

Create a Histogram of Residuals in R

Authored by
Mohammed loot

November 6, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Create a Histogram of Residuals in R*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=11370>

The Critical Role of Residual Normality in Regression Analysis

One of the foundational requirements for employing inferential statistics in many procedures, especially the standard [linear regression model](#) (LRM), is the assumption that the errors or [residuals](#)--the differences calculated between the observed data points and the values predicted by the model--are independently and identically distributed following a [Normal Distribution](#) (often called the assumption of normality). This assumption is crucial because violations can severely compromise the accuracy of subsequent statistical inferences, leading to incorrect calculations of **confidence intervals** and yielding unreliable results from hypothesis tests. While the parameter estimates (the beta coefficients) often remain unbiased even when normality is violated, the ability to correctly assess the statistical significance of those parameters is jeopardized.

A highly effective and accessible method for evaluating this critical assumption is the visual construction of a [histogram](#) specifically for the calculated [residuals](#). We are looking for the resulting frequency distribution to closely mimic the classic, symmetrical "bell-shape" that characterizes the [Normal Distribution](#). This visual diagnostic provides an immediate, intuitive check on the model's adherence to one of its core statistical requirements, complementing other diagnostics like the Q-Q plot.

This comprehensive tutorial provides a rigorous, step-by-step methodology for calculating and visualizing the [residuals](#) derived from a regression analysis. We will utilize the robust [R programming environment](#) and leverage the extensive capabilities of the popular **ggplot2** package to generate high-quality, informative residual histograms.

Step 1: Setting Up the Environment and Generating Synthetic Data

The prerequisite for fitting any statistical model is a well-defined dataset. To ensure that our steps are fully transparent and reproducible for any user, we will begin by generating synthetic data within the [R programming environment](#). This practice is standard in statistical tutorials, allowing us to demonstrate the entire workflow without relying on external files. A cornerstone of generating random data in R is the use of the `set.seed()` function; setting a specific seed guarantees that the sequence of generated random numbers remains identical every time the code is executed, thus ensuring **reproducibility**.

In this demonstration, we construct a simple dataset comprising 100 observations. We define two synthetic predictor variables, x_1 and x_2 , and one response variable, y . For simplicity and to satisfy the assumptions of our eventual model, all variables are intentionally drawn from distinct [Normal Distributions](#) characterized by specific means and standard deviations. This process culminates in a structured R `data.frame`, which is ideally suited for fitting a multiple regression analysis.

The code below executes the data generation process, ensuring that the `data` object is ready for

model fitting. We conclude this step by inspecting the initial rows of the generated data frame using `head(data)` to confirm the expected structure and variable types.

Ensure the example is fully reproducible by setting a fixed seed value

set.seed(0)

```
# Create synthetic data with 100 observations
```

```
x1 <- rnorm(n=100, 2, 1)
```

```
x2 <- rnorm(100, 4, 3)
```

```
y <- rnorm(100, 2, 3)
```

```
data <- data.frame(x1, x2, y)
```

```
# View the first six rows of the resulting data frame
```

```
head(data)
```

```
x1 x2 y
```

```
1 3.262954 6.3455776 -1.1371530
```

```
2 1.673767 1.6696701 -0.6886338
```

```
3 3.329799 2.1520303 5.8081615
```

```
4 3.272429 4.1397409 3.7815228
```

```
5 2.414641 0.6088427 4.3269030
```

```
6 0.460050 5.7301563 6.6721111
```

Step 2: Executing the Multiple Linear Regression

Once the data preparation is complete, the next logical step is to fit the chosen statistical model. For this analysis, we employ the powerful and ubiquitous `lm()` function native to R, which is the standard implementation for fitting ordinary least squares [linear regression models](#). We specify the model structure by predicting the response variable y as a function of the two predictor variables, x_1 and x_2 .

The execution of `lm(y ~ x1 + x2, data=data)` returns a model object, which we store in the variable `model`. This object is comprehensive, containing all necessary statistical components, including coefficient estimates, R-squared values, and, most pertinent to our current task, the automatically calculated vector of [residuals](#). These residuals represent the raw, unstandardized errors of prediction for every single observation in the dataset.

This single line of code efficiently performs the complex mathematical calculation required to minimize the sum of squared errors, thereby generating the optimal regression line and providing us with the necessary error data for subsequent visualization.

```
# Fit the multiple linear regression model using predictors x1 and x2
model <- lm(y ~ x1 + x2, data=data)
```

Step 3: Visualizing Residuals with ggplot2

To move beyond raw numerical outputs and visualize the distribution of errors, we rely on the industry-standard [ggplot2](#) package. Known for its declarative syntax and capacity to produce aesthetically pleasing and highly customizable statistical graphics, **ggplot2** is the preferred tool within the [R programming environment](#) for this task. Before generating the plot, we must ensure the package is loaded into the current R session using the `library(ggplot2)` command.

The core of the plotting command involves two primary actions. First, we access the vector of residuals directly from our fitted model object using the expression `model$residuals`. Second, we map this vector to the x-axis aesthetic within the `ggplot()` function. The frequency visualization itself is then created using the `geom_histogram()` layer. We enhance the visual quality by specifying standard colors--setting the fill to 'steelblue' and defining the bar boundaries with a black outline--and adding clear labels for the title, x-axis (Residuals), and y-axis (Frequency) using the `labs()` function.

The initial [histogram](#) generated by this code block uses the default bin width determined by **ggplot2**, providing our first visual assessment of residual normality.

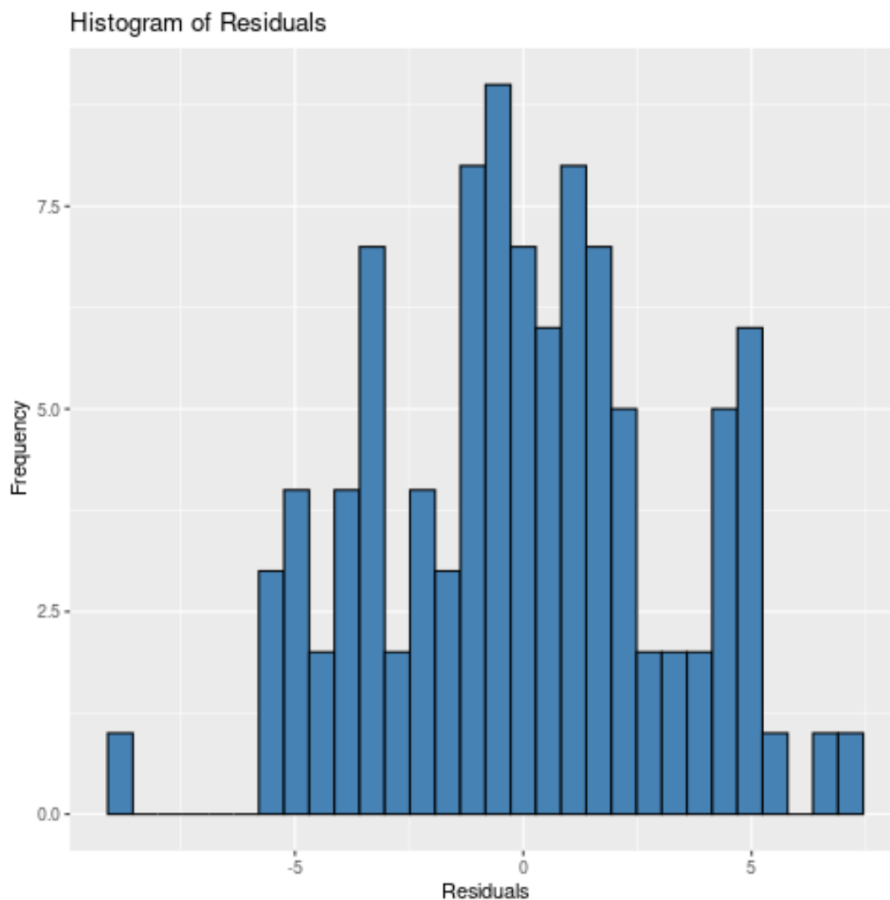
Load the ggplot2 visualization package

```
library(ggplot2)
```

```
# Create the histogram of residuals using default binning
```

```
ggplot(data = data, aes(x = model$residuals)) +
  geom_histogram(fill = 'steelblue', color = 'black') +
  labs(title = 'Histogram of Residuals', x = 'Residuals', y = 'Frequency')
```

This visualization immediately communicates the density and dispersion of the prediction errors, allowing us to judge how closely the distribution approximates the expected bell shape of a [Normal Distribution](#).



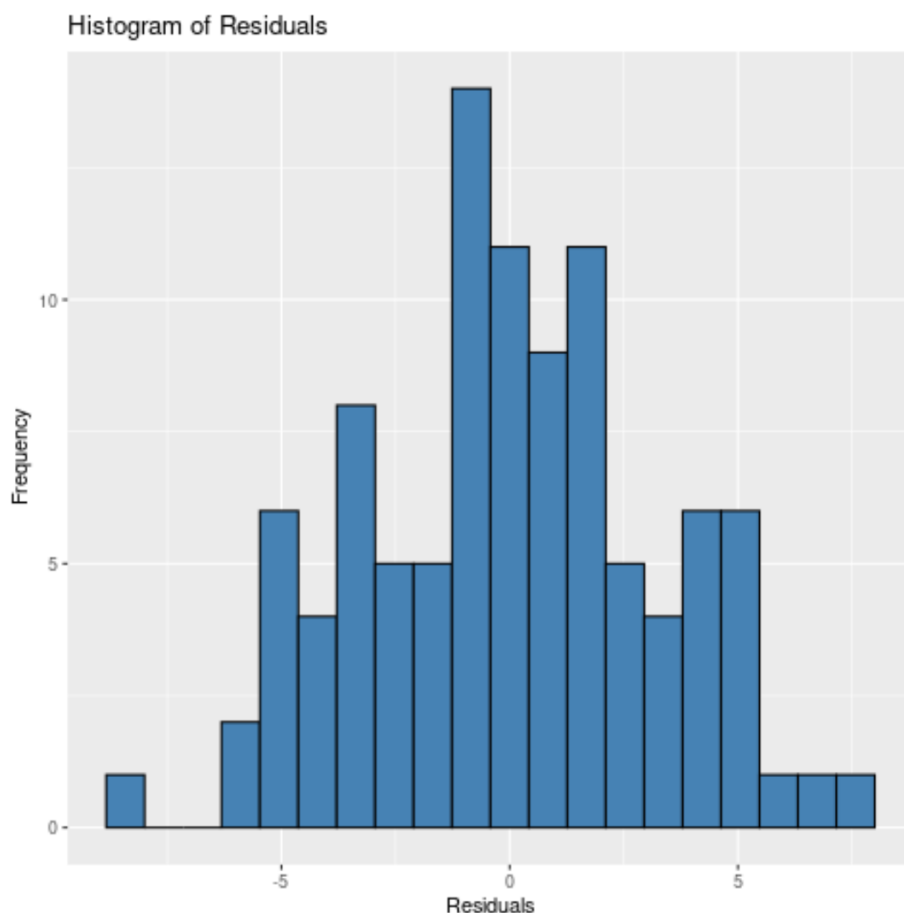
Step 4: Fine-Tuning Granularity by Adjusting Bins

A fundamental decision when constructing a [histogram](#) is determining the appropriate number of bins. Bins define the specific intervals into which the continuous data (in this case, the residuals) is categorized, and the height of the resulting bar signifies the count or frequency of data points falling within that range. Although **ggplot2** applies a statistically sound default bin count, analysts frequently need to manually adjust this parameter to obtain the clearest possible portrayal of the underlying distribution's true structure.

The choice of bin count involves a critical trade-off between detail and generalization. Employing a larger number of bins results in narrower intervals, which offers a highly detailed, high-resolution view of the data density. However, this level of detail can sometimes make the plot appear jagged or noisy, obscuring the overall pattern. Conversely, selecting fewer bins creates wider bars, yielding a smoother, more generalized representation of the distribution shape. If too few bins are chosen, crucial nuances or multimodal features within the residual distribution might be masked. We control this essential plotting parameter using the `bins` argument within the `geom_histogram()` function.

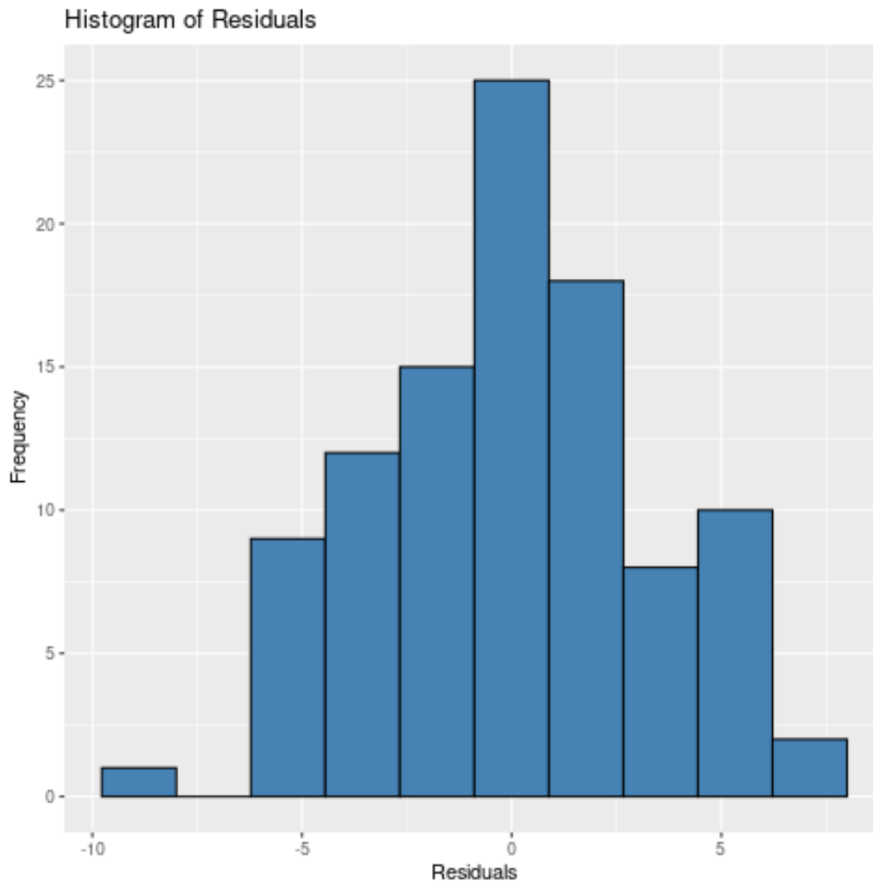
To illustrate the effect of increased granularity, we explicitly set the number of bins to **20**. This modification provides a much finer look at the distribution of the [residuals](#), often revealing subtle peaks or asymmetries that might be averaged out by the default setting:

```
# Create histogram of residuals, explicitly setting the number of bins to 20  
ggplot(data = data, aes(x = model$residuals)) +  
geom_histogram(bins = 20, fill = 'steelblue', color = 'black') +  
labs(title = 'Histogram of Residuals (20 Bins)', x = 'Residuals', y = 'Frequency')
```



Conversely, decreasing the bin count to **10** demonstrates the effect of reducing granularity. Notice how the resulting plot appears significantly coarser; the distribution shape is simplified, which might obscure fine details but clearly highlights the central tendency:

```
# Create histogram of residuals, reducing the number of bins to 10  
ggplot(data = data, aes(x = model$residuals)) +  
geom_histogram(bins = 10, fill = 'steelblue', color = 'black') +  
labs(title = 'Histogram of Residuals (10 Bins)', x = 'Residuals', y = 'Frequency')
```



Interpreting the Visual Evidence and Complementary Formal Tests

Based on the visual evidence presented across the different bin configurations, the [residuals](#) derived from our simulated [linear regression model](#) strongly suggest a rough adherence to the assumption of [normal distribution](#). Key indicators supporting this conclusion include the distinct symmetry of the shape, a singular, prominent peak centered near zero, and tails that taper off relatively evenly on both sides. This favorable outcome validates the decision to proceed with standard inferential statistical procedures for analyzing the model's coefficients.

While visual inspection via the [histogram](#) is highly intuitive and effective, statistical analysts often employ complementary formal statistical tests to quantify potential deviations from normality. The most commonly used tests available within the [R programming environment](#) include:

The [Shapiro-Wilk test](#): Excellent power for detecting non-normality, especially in smaller to medium sample sizes.

The Kolmogorov-Smirnov (Lilliefors) test: Used to compare the residual distribution against a theoretical normal distribution.

The Jarque-Bera test: Measures skewness and kurtosis relative to a normal distribution.

These formal tests provide a quantitative measure--a p-value--indicating the probability that the observed distribution significantly differs from a perfectly normal one. However, it is essential for practitioners to interpret these results cautiously, particularly in the context of large datasets.

A significant limitation of formal normality tests is their extreme sensitivity to large sample sizes. When the sample size is substantial (N is large), these tests frequently reject the null hypothesis of normality even when the deviation is statistically minute and holds no practical consequence for the stability or interpretation of the [linear regression model](#). Consequently, the consensus among robust statistical practice is to rely on a combined diagnostic approach. This involves prioritizing visual tools--specifically the [residual histogram](#) and the Q-Q plot--as the primary evidence for assessing the normality assumption, complementing them with formal tests only when necessary and interpreting the results within the context of the sample size.