

How to Test for Heteroscedasticity in R Using the Breusch-Pagan Test

Authored by
stats writer

March 11, 2026

RECOMMENDED CITATION

stats writer (2026). *How to Test for Heteroscedasticity in R Using the Breusch-Pagan Test*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=135103>

Understanding the Fundamentals of the Breusch-Pagan Test

The **Breusch-Pagan test** represents a vital diagnostic tool in the field of econometrics and statistical modeling, specifically designed to identify the presence of **heteroscedasticity** within a **linear regression** framework. In the context of **Ordinary Least Squares** (OLS) regression, one of the fundamental Gauss-Markov assumptions is **homoscedasticity**, which posits that the variance of the **residuals** remains constant across all levels of the independent variables. When this assumption is violated, the resulting **standard errors** of the regression coefficients can become biased, potentially leading to incorrect inferences regarding the statistical significance of the predictors.

The primary utility of the **Breusch-Pagan test** is its ability to determine whether the estimated variance of the **residuals** from a **regression analysis** is dependent on the values of the explanatory variables. By analyzing these relationships, researchers can ascertain if the model's predictive accuracy is consistent across the entire dataset or if the error terms exhibit a systematic pattern. This is particularly crucial in financial modeling and social sciences, where data often display increasing or decreasing variability over time or across different demographic strata.

To implement this test effectively in the **R programming language**, practitioners typically rely on specialized libraries that automate the complex matrix calculations involved. The most prominent of these is the **lmtest** package, which provides the **bptest()** function. This function streamlines the process of evaluating the **null hypothesis**, which states that the variance of the error terms is constant. By generating a **chi-squared** test statistic and a corresponding **p-value**, the test offers a rigorous mathematical basis for accepting or rejecting the assumption of constant variance.

The Statistical Implications of Heteroscedasticity

Identifying **heteroscedasticity** is not merely a theoretical exercise; it has profound practical implications for the reliability of a **linear model**. When the variance of the **residuals** is non-constant, the OLS estimators remain unbiased but are no longer the Best Linear Unbiased Estimators (BLUE). This loss of efficiency means that the confidence intervals generated by the model may be unnecessarily wide or misleadingly narrow, which can lead to Type I or Type II errors during **hypothesis testing**.

The **Breusch-Pagan test** addresses this issue by regressing the squared **residuals** from the original model onto the original independent variables. If the independent variables can significantly explain the variation in the squared residuals, it suggests that the variance is indeed a function of those predictors. This secondary **regression analysis** serves as the foundation for the test statistic, allowing for a quantified assessment of how the spread of the data points changes in relation to the model's inputs.

In modern data science, ensuring **homoscedasticity** is a prerequisite for many advanced predictive modeling techniques. While some robust algorithms can tolerate minor deviations, traditional **linear regression** requires careful validation of its underlying assumptions. Therefore, the **Breusch-Pagan test** acts as a gateway, ensuring that the insights derived from a model are grounded in statistically sound principles rather than artifacts of irregular data distribution.

Preparing the R Environment for Regression Diagnostics

Before conducting a **Breusch-Pagan test**, it is essential to prepare the **R** environment by loading the necessary datasets and computational libraries. The **mtcars** dataset, which is built into the base **R** installation, provides a classic platform for demonstrating these techniques. It contains performance data for various automobile models, making it an excellent candidate for exploring the relationships between fuel efficiency and engine characteristics through **regression analysis**.

The primary library required for the diagnostic portion of this tutorial is **lmtest**. This package is specifically curated for testing linear regression models and includes a wide array of functions for detecting autocorrelation, **heteroscedasticity**, and functional form specification. By integrating these tools into your workflow, you can move beyond simple model fitting and perform comprehensive diagnostic checks that are standard in professional **data analysis**.

In addition to **lmtest**, researchers often utilize visualization libraries like **ggplot2** to complement the numerical results of the **Breusch-Pagan test**. Plotting the **residuals** against the fitted values can provide a visual confirmation of the test's findings. A clear "fan" or "butterfly" shape in the residual plot is a classic visual indicator that **heteroscedasticity** is present, necessitating further investigation via the formal test statistic.

Step 1: Fitting the Multivariate Regression Model

The first practical step in our analysis involves fitting a **linear regression** model using the **lm()** function in **R**. In this specific example, we will treat "mpg" (miles per gallon) as the **dependent variable**, while "disp" (displacement) and "hp" (horsepower) will serve as our explanatory variables. This configuration allows us to examine how the physical attributes of an engine influence its fuel economy.

```
#load the dataset
```

```
data(mtcars)
```

```
#fit a regression model
```

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

```
#view model summary
```

```
summary(model)
```

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

Upon reviewing the **model summary**, we can observe the **coefficients**, standard errors, and the overall R-squared value. The R-squared value of approximately 0.7482 suggests that our model explains a significant portion of the variance in fuel efficiency. However, the validity of these results depends entirely on whether the underlying assumptions of the **linear model** are met, which is why the **Breusch-Pagan test** is our next priority.

It is important to note that the **standard errors** reported in this summary are calculated under the assumption of **homoscedasticity**. If our subsequent test reveals **heteroscedasticity**, these specific standard error values--and the t-statistics derived from them--could be inaccurate. This potential for inaccuracy highlights why diagnostics are an inseparable part of any robust **statistical analysis** pipeline.

Step 2: Executing the Breusch-Pagan Test in R

Once the model has been fitted and its initial parameters reviewed, we can proceed to the execution of the **Breusch-Pagan test**. This requires calling the **lmtest** library and applying the **bptest()** function directly to our model object. This function evaluates the relationship between the **residuals** and the predictors to check for non-constant variance.

```
#load lmtest library
```

```
library(lmtest)
```

```
#perform Breusch-Pagan Test
```

```
bptest(model)
```

```
studentized Breusch-Pagan test
```

```
data: model
```

```
BP = 4.0861, df = 2, p-value = 0.1296
```

The output of the **bptest()** function provides several key metrics. The **BP statistic** (4.0861) represents the calculated value used to determine the likelihood of the observed data under the **null hypothesis**. The **degrees of freedom** (df = 2) correspond to the number of explanatory variables in the model. Finally, the **p-value** (0.1296) is the critical value used to make a final determination about the presence of **heteroscedasticity**.

The "studentized" version of the test, which is the default in the **lmtest** package, is generally preferred as it is more robust to departures from **normality** in the error terms. By using this version, we ensure that our diagnostic is reliable even if the residuals are not perfectly normally distributed, which is a common occurrence in real-world datasets like **mtcars**.

Interpreting the Test Results and Statistical Significance

The interpretation of the **Breusch-Pagan test** revolves around the **p-value**. In standard statistical practice, a **significance level** (alpha) of 0.05 is typically used as the threshold for decision-making. If the p-value is less than 0.05, we reject the **null hypothesis** of **homoscedasticity**, indicating that **heteroscedasticity** is present in the model.

In our specific example, the p-value is **0.1296**. Since this value is greater than the 0.05 threshold, we fail to reject the **null hypothesis**. This result suggests that we do not have sufficient evidence to conclude that the variance of the **residuals** is non-constant. Consequently, we can proceed with the interpretation of our original **regression analysis** with a reasonable degree of confidence that our standard errors are not compromised by unequal variance.

It is important to remember that "failing to reject" the null hypothesis is not the same as "proving" that the variance is perfectly constant. It simply means that the data does not provide strong enough evidence to suggest otherwise. In many professional contexts, it is still advisable to visualize the **residuals** to ensure that no specific patterns (such as non-linearity) are being overlooked by the global test statistic.

Advanced Remediation: What to Do if Heteroscedasticity is Detected

In scenarios where the **Breusch-Pagan test** yields a significant p-value (e.g., $p < 0.05$), researchers must take corrective action to ensure the validity of their conclusions. One of the most common approaches is to use **robust standard errors**, often referred to as Huber-White standard errors. These adjusted errors provide more accurate **significance** tests without requiring changes to the regression coefficients themselves.

Another effective strategy involves re-evaluating the functional form of the model. Sometimes, **heteroscedasticity** is a symptom of a missing interaction term or a non-linear relationship that hasn't been properly specified. By refining the model to include these complexities, the pattern in the **residuals** may disappear, restoring **homoscedasticity** and improving the overall explanatory power of the analysis.

If the model is correctly specified but the variance remains non-constant, more intensive mathematical interventions may be required. These include variable transformations and specialized regression techniques that account for the variance structure directly. The goal of these interventions is to "stabilize" the variance, ensuring that every observation contributes appropriately to the estimation of the **regression** parameters.

Variable Transformation Techniques for Variance Stabilization

Transforming the response variable is a highly effective way to mitigate **heteroscedasticity**. The most common transformation is the **logarithm** transformation. By taking the natural log of the dependent variable, we compress the scale of the data, which often reduces the impact of outliers and stabilizes the variance across different levels of the predictors. This is particularly useful for variables that exhibit exponential growth or have a heavily skewed distribution.

Other common transformations include the square root transformation and the inverse transformation. The choice of transformation often depends on the nature of the data and the specific pattern of the **residuals**. For example, if the variance is proportional to the mean, a square root transformation might be appropriate. These techniques fall under the broader umbrella of **Power Transformations**, such as the Box-Cox transformation, which can automatically determine the optimal mathematical shift to achieve **homoscedasticity**.

While transformations can solve the statistical issue of non-constant variance, they also change the interpretation of the **regression coefficients**. For instance, in a log-linear model, the coefficients represent the percentage change in the dependent variable for a unit change in the independent variable. Researchers must be prepared to adjust their narrative and reporting to reflect these transformed relationships accurately.

Implementing Weighted Least Squares (WLS) Regression

When simple transformations are insufficient, **Weighted Least Squares** (WLS) regression offers a sophisticated alternative to standard OLS. In WLS, each observation is assigned a weight that is inversely proportional to its variance. This means that observations with higher variance (which are less "reliable") are given less influence in determining the **regression** line, while observations with lower variance are given more weight.

Implementing WLS in **R** involves calculating the weights based on the **residuals** of an initial OLS model. Typically, the weights are defined as the reciprocal of the squared residuals or a function of the independent variables. By applying these weights within the **lm()** function using the **weights** argument, the **linear regression** model is re-estimated to minimize the weighted sum of squares.

The use of **Weighted Least Squares** effectively eliminates the problems caused by **heteroscedasticity**, resulting in estimators that are once again BLUE. However, WLS requires the researcher to know or accurately estimate the structure of the variance, which can be challenging in complex datasets. Despite this, it remains a powerful tool in the statistician's arsenal for ensuring that **regression models** remain robust and accurate in the face of non-constant variance.

ARABPSYCHOLOGY.COM