

Conducting and Presenting Multiple Linear Regression Analysis Using R

Maria Paula Rojas Conejo 400962585
Fresenius University of Applied Science
Data Science and Data Analytics (WS 2025/26)
Prof. Dr. Stephan Huber
2025-12-10

Author Note

Correspondence concerning this article should be addressed to Maria Paula Rojas Conejo 400962585, Email: rojas_conejo.maria@stud.hs-fresenius.de

Abstract

This document outlines the requirements for successfully contributing to the course through a presentation and a project report. It introduces the process, establishes the standards, and provides tips how students can effectively meet the expectations.

Conducting and Presenting Multiple Linear Regression Analysis Using R

Word count: 1406

1 Introduction

According to Alexopoulos (2010) when it comes to data analysis one of the main goals is to extract from raw information the accurate estimation. One of the common questions in data analysis is the statistical relationship between a response variable (Y) and explanatory variables (Xi). Alexopoulos (2010) continues to explain that an option to answer this question is to employ regression analysis in order to model its relationship. In this project we examine the fundamental principles of multiple linear regression, a widely used statistical method as the author mention. At the same time it is intended to clarify how the method works and when it is appropriate to use.

2 What is Multiple Linear Regression?

The author Uyanık and Güler (2013) explains that regression analysis is used to determine the relationships between two or more variables that exhibit cause–effect patterns, and to make predictions based on those relationships. Similarly, Marill (2004) states that multiple linear regression is a generalization of simple linear regression in which more than one predictor variable is included. The formula is the following:

$$Y_i = b_0 + b_1 X_{1i} + b_2 X_{2i} + \dots + b_n X_{ni} + u_i$$

Where - Y_i — dependent variable

- b_0 — intercept
- $b_1 \dots b_n$ — regression coefficients
- $X_{1i} \dots X_{ni}$ — independent variables
- u_i — disturbance error

3 Usefull aplications

Multiple linear regression is widely applied across numerous fields due to its ability to model relationships between multiple predictors and a continuous outcome. In business and economics for example, it is commonly used for demand forecasting as predicting sales based on factors like price, advertising, etc. As well as for financial risk modeling and market research aimed at identifying the drivers of customer satisfaction or spending. Moreover, social sciences and psychology rely heavily on multiple regression to explore how variables

such as education, income, and demographics influence life satisfaction for example. Another example is to predict academic performance based on study habits and cognitive factors.

In medicine, regression models can help predict patient outcomes from characteristics such as age, BMI, etc. Finally, in machine learning and data science, multiple linear regression serves both as a baseline predictive model for continuous outcomes and as a tool for feature importance analysis, often forming a foundational step before applying more advanced modeling techniques such as regularization or generalized linear models.

4 Assumptions

According to Ernst and Albers (2017) there are 4 assumptions in the linear regression model, violating these assumptions can create several kinds of problems. First, the resulting estimates may be biased, meaning they do not reflect the true population values. Second, the estimators might become inconsistent, so increasing the sample size no longer estimates the true parameters. Third, may lose its efficiency, meaning it can yield less precise estimates. Fourth, hypothesis tests and confidence intervals may become unreliable. Note that these assumptions specifically apply to estimation using the Ordinary Least Squares (OLS) method, which is the default in many statistical software programs.

1. **Linearity:** The model assumes a linear and additive relationship between each predictor and the mean of the dependent variable, with errors averaging zero for any combination of predictors. Problems arise if relationships are nonlinear or if measurement error is present. It is important to mention that multicollinearity is not an assumption, but low overlap between predictors is desirable. (Ernst & Albers, 2017)
2. **Normality:** The model errors are assumed to follow a normal distribution centered at zero. (Ernst & Albers, 2017)
3. **Homoscedasticity:** The errors are expected to have constant variance across all values of the predictors. If the spread of errors changes (heteroscedasticity), the assumption is violated. (Ernst & Albers, 2017)
4. **Independence:** Errors must be independent of each other, meaning their covariances are zero. This requires proper random sampling. Residual plots or autocorrelation checks help detect issues. (Ernst & Albers, 2017)

4.1 Multicollinearity

According to Daoud (2017) multicollinearity occurs when two or more predictor variables are highly correlated, which leads to an increase in the standard errors of the estimated coefficients. When standard errors become inflated, some or all coefficients may no longer appear significantly different from zero, meaning variables that should be statistically meaningful can incorrectly seem insignificant.

5 Linear Model or `lm()` in R

Basic Syntax

```
model <- lm(formula, data = dataset)
```

5.1 Preparing Data and Fitting the Model in R

```
# Load data
data(mtcars)
# Quick overview
summary(mtcars)
```

5.2 Fitting a multiple linear regression with `lm()`

```
# Fit the model
model <- lm(mpg ~ wt + hp + cyl, data = mtcars)
# Basic summary output
summary(model)
```

To understand the output of this code James et al. (2013) offers further explanation:

- Residuals: the difference between what the model predicts and what actually happened.
- Intercept: is the model's prediction of the outcome when all predictors are 0
- Estimated coefficients (Estimate): This tells you how much the outcome changes when a predictor increases by 1 unit, while the others stay the same.
- t-values and p-values: These are used to test if an estimated coefficient is statistically different from zero. The t-value indicates how far the estimate is from zero in terms of standard errors, and the p-value shows the probability of obtaining such a result if the true coefficient were actually zero.

- Residual standard error: This is the average size of prediction errors.
- Multiple R-squared: measures how well the regression model explains the variability in the dependent variable.
- F-statistic: checks if the whole model (all predictors together) is useful.

6 Model Diagnostics in R

6.1 Residual plots and normality

```
# Basic diagnostic plots (4-panel layout)
par(mfrow = c(2, 2))
plot(model)
par(mfrow = c(1, 1))
```

Key plots:

- Residuals vs Fitted: Visualize whether the model is missing any non-linear relationships. When the residuals are evenly scattered around a flat horizontal line with no clear shape or trend, it suggests that the model fits well. (Kim, 2015)
- Normal Q-Q: This plot shows if residuals are normally distributed. Points should follow the diagonal line if residuals are approximately normal. (Kim, 2015)
- Scale-Location: shows if residuals are spread equally along the ranges of predictors. This is how you can check homoscedasticity. It's good if you see a horizontal line with randomly spread points. (Kim, 2015)
- Residuals vs Leverage: it spots influential data points that can strongly affect the regression results. The focus is on points that appear in the upper-right or lower-right areas of the plot or fall outside the dashed lines. (Kim, 2015)

6.2 Checking multicollinearity with VIF

The variance inflation factor (VIF) is a metric that quantifies how much multicollinearity is in the model, values above about 5 are flagged as problematic.

```
# install.packages("car") # run once if not installed
library(car)
vif(model)
```

7 Creating and presenting predictions

7.1 Creating predictions

```
# Predictions for the observed data
mtcars$pred_mpg <- predict(model)

# Create a small new data set for scenario-based predictions
new_cars <- data.frame(
  wt  = c(2.5, 3.0, 3.5),
  hp  = c(100, 150, 200),
  cyl = c(4, 6, 8)
)

predictions<- predict(model, newdata = new_cars)
cbind(new_cars, predictions) #column bind
```

7.2 Presenting prediction

One way to present the predictions is by using the `ggplot2` package, an R library designed for creating statistical and data visualizations. It provides a simple and flexible way to format and display data graphically. Another useful option is the `broom` package, which helps present model results more cleanly by converting output into tidy tibbles.

```
# install.packages("ggplot2") # run once
library(ggplot2)
ggplot(mtcars, aes(x = pred_mpg, y = mpg)) +
  geom_point() +
  geom_abline(intercept = 0, slope = 1, linetype = "dashed") +
  labs(
    x = "Predicted MPG",
    y = "Observed MPG",
    title = "Observed vs Predicted Fuel Efficiency"
  )
```

8 References

Alexopoulos, E. C. (2010). Introduction to multivariate regression analysis. *Hippokratia*, 14(Suppl 1), 23.

- Daoud, J. I. (2017). Multicollinearity and regression analysis. *Journal of Physics: Conference Series*, 949, 012009.
- Ernst, A. F., & Albers, C. J. (2017). Regression assumptions in clinical psychology research practice—a systematic review of common misconceptions. *PeerJ*, 5, e3323.
- James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning: With applications in r*. Springer. <https://doi.org/10.1007/978-1-4614-7138-7>
- Kim, B. (2015). *Understanding diagnostic plots for linear regression analysis*. University of Virginia Library; <https://library.virginia.edu/data/articles/diagnostic-plots>.
- Marill, K. A. (2004). Advanced statistics: Linear regression, part II: Multiple linear regression. *Academic Emergency Medicine*, 11(1), 94–102.
- Uyanık, G. K., & Güler, N. (2013). A study on multiple linear regression analysis. *Procedia-Social and Behavioral Sciences*, 106, 234–240.

Appendix**Affidavit**

I hereby affirm that this submitted paper was authored unaided and solely by me. Additionally, no other sources than those in the reference list were used. Parts of this paper, including tables and figures, that have been taken either verbatim or analogously from other works have in each case been properly cited with regard to their origin and authorship. This paper either in parts or in its entirety, be it in the same or similar form, has not been submitted to any other examination board and has not been published. I acknowledge that the university may use plagiarism detection software to check my thesis. I agree to cooperate with any investigation of suspected plagiarism and to provide any additional information or evidence requested by the university.

Checklist:

- ☒ The handout contains 3-5 pages of text.
- ☒ The submission contains the Quarto file of the handout.
- ☒ The submission contains the Quarto file of the presentation.
- ☒ The submission contains the HTML file of the handout.
- ☒ The submission contains the HTML file of the presentation.
- ☒ The submission contains the PDF file of the handout.
- ☒ The submission contains the PDF file of the presentation.
- ☒ The title page of the presentation and the handout contain personal details (name, email, matriculation number).
- ☒ The handout contains a bibliography, created using BibTeX with an APA citation style.
- ☒ Either the handout or the presentation contains R code that proves the expertise in coding.
- ☒ The filled out Affidavit.
- ☒ The link to the presentation and the handout published on GitHub.

- ☒ In group work, each student's contribution is clearly defined and individual performance can be assessed using specified sections, page numbers or other objective criteria.

Maria Rojas, 10.12.2025, Köln, Germany