RESEARCH ARTICLE              OPEN ACCESS

# Fitting Regression Models and Experimental Designs for Accurate Data Analysis

Kannoju Divya*, Viraja Mukthavarapu

*( Department of Statistics, Kakatiya University, Warangal

Email: sree2846@gmail.com)

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

---

## Abstract:

Data analysis helps us understand the world. Data analysis uses tools like regression models and experimental designs. Regression models reveal how different factors influence one another while experimental designs ensure ideas undergo fair and unbiased testing. This study explores how these two methods complement each other to enable accurate predictions and informed decisions. Simple tools explain complex math concepts like variance and $R^2$ and solve real-world problems in India. These approaches affect farming, healthcare, education and online shopping. In the final section the conclusion explores the future of applied analytics and suggests ways in which researchers might further improve their methods.

*Keywords* — Covariance analysis, Regression analysis, Goodness of Fit ($R^2$), Experimental Designs, Predictive Modelling.

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

---

## I. INTRODUCTION

Data drive the modern world. Each digital transaction and farm harvest generates new data. In India this data explosion is occurring faster than ever. The market now offers more smartphones and cheaper internet access. People leave digital footprints but raw data remain useless until analysis adds value. Data analysis converts raw numbers into useful insights.

Different sectors in India urgently require accurate data analysis. Farmers can determine the correct amount of fertiliser required to produce the best crop yield. Teachers can identify the teaching methods that help students learn more quickly. Businesses identify where to allocate advertising funds to achieve the highest sales. Proven mathematical methods provide reliable answers to these questions whereas guessing alone is insufficient.

The regression model is a powerful tool that works like a line drawn through a scatterplot. It positions this line as close as possible to all points and maintains a straight line. By drawing this line through the data points the regression model predicts outcomes and shows the relationship between two or more variables. The "dependent variable" is the result that the regression model predicts in cases such as test scores or total sales. The "independent variable" includes factors that the regression model considers as causes of changes in the outcome in situations such as hours spent studying or money spent on advertising.

The occurrence of things at the same time does not prove that one thing causes the other. For example ice cream sales and sunburns both increase during the summer. However, ice cream does not cause sunburn because the hot summer sun is the cause of both. Experimental designs demonstrate when one thing causes another and establish strict rules to prevent external factors from affecting the results.

---

## II. ADVANCED REGRESSION MODELS

A single regression model cannot solve every problem. Researchers use different types based on the data they have.

### 2.1 Simple and Multiple Linear Regression

Simple Linear Regression analyses the effect of a one cause on one outcome. The mathematical formula is as follows:

$$Y = \beta_0 + \beta_1 X + \epsilon$$

Here $Y$ is the outcome $X$ is the cause $\beta_0$ is the starting point $\beta_1$ is the slope of the line and $\epsilon$ is the error or noise in the data.

Real life is often complex and in most cases more than one factor affects an outcome. Multiple regression accounts for several factors such as size or location or age when predicting the price of a house. The formula for Multiple Linear Regression is a complex matrix equation:

$$Y = X\beta + \epsilon$$

Where $Y$ is a vector of observations $X$ is a design matrix of predictors $\beta$ is a vector of coefficients and $\epsilon$ represents the error vector. The Ordinary Least Squares estimator minimizes the sum of squared residuals to find the optimal coefficients:

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

### 2.2 Logistic Regression

Sometimes the outcome is not a number but rather a category such as whether a customer will click on an ad. Logistic regression is used for these situations. The model estimates the probability using a sigmoid link function:

$$P(Y = 1|X) = \frac{1}{1 + e^{-X\beta}}$$

### 2.3 Mathematical Assumptions and Diagnostics

The data must satisfy specific assumptions for a regression model to work correctly. If the data does not satisfy these assumptions it causes the regression model to produce incorrect predictions.

**Linearity:** The relationship between the variables is illustrated most clearly by a straight line.

**Independence:** Each data point must remain distinct from the others.

**Homoscedasticity:** The spread of errors in the predictions maintains consistency across all values of the independent variable. When the spread of errors becomes excessive the spread of errors reduces the reliability of the results. The spread of errors uses the application of logarithms in order to maintain control.

**Normality:** The errors should form a bell curve.

Researchers must check for unusual data points. Outliers are data points which are significantly higher or lower than the others and these can skew the results so researchers need to check them carefully. High-leverage points are extreme values among the cause variables. For instance, when there is a massive or unusual ad spend, it will forcefully pull the regression line towards it and distort the model.

### 2.4 Mathematical Assumptions and Diagnostics

Tests like the p-value show if results happen by chance. A low p-value (less than 0.05) means the relationship is real not just luck. Confidence intervals indicate the range in which the true value is likely to be found.

The R² value demonstrates how well the model fits the data by explaining the change in the outcome. The formula is:

$$R^2 = 1 - \frac{\sum(y_i - \hat{y}_i)^2}{\sum(y_i - \bar{y})^2}$$

A higher $R^2$ means the model is more accurate.

Table 1 demonstrates this using an Indian retail company predicting sales based on its advertising budget.

TABLE 1
*MARKETING BUDGET VS. SALES REGRESSION DATA*

| Advertising Budget (₹) ($X$) | Actual Sales (Units) ($Y$) | Predicted Sales ($\hat{Y}$) | Residuals (Error) ($\epsilon$) |
|---|---|---|---|
| 500 | 200 | 180 | +20 |
| 700 | 240 | 250 | -10 |
| 1000 | 350 | 330 | +20 |
| 1500 | 500 | 460 | +40 |
| 2000 | 600 | 590 | +10 |

*Note: Small residuals mean the model is doing a good job.*

## III. EXPERIMENTAL DESIGNS

Regression processes data while experimental designs serve as the rulebook for gathering data. If data is collected poorly the math remains correct while the conclusions become incorrect. To create a fair environment scientists follow three basic rules:

**Randomisation:** Assigning subjects to groups at random helps to eliminate hidden bias.

**Replication:** Testing something once is never sufficient. Repeating the experiment on a wide range of subjects ensures reliable results.

**Blocking:** Grouping similar subjects together before the test begins to keep things fair.

Researchers use different layouts depending on the goal. Researchers find the Randomised Complete Block Design (RCBD) is popular in farming. Researchers use Split-Plot Design when treatments are hard to apply to small areas and they often use Cluster-Randomised Trials in schools and hospitals.

Analytical power tells the chance of finding a real effect. When you have a larger sample size the power of a test increases in the same way that using a stronger magnifying glass lets you observe details more clearly. To mathematically ensure an experiment has enough strength to detect a difference (Δ), researchers use standard sample size estimation equations based on variance ($\sigma^2$):

$$n = \frac{2(Z_{\alpha/2} + Z_\beta)^2 \sigma^2}{\Delta^2}$$

.

TABLE 2
*RELATIONSHIP BETWEEN ANALYTICAL POWER AND SAMPLE VOLUMES*

| Desired Power Level ($1 - \beta$) | Required Sample Size ($n$ per group) | Minimum Detectable Effect |
|---|---|---|
| 80% | 64 | Medium |
| 90% | 85 | Medium |
| 95% | 105 | Small |

## IV. INTEGRATING REGRESSION AND EXPERIMENTAL DESIGNS

Regression and experimental designs are closely linked and form different aspects of the same process. This integration happens through Variance Partitioning. Analysts slice the pie chart of variance into the variation caused by the treatment and the leftover unexplained noise. Mathematically:

$$SST = SSTr + SSE$$

### 4.1 Analysis of Covariance (ANCOVA)

Sometimes experimental subjects are not equal when the test starts. Analysis of Covariance (ANCOVA) remains an excellent approach drawing on both traditional testing and regression methods. When subjects have a pre-existing trait called a covariate that affects the outcome, ANCOVA mathematically adjusts the final scores to level the playing field.

The advanced ANCOVA mathematical model is:

$$Y_{ij} = \mu + \tau_i + \beta(X_{ij} - \bar{X}) + \epsilon_{ij}$$

Table 3 shows a clinical trial where dosage is the main treatment and age is the covariate.
.

TABLE 3
*CLINICAL TRIAL RECOVERY TIME (ANCOVA MODEL)*

| Group | Dosage (mg) | Age (years) | Recovery Time | Predicted Recovery | Residuals |
|---|---|---|---|---|---|
| Low Dose | 50 | 45 | 10 days | 10.5 days | -0.5 |
| High Dose | 100 | 60 | 8 days | 7.5 days | +0.5 |

*Note: ANCOVA mathematically adjusts for age, stripping away the penalty so doctors can see the pure unclouded effect of the medicine itself.*

## 5. CASE STUDIES AND APPLICATIONS IN INDIA

### 5.1 Agriculture: Zero-Tillage Farming

The Indian Council of Agricultural Research (ICAR) leads the way in testing new farming methods. A major test involved "Zero-Tillage" farming in Punjab and Haryana which means planting seeds directly without ploughing the soil to save water and diesel.

ICAR used an RCBD, dividing large fields into blocks based on soil salt levels. Inside each block, plots were randomly chosen for traditional ploughing and Zero-Tillage. Multiple regression was used to mathematically remove the effects of the salty soil, isolating the exact yield boost caused by Zero-Tillage.

TABLE 3
*YIELD EVALUATION USING MULTIPLE REGRESSION IN ZERO-TILLAGE*

| Treatment | Nitrogen (kg/ha) | Soil Salinity Index | Mean Yield (tons/ha) | Adjusted Yield |
|---|---|---|---|---|
| Traditional | 120 | 1.4 | 4.1 | 4.2 |
| Zero-Tillage | 120 | 1.8 | 4.5 | 4.8 |

### 5.2 Healthcare: Generic Medicines

The National Health Mission (NHM) runs clinical trials to prove that cheaper generic drugs are effective. Researchers conducting Randomised Controlled Trials (RCTs) allocate patients to receive either a branded drug or a generic one and in this way researchers ensure each group receives fair treatment for comparison. ANCOVA addresses differences in patient characteristics by adjusting for covariates such as age and baseline health. The regression model mathematically proves that no difference in healing exists between generic and branded drugs.

### 5.3 Education: "Teaching at the Right Level"

The NGO Pratham created a teaching method called "Teaching at the Right Level" (TaRL) grouping kids by reading ability rather than age. Researchers ran huge cluster-randomized trials across schools in Uttar Pradesh. Researchers used multiple regression to analyse test scores taking into account student attendance and family wealth. The math proved that TaRL dramatically improved reading scores no matter how poor the student's background was.

### 5.4 Marketing: Vernacular E-Commerce

Companies like Flipkart and Amazon use A/B testing as a digital experiment to see if changing apps to local languages like Hindi or Tamil will increase sales. A logistic regression model controls for covariates like internet speed, proving exactly how much the local language boosted sales to help companies decide where to invest.

**CONCLUSIONS**

Regression acts as the math engine helping predict the future by drawing lines through past data. Experimental designs are the rulebook ensuring tests are fair and free from hidden biases. When combined using tools like ANCOVA the messy noise of the real world can be stripped away to find clear facts. These analytical tools solve massive problems in India today by helping farmers grow more food, enabling doctors to prescribe

cheaper medicine and allowing businesses to reach new customers.

Room for improvement remains in India. Researchers must be forced by academic journals to check assumptions; publishing an $R^2$ score without checking for outliers or errors makes a model useless. Indian colleges must change teaching methods so students learn experimental designs and predictive regression simultaneously using hands-on practice with real-world Indian data. India can lead the world in data analysis by demanding better diagnostic checks and modernising how methods are taught. Mastering these tools will drive innovation, increase efficiency and help leaders make smart decisions for the future.

## REFERENCES

[1] Box, G. E. P., Hunter, J. S., and Hunter, W. G. (2005). Quantitative Methods for Experimenters: Design, Innovation, and Discovery, 2nd edition. Wiley-Interscience.

[2] Brook, R. J., and Arnold, G. C. (1985). Applied Regression Analysis and Experimental Designs, Marcel Dekker.

[3] Chatterjee, S., and Hadi, A. S. (2012). Regression Analysis by Example, 5th edition. Wiley.

[4] Das, M. N., and Giri, N. C. (1986). Design and Analysis of Experiments, 2nd edition. New Age International.

[5] Draper, N. R., and Smith, H. (1998). Applied Regression Analysis, 3rd edition. Wiley.

[6] Fisher, R. A. (1935). The Design of Experiments. Oliver and Boyd.

[7] Gelman, A., and Hill, J. (2007). Data Analysis Using Regression and Multilevel/Hierarchical Models. Cambridge University Press.

[8] Goon, A. M., Gupta, M. K., and Dasgupta, B. (2002). Fundamentals of Data Analytics, Volumes 1 and 2. The World Press.

[9] Gupta, S. C., and Kapoor, V. K. (2014). Fundamentals of Applied Analytics, 4th edition. Sultan Chand & Sons.

[10] Kothari, C. R., and Garg, G. (2014). Research Methodology: Methods and Techniques, 3rd edition. New Age International.

[11] Montgomery, D. C. (2019). Design and Analysis of Experiments, 10th edition. John Wiley & Sons.

[12] Neter, J., Kutner, M. H., Nachtsheim, C. J., and Wasserman, W. (1996). Applied Linear Mathematical Models, 4th edition. Irwin.

[13] Panse, V. G., and Sukhatme, P. V. (1985). Quantitative Methods for Agricultural Workers. Indian Council of Agricultural Research.

[14] Rao, C. R. (1973). Linear Statistical Inference and Its Applications, 2nd edition. Wiley.

[15] Sancheti, D. C., and Kapoor, V. K. (2010). Analytics: Theory, Methods, and Application. Sultan Chand & Sons.