the iterative process required to navigate the multi-dimensional space of the log-likelihood function and converge on the maximum likelihood estimates for the parameters.

**Ensuring Maximum:** Once we find the solutions for $\beta_0$, $\beta_1$, and $\sigma^2$, we need to ensure that these provide a maximum to the log-likelihood function. For this, we evaluate the second-order conditions, which involve the Hessian matrix—comprising the second partial derivatives of the log-likelihood function:

$$H = \begin{bmatrix} \frac{\partial^2 \ell}{\partial \beta_0^2} & \frac{\partial^2 \ell}{\partial \beta_0 \partial \beta_1} & \frac{\partial^2 \ell}{\partial \beta_0 \partial \sigma^2} \\ \frac{\partial^2 \ell}{\partial \beta_1 \partial \beta_0} & \frac{\partial^2 \ell}{\partial \beta_1^2} & \frac{\partial^2 \ell}{\partial \beta_1 \partial \sigma^2} \\ \frac{\partial^2 \ell}{\partial \sigma^2 \partial \beta_0} & \frac{\partial^2 \ell}{\partial \sigma^2 \partial \beta_1} & \frac{\partial^2 \ell}{\partial \sigma^2 \partial \sigma^2} \end{bmatrix}.$$

The Hessian must be negative definite at the solution for the solution to correspond to a maximum. This means all eigenvalues of the Hessian matrix should be negative, indicating that the log-likelihood function is concave down at the maximum point.

This optimization is a critical part of MLE, as it enables the estimation of model parameters that are most likely to have generated the observed data. The MLE method is powerful in statistical inference, providing estimators with desirable properties such as consistency and efficiency under certain conditions.

### 9.5.2 Advantages of MLE

MLE offers several advantages, including consistency (estimators converge to the true parameter values as the sample size increases), efficiency (MLE achieves the lowest possible variance among all unbiased estimators), and invariance (function transformations of MLE estimators are also MLE estimators).

While MLE can be more computationally intensive than OLS, especially for complex models, it provides a flexible framework for estimation that is applicable to a wide range of statistical models, not limited to linear regression. Its property of maximizing the probability of observing the sample data makes it a powerful tool in statistical inference.

## 9.6   Gradient Descent

Gradient Descent (GD) is an optimization method in machine learning, particularly effective for minimizing cost or loss functions in various models, including neural networks [7]. It operates by iteratively adjusting parameters to find the function's minimum, utilizing the gradient's direction to guide these adjustments.

### 9.6.1   Mathematical Foundation of Gradient Descent

The mathematical foundation of Gradient Descent lies in its iterative process to minimize the cost function $J(\theta)$ of a model. It updates parameters $\theta$ by computing the gradient $\nabla_\theta J(\theta)$, which points to the steepest ascent, and then moves in the opposite direction:

$$\theta_{\text{new}} = \theta_{\text{old}} - \eta \cdot \nabla_\theta J(\theta),$$

where $\theta$ represents the model parameters, $\eta$ denotes the learning rate, and $\nabla_\theta J(\theta)$ is the gradient of the cost function $J(\theta)$ at the parameters $\theta$. The learning rate $\eta$ controls the step size in each iteration, balancing the speed of convergence against the risk of overshooting the minimum. The process continues until convergence, typically when the gradient is close to zero, indicating that the minimum of $J(\theta)$ has been reached or when improvements become negligible.

In the context of machine learning, model parameters ($\theta$) are the variables in the model that are learned from the training data. For example, in linear regression, $\theta$ could represent the coefficients of the model, including the slope ($\beta_1$) and intercept ($\beta_0$) in a simple linear regression ($Y = \beta_0 + \beta_1 X + \varepsilon$). These parameters define how the input data ($X$) is mapped to the predicted output ($Y$).

When we talk about optimizing these parameters using Gradient Descent, we are iteratively adjusting them to minimize the cost function ($J(\theta)$), which measures how well the model's predictions correspond to the actual data. Each iteration updates the parameters to move closer to the set of values that result in the lowest possible cost, which would represent the best-fitting model.

### 9.6.2 Learing Rate

Learning rate is a critical parameter in gradient descent optimization, it is set as a positive scalar determining the size of the step taken in the direction opposite to the gradient. A too large learning rate can lead to overshooting the minimum, while a too small rate can result in a long convergence time.

Adaptive learning rate techniques, like AdaGrad, RMSprop, and Adam, adjust the learning rate during training to improve convergence. AdaGrad adapts the learning rate to parameters, making it smaller for frequently occurring features. RMSprop modifies the learning rate based on recent gradients to prevent it from drastically changing. Adam combines RMSprop's benefits with momentum, updating learning rates based on first and second moments of gradients, enhancing the efficiency of the learning process [10].

## 9.7 Convergence Properties

The convergence properties of Gradient Descent are heavily influenced by the characteristics of the objective function $f(\theta)$ and the choice of the learning rate $\eta$. For convex functions, and under appropriate conditions on the learning rate, Gradient Descent is guaranteed to converge to the global minimum. For non-convex functions, Gradient Descent may converge to a local minimum or a saddle point.

## 9.8 Definition of Gradient

For a function $f : \mathbb{R}^n \to \mathbb{R}$ that takes an $n$-dimensional vector $\boldsymbol{\theta} = (\theta_1, \theta_2, \ldots, \theta_n)$ as input and produces a scalar output, the gradient of $f$ at $\boldsymbol{\theta}$, denoted $\nabla f(\boldsymbol{\theta})$, is the vector of partial derivatives of $f$ with respect to each component of $\boldsymbol{\theta}$. Mathematically, it is defined as:

$$\nabla f(\boldsymbol{\theta}) = \left( \frac{\partial f}{\partial \theta_1}, \frac{\partial f}{\partial \theta_2}, \ldots, \frac{\partial f}{\partial \theta_n} \right)$$

### 9.8.1 Computing the Gradient

The computation of each component of the gradient vector involves taking the partial derivative of $f$ with respect to the corresponding component of $\boldsymbol{\theta}$. In linear regression models, $\boldsymbol{\theta}$ represents the parameters (or weights) that we are trying to learn, which define the line (or hyperplane in higher dimensions) that best fits our data. The partial derivative $\frac{\partial f}{\partial \theta_i}$ measures how $f$ changes as $\theta_i$ varies, holding all other components of $\boldsymbol{\theta}$ constant. This is calculated as follows for each $i$ from 1 to $n$:

1. Identify the functional form of $f$ with respect to $\theta_i$: This involves expressing $f$ explicitly in terms of $\theta_i$, while considering the other components of $\boldsymbol{\theta}$ as constants.

2. Apply differentiation rules: Use standard differentiation rules (e.g., power rule, product rule, chain rule) to find the derivative of $f$ with respect to $\theta_i$.

3. Evaluate at $\boldsymbol{\theta_t}$: Once the derivative $\frac{\partial f}{\partial \theta_i}$ is obtained, evaluate it at the current parameter vector $\boldsymbol{\theta_t}$ to get the gradient's component for the current iteration.

### 9.8.2 Example: Gradient of a Quadratic Function

Consider a simple quadratic function in two variables, $f(\boldsymbol{\theta}) = \theta_1^2 + 3\theta_2^2 + 2\theta_1\theta_2$, where $\boldsymbol{\theta} = (\theta_1, \theta_2)$. The gradient of $f$ at $\boldsymbol{\theta}$ is:

$$\nabla f(\boldsymbol{\theta}) = \left( \frac{\partial f}{\partial \theta_1}, \frac{\partial f}{\partial \theta_2} \right)$$

Computing each partial derivative:

$$\frac{\partial f}{\partial \theta_1} = 2\theta_1 + 2\theta_2, \quad \frac{\partial f}{\partial \theta_2} = 6\theta_2 + 2\theta_1$$

So, the gradient vector is $\nabla f(\boldsymbol{\theta}) = (2\theta_1 + 2\theta_2, 6\theta_2 + 2\theta_1)$.

### 9.8.3  Importance in Gradient Descent

In Gradient Descent, the gradient $\nabla f(\boldsymbol{\theta_t})$ points in the direction of the steepest increase of $f$ at $\boldsymbol{\theta_t}$. By moving in the opposite direction (i.e., $-\nabla f(\boldsymbol{\theta_t})$), we iteratively adjust $\boldsymbol{\theta}$ towards the minimum of $f$. The magnitude of each step is controlled by the learning rate $\eta$, ensuring that the updates are proportional to the steepness of the function.

**Variants:**

- *Vanilla Gradient Descent:* Applies updates after computing the gradient using the entire dataset. This method is best suited for smaller datasets where the computational overhead of processing the entire dataset at once is manageable. It ensures a smooth and stable convergence to the global minimum for convex error surfaces, making it ideal for datasets where precision is key and computational resources are not a limiting factor.

- *Stochastic Gradient Descent (SGD):* Updates parameters for each data point, enhancing efficiency and potentially escaping local minima. SGD is particularly effective for large datasets and situations where the data may not fit entirely in memory. Its ability to update weights incrementally makes it suitable for online learning scenarios and applications requiring real-time updates, such as recommendation systems or dynamic pricing models. The stochastic nature of SGD helps it to escape local minima, making it a good choice for non-convex optimization problems.

- *Mini-batch SGD:* Strikes a balance between Vanilla GD and SGD by using subsets of the dataset (mini-batches) for updates, optimizing performance and convergence. This variant is best utilized in scenarios where one needs to efficiently handle moderately large datasets with the advantage of vectorized implementation for faster computation. Mini-batch SGD combines the stability of Vanilla GD with the efficiency and ability to escape local minima of SGD, making it highly versatile for a wide range of applications from image processing to natural language processing tasks, where data size and computational efficiency are critical concerns.

### 9.8.4  Gradient Descent Algorithm

The gradient descent algorithm can be summarized in the following steps:

1. **Initialization:** Start with initial guesses for the parameters to be optimized.

   At the beginning of the gradient descent process, parameters $(\theta)$ are initialized. This can be expressed as:
   $$\theta^{(0)}$$
   where the superscript $(0)$ denotes the initial state before optimization begins. The initialization might involve setting all parameters to zero $(\theta^{(0)} = 0)$ or assigning random values, which can influence the convergence speed and quality of the algorithm.

2. **Compute the Gradient:** For the function $f(\boldsymbol{\theta})$ representing the cost or loss, the gradient at step $t$ is computed as:
   $$\nabla f(\boldsymbol{\theta}(t)) = \left( \frac{\partial f}{\partial \theta_1}, \frac{\partial f}{\partial \theta_2}, \ldots, \frac{\partial f}{\partial \theta_n} \right)(t)$$

   This represents the direction of steepest ascent at $\boldsymbol{\theta}(t)$, with each component indicating how the cost changes with respect to each parameter.

3. **Update Parameters**

   The parameter update rule at step $t$ uses the gradient to adjust $\boldsymbol{\theta}$ in the direction of steepest descent:
   $$\boldsymbol{\theta}(t+1) = \boldsymbol{\theta}(t) - \eta \nabla f(\boldsymbol{\theta}(t))$$

   Here, $\eta$ is the learning rate, and the subtraction indicates movement against the gradient, aiming to reduce the cost.

4. **Iteration**

   This step is the loop where steps 2 (Compute the Gradient) and 3 (Update Parameters) are repeated:

   Repeat until convergence:

   - Step 2

   - Step 3

5. **Convergence Check**

   The process continues until a stopping criterion is satisfied, which could be one of several conditions, such as:

   - A maximum number of iterations is reached: $t = T_{\max}$

   - The change in cost function is below a threshold: $|f(\boldsymbol{\theta}(t+1)) - f(\boldsymbol{\theta}(t))| < \epsilon$

   - The magnitude of the gradient is below a threshold: $\|\nabla f(\boldsymbol{\theta}(t))\| < \epsilon$

   Here, $\epsilon$ is a small positive value indicating the tolerance for stopping.

Each step aims to progressively reduce the cost function value, moving the parameters towards the optimal set that minimizes the cost.

## 9.9 Challenges with Learning Rate

### 9.9.1 Challenges:

- **Too Small Learning Rate:** Leads to slow convergence, requiring many iterations to reach the optimal solution, which can be computationally expensive and time-consuming.

- **Too Large Learning Rate:** Can cause the algorithm to overshoot the minimum, leading to divergence or oscillation around the optimal solution, rather than converging.

### 9.9.2 Solutions:

- **Trial and Error:** Start with a small learning rate and gradually increase it to find a balance between convergence speed and stability.

- **Learning Rate Scheduling:** Adjust the learning rate dynamically during training. For example, decrease the learning rate as the number of iterations increases to allow for finer adjustments as the solution approaches the optimum.

- **Adaptive Learning Rate:** Use optimization algorithms like AdaGrad, RMSProp, or Adam, which automatically adjust the learning rate for each parameter based on the historical gradient information, although these are more commonly applied in more complex models beyond simple linear regression .

## 9.10 Challenges in Convergence Criteria

### 9.10.1 Challenges:

- **Determining When to Stop:** It can be challenging to set a precise rule for when the algorithm has converged to the optimal solution. Stopping too early might mean the solution is not optimal, while stopping too late might waste computational resources.

- **Oscillation:** Near the optimum, especially with a learning rate that's too high, the updates might oscillate, making it hard to determine convergence.

### 9.10.2   Solutions:

- **Convergence Threshold:** Define a threshold for the change in the cost function between iterations. If the change is below this threshold, the algorithm can be considered to have converged.

- **Gradient Norm:** Another approach is to monitor the norm of the gradient. If the norm falls below a certain threshold, it indicates that the slope of the cost function is very gentle, suggesting convergence.

- **Maximum Iterations:** Set a maximum number of iterations to ensure the algorithm terminates even if the above criteria are not met. This acts as a fail-safe to prevent endless looping in cases where convergence is slow or oscillations occur.

## 9.11   Gradient Descent Diagnostics

To effectively manage these challenges, specific diagnostic strategies can be employed:

- **Plot Learning Curves:** Graph the cost function value over iterations. A properly tuned learning rate should show a smooth and steady decline towards the minimum. If the cost function value increases or exhibits erratic behavior, this might indicate a learning rate that's too high.

- **Adjustment in Small Increments:** When tuning the learning rate, make adjustments in small increments and observe the impact on convergence behavior.

- **Gradient Checking:** Although more common in debugging neural networks, gradient checking can be used in linear regression to ensure that the gradient computation is correct. This involves comparing the computed gradient against a numerical approximation of the gradient.

- **Early Stopping:** Monitor performance on a validation set (if available) during training. If performance on the validation set begins to worsen, it may be a signal to stop training to prevent overfitting. While more relevant for complex models, it can also be a consideration in linear regression when regularization is used.

By carefully managing learning rate and convergence criteria, and employing diagnostic strategies, challenges in using Gradient Descent for linear regression models can be effectively addressed, leading to efficient and successful model optimization.

**Contextual Applications:** GD's flexibility makes it suited for a wide range of applications, from simple linear regression models to complex neural networks, demonstrating its foundational role in the optimization landscape of machine learning.

# 10   Conclusion

This paper has presented a collection of concepts, methods, and applications of linear regression in the context of machine learning. Starting with a fundamental definition, the paper outlined the importance of linear regression in predictive modeling and highlighted a variety of its applications.

In reviewing the theoretical background, we examined the foundational concepts of linear regression, the types of data it can be applied to, and the intrinsic assumptions upon which the model's validity relies. The exploration of simple linear regression provided clarity on its mathematical formulation and practical guidance on its appropriate data contexts.

The discussion extended to multiple linear regression, where the inclusion of multiple predictors was shown to enhance the model's complexity and applicability. Similarly, multivariate linear regression was explored, demonstrating its capacity to model multiple dependent variables simultaneously.

The treatment of polynomial regression addressed the incorporation of non-linear relationships, offering a methodological bridge to modeling more complex data structures. This section emphasized the identification of suitable scenarios for polynomial regression and the interpretation of its coefficients.

Various estimation techniques were briefly introduced. The Ordinary Least Squares method was detailed as the most common approach, with Maximum Likelihood Estimation and Gradient Descent providing alternative frameworks under different circumstances. Each method's mathematical basis and practical implications for model estimation were evaluated.

In conclusion, linear regression emerges from this exposition as a versatile and robust statistical tool in machine learning. As a foundational statistical method, linear regression sets the stage for further research and application in more complex and nuanced machine learning models.

# References

[1] Linear regression. `http://www.stat.yale.edu/Courses/1997-98/101/linreg.htm`, 1998. Accessed: 2024-02-12.

[2] Lesson 1.2: Introduction to linear regression analysis. `https://online.stat.psu.edu/stat415/lesson/1/1.2`, Access Year. Accessed: [Insert Access Date Here].

[3] Business Econometrics. Explaining Error Term in Regression Equation, 2021. YouTube video, published on Jun 3, 2021. Available from: https://www.youtube.com/watch?v=83O3J0nNr5s. Accessed on: [Your Access Date Here].

[4] David A. Freedman. What is the error term in a regression equation? `https://www.stat.berkeley.edu/~census/epsilon.pdf`, Year. Accessed: Your Access Date Here.

[5] Damodar N. Gujarati and Dawn C. Porter. *Basic Econometrics*. Tata McGraw Hill, 5 edition, 2007.

[6] Bertha Hidalgo and Melody Goodman. Multivariate or multivariable regression? *American Journal of Public Health*, 103(1):39–40, Jan 2013.

[7] Jun Lu. Gradient descent, stochastic optimization, and other tales, 2024.

[8] BUŞE Lucian, Assist Mirela Ganea, Lect Daniel CÎRCIUMARU, et al. Using linear regression in the analysis of financial-economic performances. *Annals of University of Craiova-Economic Sciences Series*, 2(38):32–43, 2010.

[9] Aleksandar Pečkov. *A Machine Learning Approach to Polynomial Regression*. Doctoral dissertation, Jožef Stefan International Postgraduate School, Ljubljana, Slovenia, October 2012. Evaluation Board: Prof. Dr. Bogdan Filipič, Prof. Dr. João Gama, Prof. Dr. Djani Juričić.

[10] Sebastian Ruder. An overview of gradient descent optimization algorithms. *CoRR*, abs/1609.04747, 2016.

[11] Ramaswamy Seethalakshmi. Analysis of stock market predictor variables using linear regression. *International Journal of Pure and Applied Mathematics*, 119(15):369–378, 2018.

[12] Joram Soch, The Book of Statistical Proofs, Maja, Pietro Monticone, Thomas J. Faulkenberry, Alex Kipnis, Kenneth Petrykowski, Carsten Allefeld, Heiner Atze, Adam Knapp, Ciarán D. McInerney, Lo4ding00, and amvosk. Statproofbook/statproofbook.github.io: Statproofbook 2023, January 2024. `https://doi.org/10.5281/zenodo.4305949`.

[13] Teresa Angela Trunfio, Arianna Scala, Cristiana Giglio, Giovanni Rossi, Anna Borrelli, Maria Romano, and Giovanni Improta. Multiple regression model to analyze the total los for patients undergoing laparoscopic appendectomy. *BMC Medical Informatics and Decision Making*, 22(1):141, 2022.