

# Understanding Linear (lm) and Generalized Linear (glm) Models in R

Authored by  
**Mohammed loot**

November 5, 2025

## RECOMMENDED CITATION

Mohammed loot (2025). *Understanding Linear (lm) and Generalized Linear (glm) Models in R*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=10561>

The [R programming language](#) serves as the foundational environment for sophisticated statistical computation and data analysis utilized by researchers and data scientists globally. Within R's extensive toolkit, two functions dominate the field of relationship modeling between variables: `lm()` and `glm()`. Although their usage appears superficially similar, mastering the subtle yet profound distinctions between them is absolutely critical for developing statistically sound and appropriate models. Misunderstanding their core assumptions can lead to invalid inferences and flawed conclusions.

This article will provide a detailed, comparative exploration of the syntax, underlying statistical theory, and practical application scope of both functions. We will clearly establish why `lm()` is the specialized function for standard [linear models](#), adhering to strict distributional assumptions, while `glm()` offers the necessary robustness and flexibility to accommodate a vastly wider spectrum of data types and error distributions encountered in complex, real-world datasets.

## The Foundational Role of the `lm()` Function: Ordinary Least Squares

The `lm()` function, standing for **Linear Model**, is purpose-built for fitting classical linear relationships, primarily leveraging the technique of **Ordinary Least Squares (OLS)** estimation. OLS seeks to find the line of best fit that minimizes the sum of the squared vertical distances (residuals) between the observed data points and the predicted values. This method is mathematically elegant and computationally efficient, making `lm()` the default starting point for many regression analyses.

The validity of the results produced by `lm()` rests heavily on a set of stringent statistical assumptions regarding the error term. Crucially, `lm()` assumes that the response variable is **continuous** and that the errors--the differences between observed and predicted outcomes--follow a [Gaussian distribution](#), often referred to as the normal distribution. Furthermore, it assumes homoscedasticity (constant variance of errors) and independence of observations. When these assumptions hold true, `lm()` provides the Best Linear Unbiased Estimator (BLUE), making it the optimal tool for standard [linear regression model](#) analysis.

The structure of the `lm()` function is straightforward, requiring only the definition of the model structure and the data source. Its simplicity reflects its focused statistical purpose.

The standard syntax for the `lm()` function is concise and focused:

```
lm(formula, data, ...)
```

The essential arguments required for its operation are:

**formula:** This argument defines the relationship, specifying the dependent variable followed by the

tilde (~) and the independent variables (e.g., `y ~ x1 + x2`).

**data:** This specifies the name of the [data frame](#) or matrix containing the variables referenced in the formula.

## Introducing `glm()`: The Power of Generalized Linear Models

The `glm()` function represents a fundamental extension of classical regression, enabling the fitting of **Generalized Linear Models (GLMs)**. GLMs are a highly versatile framework that allows modeling response variables whose error distribution belongs to the **exponential family** of distributions, which includes not only the Gaussian but also distributions like Poisson, Binomial, and Gamma. This framework is essential when the data structure violates the normality or homoscedasticity assumptions inherent in `lm()`.

A GLM is defined by three key components: the **Random Component** (specifying the probability distribution of the response variable), the **Systematic Component** (the linear predictor, similar to the right side of an `lm()` formula), and the **Link Function**. This link function is the critical element that relates the expected value of the response variable to the linear predictor. For example, in logistic regression, the link function is the logit, ensuring that the predicted probabilities remain between 0 and 1, a constraint that standard OLS cannot enforce.

This inherent flexibility means that `glm()` is not confined to analyzing continuous, normally distributed data. It is the appropriate tool for response variables that are counts (e.g., using the [poisson](#) family), binary outcomes or proportions (e.g., using the [binomial](#) family), or time-to-event data. This adaptability positions `glm()` as a significantly more versatile and necessary tool for complex, advanced statistical analysis where data types vary widely.

The syntax for the `glm()` function includes a pivotal additional argument--the `family` argument--that dictates the model's distributional assumptions and the associated link function:

```
glm(formula, family=gaussian, data, ...)
```

The arguments for `glm()` include:

**formula:** Specifies the linear relationship between variables.

**family:** This essential parameter defines the statistical distribution used to fit the model and implicitly sets the link function. The default is `gaussian` (which assumes normal errors and an identity link function), making it functionally equivalent to `lm()`. Other common options include [binomial](#) for dichotomous or proportional outcomes, `Gamma`, and [poisson](#) for modeling event counts, among a comprehensive list of possibilities.

`data`: The name of the data frame or object containing the relevant analysis variables.

## The Key Distinction: The Definitive Role of the `family` Argument

The defining difference between these two powerful R functions lies entirely in the presence and utilization of the `family` argument within `glm()`. By explicitly defining the error distribution and the required link function, `glm()` gains the capacity to model data structures that fundamentally violate the core assumptions of classical OLS regression. This adaptability is what allows statisticians to accurately handle non-continuous or constrained response variables.

It is crucial to understand that `lm()` is, statistically speaking, merely a specialized, nested case within the much broader theoretical framework provided by the Generalized Linear Model. When the `glm()` function is invoked without explicitly defining the `family` argument, it automatically defaults to `family=gaussian`. This setting mandates a normal error distribution and an identity link function--the exact specifications implemented by the `lm()` function.

Consequently, if you employ the `lm()` function or explicitly specify `glm(..., family=gaussian)` to fit a standard [linear regression model](#), **the resulting coefficient estimates and statistical inferences will be mathematically identical**. However, despite this equivalence, utilizing `lm()` is generally recommended for purely linear models because it is computationally optimized for OLS calculation, often resulting in slightly faster execution and producing a summary output that is specifically tailored and simpler for interpreting classical linear regression diagnostics (like R-squared).

## Practical Equivalence: Demonstrating Results Consistency

To solidify the conceptual overlap between the two functions when applied to classical linear modeling, we will utilize a practical demonstration. We use R's well-known `mtcars` dataset, aiming to model the relationship for miles per gallon (`mpg`) based on two predictors: engine displacement (`disp`) and horsepower (`hp`). Since `mpg` is a continuous response variable, both functions should yield the same coefficients.

## Example of Using the `lm()` Function

The output generated by `lm()` is highly readable and includes key diagnostic statistics specifically tailored for linear regression, such as the R-squared value and the residual standard error.

```
#fit multiple linear regression model
model <- lm(mpg ~ disp + hp, data=mtcars)
```

```
#view model summary
```

```
summary(model)
```

Call:

```
lm(formula = mpg ~ disp + hp, data = mtcars)
```

Residuals:

```
Min 1Q Median 3Q Max
```

```
-4.7945 -2.3036 -0.8246 1.8582 6.9363
```

Coefficients:

```
Estimate Std. Error t value Pr(>|t|)
```

```
(Intercept) 30.735904 1.331566 23.083 < 2e-16 ***
```

```
disp -0.030346 0.007405 -4.098 0.000306 ***
```

```
hp -0.024840 0.013385 -1.856 0.073679 .
```

```
---
```

```
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Residual standard error: 3.127 on 29 degrees of freedom

Multiple R-squared: 0.7482, Adjusted R-squared: 0.7309

F-statistic: 43.09 on 2 and 29 DF, p-value: 2.062e-09

## Example of Using the `glm()` Function for Linear Regression

By fitting the identical [linear regression model](#) using `glm()` without specifying the `family`, we implicitly rely on the default `family=gaussian` setting, forcing the model into the OLS framework.

```
#fit multiple linear regression model
```

```
model <- glm(mpg ~ disp + hp, data=mtcars)
```

```
#view model summary
```

```
summary(model)
```

Call:

```
glm(formula = mpg ~ disp + hp, data = mtcars)
```

Deviance Residuals:

```
Min 1Q Median 3Q Max
```

```
-4.7945 -2.3036 -0.8246 1.8582 6.9363
```

Coefficients:

```
Estimate Std. Error t value Pr(>|t|)
```

```
(Intercept) 30.735904 1.331566 23.083 < 2e-16 ***
```

```

disp -0.030346 0.007405 -4.098 0.000306 ***
hp -0.024840 0.013385 -1.856 0.073679 .
---
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

(Dispersion parameter for gaussian family taken to be 9.775636)

Null deviance: 1126.05 on 31 degrees of freedom  
 Residual deviance: 283.49 on 29 degrees of freedom  
 AIC: 168.62

Number of Fisher Scoring iterations: 2

A comparison of the output confirms that the coefficient estimates (Estimate column) and their associated standard errors are precisely the same across both functions. The differences observed in the summaries stem from `glm()` including metrics like Deviance Residuals, AIC (Akaike Information Criterion), and the Dispersion parameter. These generalized metrics are standard for all GLMs but are often redundant or interpreted differently in the context of standard OLS, where R-squared and F-statistics are traditionally prioritized.

## Beyond Linearity: Unlocking Complex Modeling with `glm()`

The necessity and true power of the `glm()` function become unequivocally clear when the response variable does not satisfy the assumption of a [Gaussian distribution](#). By correctly specifying a non-Gaussian statistical family, we can accurately model outcomes that are inherently constrained or discrete, such as binary survey responses, counts of events, or categorical classifications. This is where `glm()` moves from equivalence to necessity, providing methodologies that `lm()` cannot handle.

For instance, attempting to model a binary outcome (0 or 1) using `lm()` would violate the normality assumption, potentially leading to predicted values outside the range, and resulting in heteroscedastic residuals. `glm()` overcomes this by utilizing a distribution and link function appropriate for the data type, transforming the linear predictor to match the scale of the expected response.

The most common non-Gaussian applications of `glm()` include:

**Logistic Regression:** Achieved using `family=binomial`. This model employs the logit link function to ensure the predicted outcome is a probability (between 0 and 1), making it indispensable for binary outcome modeling.

**Poisson Regression:** Achieved using `family=poisson`. This model uses the log link function, which is necessary when modeling count data, where outcomes must be non-negative integers.

### Example of Logistic Regression using `family=binomial`

Here, we demonstrate a practical application of `glm()` to model a binary outcome. We model the probability of a car having a manual transmission (`am`: 1 for manual, 0 for automatic) based on engine characteristics. This requires the specification of `family=binomial` to correctly fit a [logistic regression model](#). Notice the output now includes z-values instead of t-values, a standard feature of GLMs relying on maximum likelihood estimation.

#### #fit logistic regression model

```
model <- glm(am ~ disp + hp, data=mtcars, family=binomial)
```

```
#view model summary
```

```
summary(model)
```

Call:

```
glm(formula = am ~ disp + hp, family = binomial, data = mtcars)
```

Deviance Residuals:

Min 1Q Median 3Q Max

```
-1.9665 -0.3090 -0.0017 0.3934 1.3682
```

Coefficients:

Estimate Std. Error z value Pr(>|z|)

```
(Intercept) 1.40342 1.36757 1.026 0.3048
```

```
disp -0.09518 0.04800 -1.983 0.0474 *
```

```
hp 0.12170 0.06777 1.796 0.0725 .
```

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 43.230 on 31 degrees of freedom

Residual deviance: 16.713 on 29 degrees of freedom

AIC: 22.713

Number of Fisher Scoring iterations: 8

### Example of Poisson Regression using `family=poisson`

If our analysis involved modeling count data--such as the number of customer visits or the frequency of manufacturing defects--the appropriate statistical choice would be the Poisson family. Although the `am` variable is binary in this dataset, we demonstrate the syntax below to illustrate how the Poisson family structure is invoked, which applies a logarithmic link function to the linear predictor.

#### #fit Poisson regression model

```
model <- glm(am ~ disp + hp, data=mtcars, family=poisson)
```

```
#view model summary
```

```
summary(model)
```

Call:

```
glm(formula = am ~ disp + hp, family = poisson, data = mtcars)
```

Deviance Residuals:

Min 1Q Median 3Q Max

```
-1.1266 -0.4629 -0.2453 0.1797 1.5428
```

Coefficients:

Estimate Std. Error z value Pr(>|z|)

```
(Intercept) 0.214255 0.593463 0.361 0.71808
```

```
disp -0.018915 0.007072 -2.674 0.00749 **
```

```
hp 0.016522 0.007163 2.307 0.02107 *
```

```
---
```

```
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

(Dispersion parameter for poisson family taken to be 1)

Null deviance: 23.420 on 31 degrees of freedom

Residual deviance: 10.526 on 29 degrees of freedom

AIC: 42.526

Number of Fisher Scoring iterations: 6

### Summary of Choice and Application

Ultimately, the decision regarding which R function to use--`lm()` or `glm()`--is determined by the inherent distributional properties of your response variable and the underlying assumptions you can make about the model errors. If your response variable is continuous, unbounded, and the

errors are safely assumed to follow a [Gaussian distribution](#), then `lm()` is the specialized, computationally efficient choice for classical [linear models](#).

However, if your response variable is discrete (such as counts or binary data), or if the underlying distribution falls outside the normal family (e.g., exponential, Gamma), then `glm()` is the mandatory tool. While `glm(family=gaussian)` provides the capability to replicate `lm()` results, reserving `lm()` for purely OLS regression maintains code clarity and utilizes the function optimized for that specific task.

By mastering the versatility of `glm()`, analysts gain the capacity to accurately model virtually any type of regression problem encountered in serious statistical research and data science, moving beyond the restrictive assumptions of classical linear regression and into the realm of modern statistical modeling.

## Additional Resources

To continue deepening your knowledge of these core modeling techniques in R, we highly recommend consulting the official R documentation, specialized textbooks on generalized linear models, and academic papers detailing the theoretical underpinnings of the exponential family distributions.