

# How to Perform a Box-Cox Transformation in SAS

Authored by  
**stats writer**

November 19, 2025

## RECOMMENDED CITATION

stats writer (2025). *How to Perform a Box-Cox Transformation in SAS*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=97019>

The Box-Cox Transformation is a powerful technique utilized in statistical modeling, particularly within the realm of linear regression, designed to address violations of key model assumptions. Primarily, this transformation serves to stabilize the variance (achieving homoscedasticity) and render the data distribution closer to a standard normal distribution. When the response variable in a model displays severe skewness or heteroscedasticity, applying the appropriate Box-Cox transformation can significantly improve the reliability and interpretability of the resulting statistical inferences. It is fundamentally a family of power transformations parameterized by a value known as lambda ( $\lambda$ ).

The core philosophy behind the Box-Cox method involves estimating the optimal value of the parameter  $\lambda$  (lambda) that maximizes the log-likelihood function for the transformed data, effectively making the resulting dataset approximate a Gaussian shape as closely as possible. This estimation process is often iterative and data-driven. By transforming the dependent variable, we aim to satisfy the underlying assumptions required for robust modeling techniques, such as the assumption that the errors are independent and identically distributed, following a normal distribution with a mean of zero and constant variance. Failing to meet these criteria can lead to biased coefficient estimates and incorrect standard errors, undermining the validity of hypothesis testing.

## Understanding the Box-Cox Transformation

A Box-Cox transformation is a widely accepted statistical method for transforming non-normally distributed datasets into a more manageable, normally distributed form. This transformation is particularly valuable when the variability of the data changes systematically across the range of values, a condition known as heteroscedasticity. By applying a power transformation, the method attempts to stabilize the variance and mitigate the effects of skewness, thereby satisfying the classical assumptions of linear models.

The implementation of the Box-Cox technique involves finding an optimal transformation parameter,  $\lambda$ , which is estimated from the data itself. This estimated  $\lambda$  determines the specific power to which the variable is raised. Unlike arbitrary transformations (such as log or square root), the Box-Cox method provides a systematic, mathematically derived solution that maximizes the likelihood of achieving normality for the transformed variable. In the SAS environment, this complex estimation process is handled efficiently through specialized statistical procedures.

## Why Data Transformation is Necessary

Many widely used statistical procedures, especially those based on the Ordinary Least Squares (OLS) method for regression, rely heavily on the assumption that the error terms (residuals) are normally distributed and exhibit constant variance across all levels of the predictor variables. When

raw data deviates significantly from this ideal, typically manifesting as heavy positive or negative skewness, the standard OLS framework becomes inefficient and potentially misleading. A visual inspection of the data, such as histograms or Q-Q plots of the residuals, often reveals these issues, signaling the need for corrective measures.

While simpler transformations like the natural logarithm or the square root transformation are commonly applied to skewed data, they represent fixed choices ( $\lambda=0$  or  $\lambda=0.5$ , respectively) and may not be the optimal fit for a specific dataset. The power of the Box-Cox transformation lies in its generality; it searches across a continuum of potential power values ( $\lambda$ ) to find the single transformation that best stabilizes the variance and induces normality. This flexibility makes it a preferred tool for statisticians seeking the most appropriate data preprocessing technique before fitting complex models.

## The Mathematical Foundation: Box-Cox Formula

The Box-Cox transformation defines a family of power transformations, which are mathematically expressed using a piecewise function dependent on the value of  $\lambda$ . This parameter  $\lambda$  dictates the shape and extent of the transformation applied to the variable, denoted as  $y$ . Understanding this formula is crucial for correctly interpreting and applying the results derived from the transformation procedure.

The generalized formula for transforming a variable  $y$  (where  $y$  must be strictly positive) is defined as follows:

$$y(\lambda) = (y^\lambda - 1) / \lambda, \text{ if } \lambda \neq 0$$

$$y(\lambda) = \log(y), \text{ if } \lambda = 0$$

The second case, where  $\lambda$  equals zero, represents the natural logarithm transformation. Although the formula appears discontinuous at  $\lambda=0$ , it can be shown through L'Hôpital's Rule that the limit of the first expression as  $\lambda$  approaches zero is indeed  $\log(y)$ . This mathematical elegance ensures that the transformation is smooth across all values of  $\lambda$ , providing a continuous search space for the optimal parameter during the estimation phase.

## Identifying the Optimal Lambda ( $\lambda$ ) in SAS

In the SAS statistical software environment, the most straightforward and reliable method for determining the optimal  $\lambda$  parameter is by utilizing the **PROC TRANSREG** procedure. PROC TRANSREG stands for Transformation and Regression and is specifically designed to handle complex data transformations, including optimal scaling and power transformations like Box-Cox. This procedure iteratively searches through the possible values of  $\lambda$ , typically between -2 and 2, identifying the value that maximizes the goodness-of-fit statistic, usually associated with

maximizing the likelihood function of the transformed data being normally distributed.

The strength of **PROC TRANSREG** lies in its ability to integrate the transformation step directly within a regression context. By specifying the desired transformation type (e.g., `BOXCOX(Y)`) within the model statement, SAS automatically executes the complex iterative process required to estimate  $\lambda$ . The output then provides the precise estimate for  $\lambda$ , which is the value we use to manually calculate the new, transformed response variable for subsequent analyses using procedures like **PROC REG**.

## Step-by-Step Example: Initial Data Setup

To illustrate the practical application of the Box-Cox transformation, we will use a small dataset in SAS. This example walks through the process of diagnosing non-normality, identifying the optimal transformation parameter, and applying that parameter to improve the regression model's compliance with core statistical assumptions. The following code creates and displays our sample dataset, `my_data`, which includes the predictor variable `x` and the response variable `y`.

```
/*create dataset*/
```

```
data my_data;
```

```
input x y;
```

```
datalines;
```

```
7 1
```

```
7 1
```

```
8 1
```

```
3 2
```

```
2 2
```

```
4 2
```

```
4 2
```

```
6 2
```

```
6 2
```

```
7 3
```

```
5 3
```

```
3 3
```

```
3 6
```

```
5 7
```

```
8 8
```

```
;
```

```
run;
```

```
/*view dataset*/
```

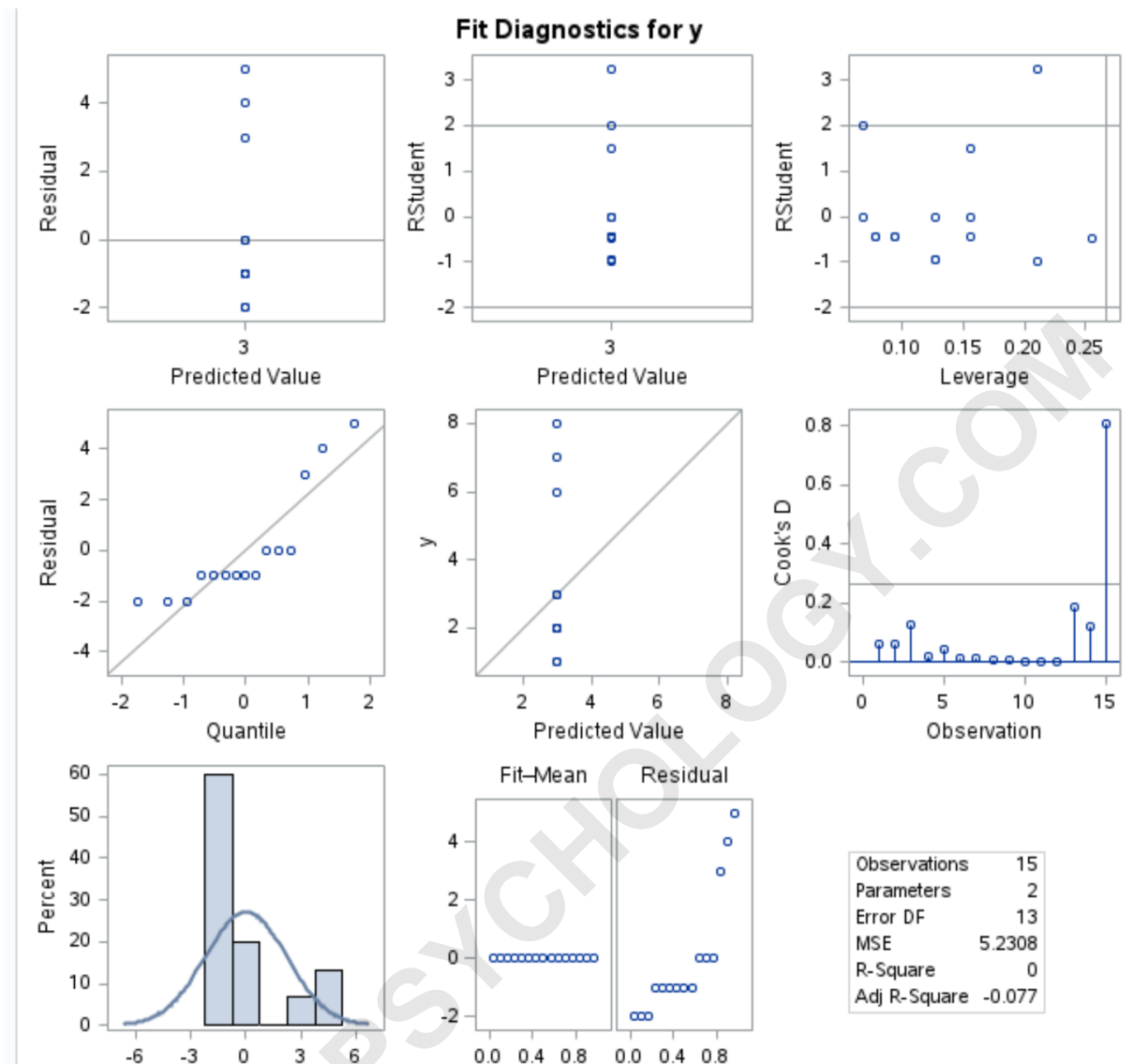
```
proc print data=my_data;
```

## Diagnosing Normality Before Transformation

Before implementing any transformation, we must first confirm that the initial model exhibits issues related to non-normality. We fit a simple linear regression model using **PROC REG**, using **x** as the predictor variable and **y** as the response variable. This procedure generates critical diagnostic plots essential for evaluating model assumptions.

```
/*fit simple linear regression model*/  
proc reg data=my_data;  
model y = x;  
run;
```

The primary diagnostic tool for assessing the normality of the error terms is the **Residual vs. Quantile plot**, often referred to as the Q-Q plot. This plot visualizes the standardized residuals against the theoretical quantiles of the normal distribution. If the residuals follow a normal distribution, the data points should align closely with the straight diagonal line running through the plot.



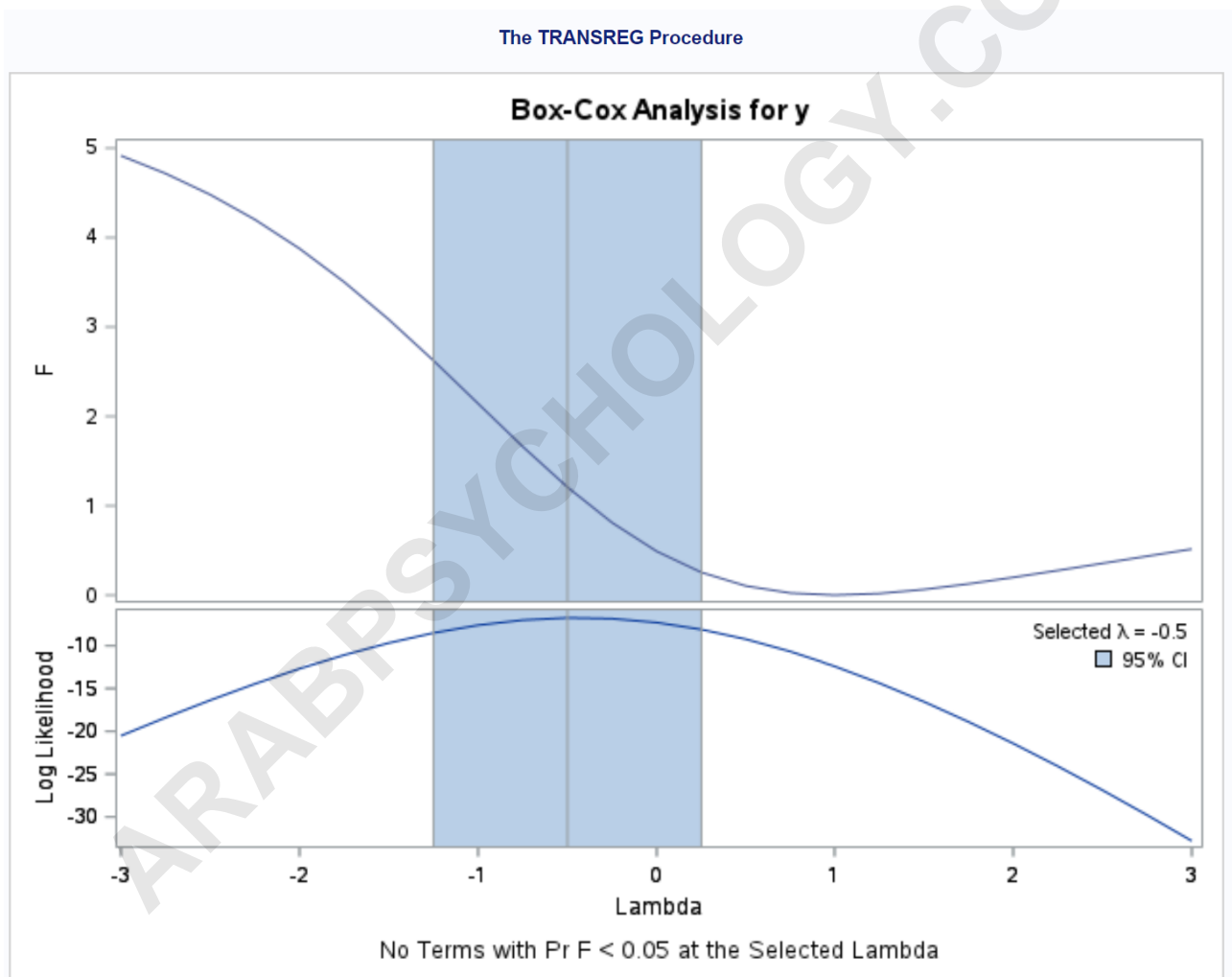
Upon reviewing the initial plot above, it is evident that the plotted points deviate significantly from the straight diagonal reference line, particularly at the tails. This severe non-linearity is a strong indication that the residuals of the original linear model are not normally distributed, suggesting that the response variable  $y$  itself is likely skewed. This confirms the necessity of applying a transformation, such as the Box-Cox method, to improve the model's adherence to assumptions.

## Executing the Transformation using PROC TRANSREG

Having established the need for transformation, we now invoke **PROC TRANSREG** to estimate the optimal value of  $\lambda$ . The syntax is straightforward: we specify the data set and use the `MODEL` statement, where the response variable  $\bar{y}$  is encased within the `BOXCOX()` function, and the predictor  $\bar{x}$  is specified using the `IDENTITY()` function, indicating that  $\bar{x}$  should not be transformed.

```
/*perform box-cox transformation*/  
proc transreg data=my_data;  
model boxcox(y) = identity(x);  
run;
```

The output generated by **PROC TRANSREG** provides the estimated optimal value for  $\lambda$ , which is determined by maximizing the likelihood of the transformed variable being normally distributed. This estimation step is the mathematical core of the Box-Cox procedure, allowing the data to dictate the best possible transformation power rather than relying on an arbitrary choice.



Based on the calculations performed by the procedure, the selected value to use for  $\lambda$  is exactly **-0.5**. This suggests that the reciprocal square root transformation is the mathematically optimal power to apply to the response variable  $y$  to achieve the desired level of normality and variance stability for subsequent regression analysis.

## Applying the New Transformed Variable and Re-running Regression

Once the optimal  $\lambda$  is identified ( $\lambda = -0.5$ ), we must manually apply the transformation formula to the original response variable  $y$  to create a new, transformed variable, which we call `new_y`. Using the Box-Cox formula where  $\lambda \neq 0$ , the transformation becomes:  $new\_y = (y^{-0.5} - 1) / -0.5$ . We incorporate this calculation within a new SAS data step to create a clean dataset ready for the final model fit.

```
/*create new dataset that uses box-cox transformation to create new y*/
```

```
data new_data;
```

```
set my_data;
```

```
new_y = (y**(-0.5) - 1) / -0.5;
```

```
run;
```

```
/*fit simple linear regression model using new response variable*/
```

```
proc reg data=new_data;
```

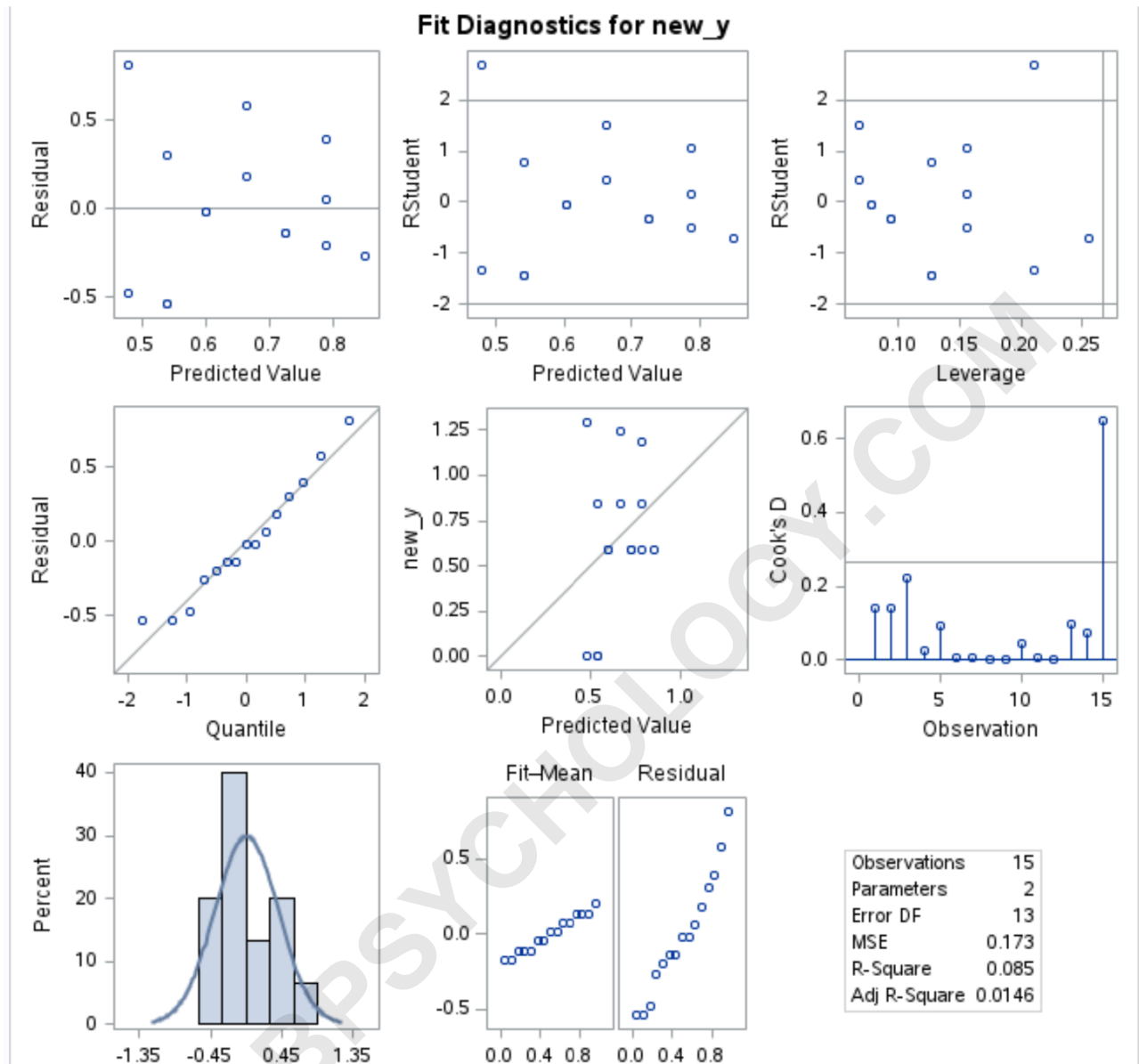
```
model new_y = x;
```

```
run;
```

We then use **PROC REG** again, this time replacing the original response variable  $y$  with the newly calculated `new_y`. This step fits the linear model to the transformed data, allowing us to evaluate whether the Box-Cox procedure successfully mitigated the non-normality issues identified earlier. The success of the transformation hinges entirely on the diagnostic plots generated by this second execution of **PROC REG**.

## Interpreting the Results Post-Transformation

A final examination of the diagnostic plots from the model utilizing `new_y` is essential. We specifically focus on the **Residual vs. Quantile plot** (Q-Q plot) to assess the new distribution of the residuals. A successful Box-Cox transformation should result in data points that now closely hug the straight diagonal line, confirming that the normality assumption is satisfied.



As clearly demonstrated by the updated diagnostic plot, the residuals now align much more closely with the theoretical normal distribution line compared to the initial plot shown in the diagnostics section. This remarkable improvement indicates that the transformation successfully corrected the heavy skewness present in the original response variable. Consequently, the new regression model ( $\text{new\_y} = x$ ) is far more robust, satisfying one of the main prerequisites for valid statistical inference in linear regression. Satisfying these core assumptions, such as the normality of errors and constant variance, ensures that the p-values and confidence intervals derived from the model are statistically reliable, providing a stronger foundation for drawing conclusions from the data.