

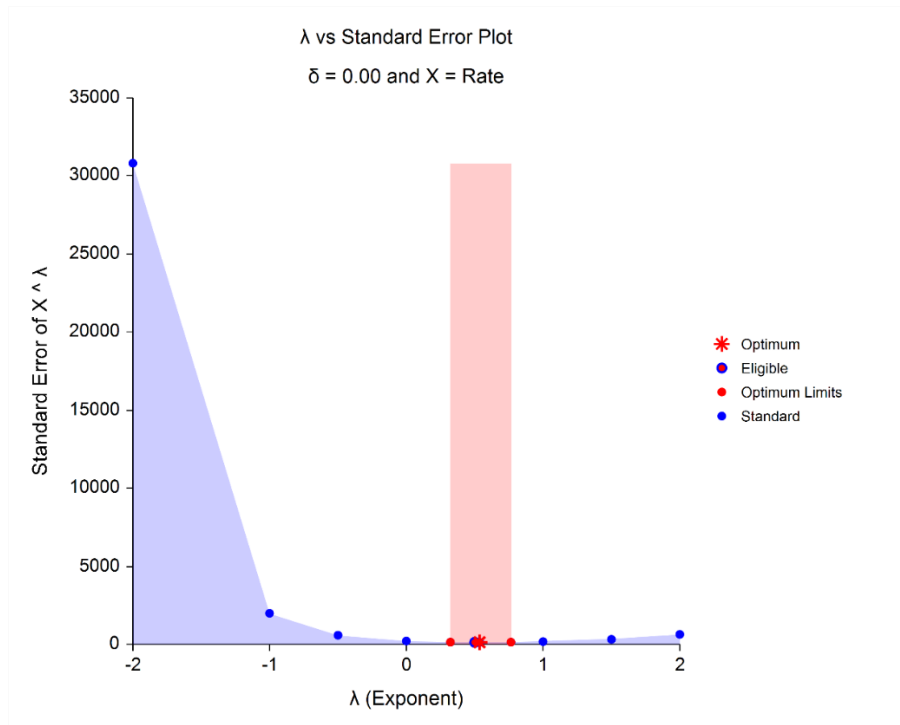
Chapter 191

Box-Cox Transformation for Two or More Groups (T-Test and One-Way ANOVA)

Introduction

This procedure finds the appropriate Box-Cox power transformation (1964) for a dataset containing a response value divided among two or more groups. The data will eventually be analyzed by a two-sample t-test (two groups) or a one-way ANOVA F-test (two or more groups). This procedure is often used to modify the distributional shape of the data so that the residuals are more normally distributed and/or the within-group variances are closer to equality. This is done so that tests and confidence limits that require normality can more appropriately be used. It cannot correct every data ill. For example, data that contain outliers may not be properly adjusted by this technique.

Example of the Box-Cox λ Plot



The Box-Cox transformation has the following mathematical form

$$Y = (X + \delta)^\lambda$$

where λ is the exponent (power) and δ is a shift amount that is added when X is zero or negative. When λ is zero, the above definition is replaced by

$$Y = \ln(X + \delta)$$

Box-Cox Transformation for Two or More Groups (T-Test and One-Way ANOVA)

Usually, the standard λ values of -2, -1.5, -1, -0.5, 0, 0.5, 1, 1.5, and 2 are investigated to determine which, if any, is most suitable. The program will also solve for the optimum value of λ using maximum likelihood estimation. The program also calculates confidence limits about the optimum value. The usual procedure is to adopt the most convenient standard value between the confidence limits. For example, if the confidence limits were 0.4 to 1.1, λ would be set to the standard value of '1' (no transformation) since this is the most convenient. Care must be used when using the confidence limits, because they are heavily dependent on the sample size.

Box-Cox Algorithm

Suppose you have a sample of n response values X_1, X_2, \dots, X_n divided into K groups. Further suppose you visually determine a value of δ that will keep all $X + \delta > 0$. Calculate a set of Z 's corresponding to the X 's using

$$Z = \begin{cases} [(X + \delta)^\lambda - 1]/[\lambda H^{\lambda-1}] & \lambda \neq 0 \\ H \ln(X + \delta) & \lambda = 0 \end{cases}$$

where H is the geometric mean of $X + \delta$. That is,

$$H = \sqrt[n]{\prod_{i=1}^n (X + \delta)}$$

Scaling by H is intended to keep the standard deviation of the Z 's approximately the same as the standard deviation of the X 's so that the standard deviations can be compared at various values of λ .

Maximum Likelihood Estimation of λ

In this case, the likelihood for a given λ is inversely proportional to the square root of the mean square error of the corresponding Z 's. The likelihood function is maximized when this value is minimized. A bracketing search algorithm is conducted that continues to tighten the boundaries until a specified precision (bracket width) is reached.

Approximate Confidence Interval for λ

An approximate confidence interval for λ is based on likelihood function which in turn is proportional to the standard deviation (SD) of Z . The confidence limits correspond to the two values of λ at which

$$SD_\lambda^2 = SD_{\hat{\lambda}}^2 \exp\left(\frac{\chi_1^2(1 - \alpha)}{n}\right)$$

where $\hat{\lambda}$ is the maximum likelihood estimate of λ and $\chi_1^2(1 - \alpha)$ is the percentage point of the chi-squared distribution with one degree of freedom.

F-Test

Note that because of the correspondence between the two-sided, two-sample t-test and one-way ANOVA F-test, only the results for the F-test are presented here. The formula for the one-way analysis of variance (ANOVA) F-test is

$$F_{K-1, n-K} = \frac{MSB}{MSE}$$

where

$$MSB = \frac{1}{K-1} \sum_{i=1}^K n_i (\bar{Z}_i - \bar{Z})^2$$

$$MSE = \frac{1}{n-K} \sum_{i=1}^K \sum_{j=1}^{n_i} (Z_{ij} - \bar{Z})^2$$

$$\bar{Z} = \frac{1}{n} \sum_{i=1}^K \sum_{j=1}^{n_i} Z_{ij}$$

$$n = \sum_{i=1}^K n_i$$

$$\bar{Z}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} Z_{ij}$$

The quantities *MSB* and *MSE* are often called the *mean square between* and *mean square error*, respectively. This *F* is compared to a tabulated *F* with *K-1* numerator degrees of freedom and *n-K* denominator degrees of freedom.

Data Structure

The data may be entered in two formats.

The first format puts the responses for each group in separate variables; that is, each variable contains all responses for a single group.

The second format arranges the data so that all responses are entered in a single variable. A second variable, the Factor Variable, contains an index that gives the group to which that row of data belongs.

In most cases, the second format is more flexible. Unless there is some special reason to use the first format, we recommend that you use the second.

Example 1 – Box-Cox Transformation for One-Way ANOVA

This section presents an example of how to run a Box-Cox transformation analysis on a set of data with 4 groups of six items each. The data used are found in the Box Cox One Way dataset.

Setup

To run this example, complete the following steps:

1 Open the BoxCoxOneWay example dataset

- From the File menu of the NCSS Data window, select **Open Example Data**.
- Select **BoxCoxOneWay** and click **OK**.

2 Specify the Box-Cox Transformation for Two or More Groups (T-Test and One-Way ANOVA) procedure options

- Find and open the Box-Cox Transformation for Two or More Groups (T-Test and One-Way ANOVA) procedure using the menus or the Procedure Navigator.
- The settings for this example are listed below and are stored in the **Example 1** settings file. To load these settings to the procedure window, click **Open Example Settings File** in the Help Center or File menu.

Variables Tab

Response Variable(s)**Rate**
 Factor Variable**Brand**

3 Run the procedure

- Click the **Run** button to perform the calculations and generate the output.

Run Summary Section

Run Summary for Rate with Factor = Brand

Parameter	Value	Parameter	Value
δ (Shift)	0	Rows Processed	24
Optimum λ (Power)	0.5356	Rows Used	24
Minimum λ Searched	-5	Geometric Mean (with δ)	329.4107
Maximum λ Searched	5	Minimum	12
Target Search Width of λ	0.0001	Maximum	1720
MLE Iterations Used	26	Max/Min (> 10?)	143.33
Max MLE Iterations	50		
Number of Groups	4		

This report summarizes the run by showing main results as well as the input settings that were used. You should pay particular attention to the *Rows* lines to make sure that they are as you expect. Also, if the number of MLE Iterations is equal to the Max MLE Iterations, the search algorithm may not have converged properly.

Box-Cox Transformation for Two or More Groups (T-Test and One-Way ANOVA)

When the ratio of the maximum to the minimum is greater than 10, the Box-Cox transformation is often useful.

Optimum (Maximum Likelihood) Estimate of λ

Optimum (Maximum Likelihood) Estimate of λ for X = Rate

Power Transformation: $Y = (X + \delta)^\lambda$

Item	Power λ	Shift δ	Standard Error of $Y =$ $(X + \delta)^\lambda$	Shapiro-Wilk Normality Test Prob Level	Levene's Equal Variance Test Prob Level
Optimum (MLE)	0.5356	0	132.7967	0.3026	0.8728
Lower 95% C. L.	0.3225	0	143.8649	0.8888	0.1573
Upper 95% C. L.	0.7651	0	143.8605	0.3129	0.1373

This report gives the results for the maximum likelihood estimation portion of the analysis.

Item

The name of item being reported on this line of the report.

Power λ

The value of λ for this item. This is the transformation exponent.

Shift δ

The value of δ , the shift value.

Standard Error of $Y = (X + \delta)^\lambda$

This is the square root of the mean square error of the transformed data values. Actually, the data have not only been shifted and raised to the indicated power, but they have also been scaled by the geometric mean so that these values are directly comparable. Note the geometric mean is not used when using the λ that has been found by this algorithm.

Shapiro-Wilk Normality Test Prob Level

The probability level of the Shapiro-Wilk normality test. Since the desire is to transform the data to be more normally distributed, you are looking for large (non-significant) values. Remember that this value is not only influenced by the normality of the data, but also by the sample size.

Levene's Equal-Group Variance Test Prob Level

Besides achieving normality, it is also hoped that the power transformation will result in group variances that are more closely identical. The Levene test is a popular way of testing how well this assumption is met. Since we want the variances to be nearly equal, we are looking for large values (close to one) of this probability.

Standard λ 's

Standard λ 's for X = Rate

Power Transformation: $Y = (X + \delta)^\lambda$

Item	Power λ	Shift δ	Standard Error of $Y = (X + \delta)^\lambda$	Shapiro-Wilk Normality Test Prob Level	Levene's Equal Variance Test Prob Level
1	-2.0000	0	30819.2299	0.0000	0.0000
2	-1.0000	0	1990.2394	0.0000	0.0000
3	-0.5000	0	586.1318	0.0001	0.0000
4	0.0000	0	214.4276	0.0527	0.0007
5	0.5000	0	133.0863	0.4141	0.8397
6	1.0000	0	176.1833	0.3814	0.0032
7	1.5000	0	322.5846	0.0746	0.0000
8	2.0000	0	639.8299	0.0067	0.0000

λ 's between the maximum likelihood confidence limits are bolded.

This report displays the results for each of the standard λ 's.

Item

The number of items being reported on this line of the report.

Power λ

The value of λ for this item. This is the transformation exponent.

Shift δ

The value of δ , the shift value.

Standard Error of $Y = (X + \delta)^\lambda$

This is the square root of the mean square error of the transformed data values. Actually, the data have not only been shifted and raised to the indicated power, but they have also been scaled by the geometric mean so that these values are directly comparable. Note the geometric mean is not used when using the λ that has been found by this algorithm.

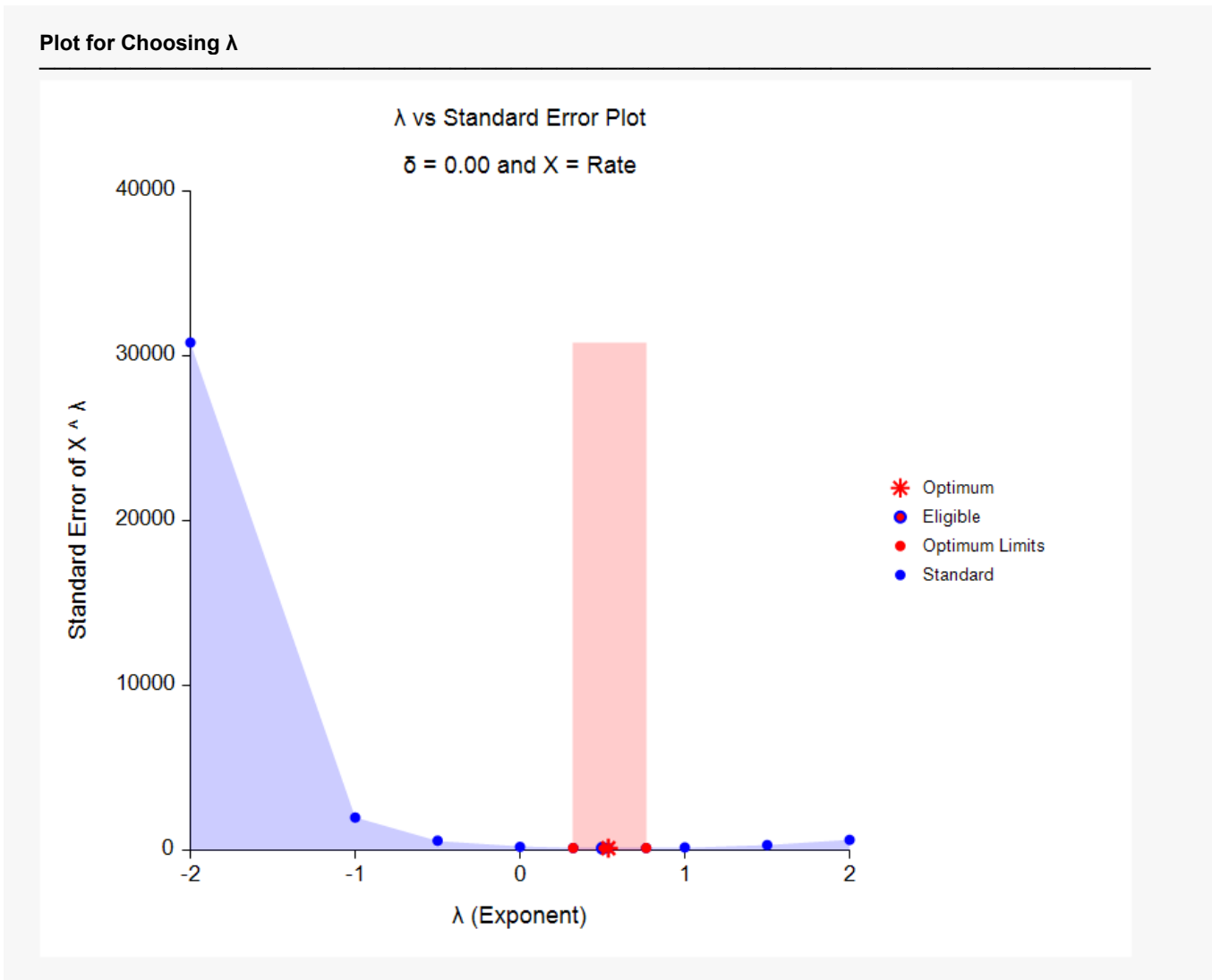
Shapiro-Wilk Normality Test Prob Level

The probability level of the Shapiro-Wilk normality test. Since the desire is to transform the data to be more normally distributed, you are looking for large (non-significant) values. Remember that this value is not only influenced by the normality of the data, but also by the sample size.

Levene's Equal Variance Test Prob Level

Besides achieving normality, it is also hoped that the power transformation will result in group variances that are more closely identical. The Levene test is a popular way of testing how well this assumption is met. Since we want the variances to be nearly equal, we are looking for large values (close to one) of this probability.

Plots

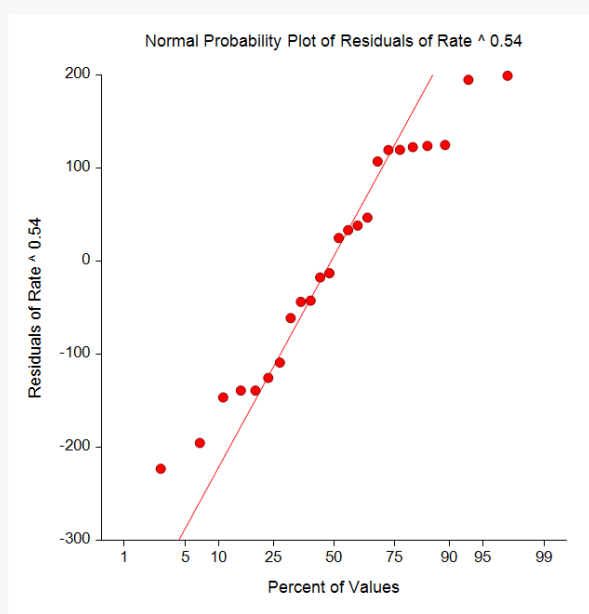
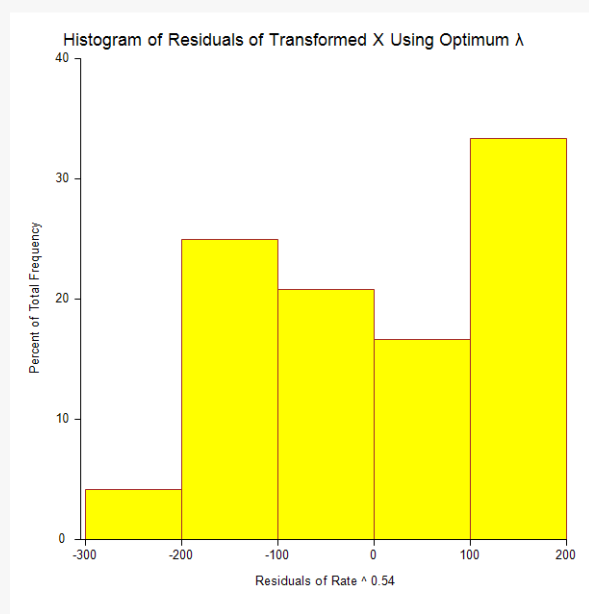
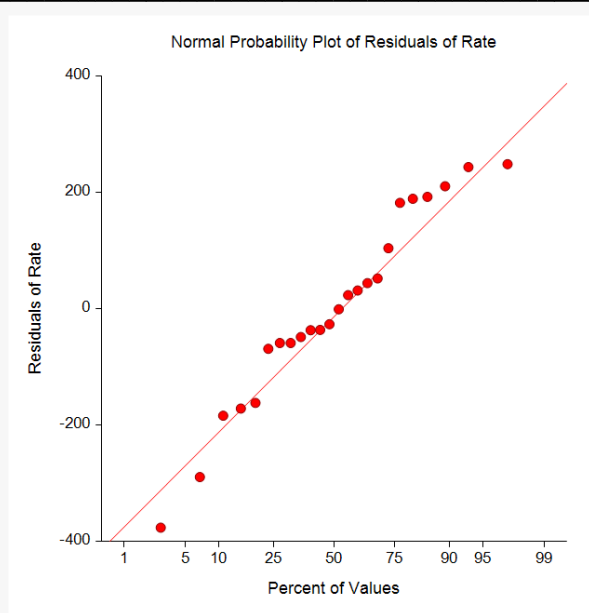
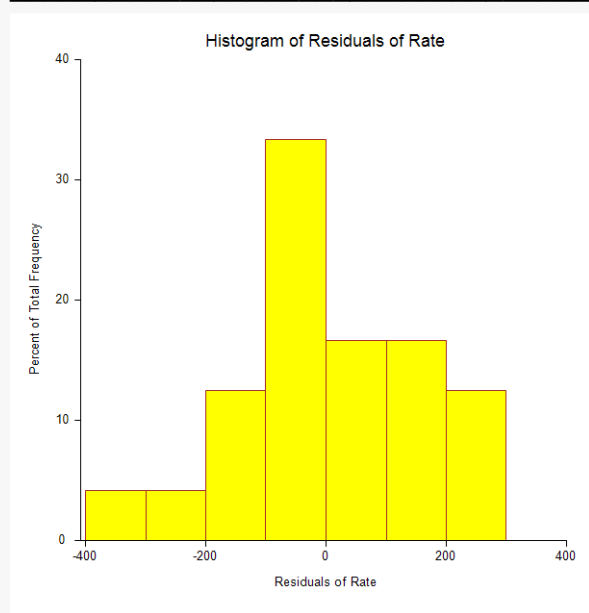


This plot gives a visual representation that will help you select the value of λ that you want to use. The optimum value found by maximum likelihood is plotted with a large, red asterisk. This value is usually inconvenient to use, so a convenient (standard) value is sought for that is close to the optimum value. These convenient values are plotted using a blue circle with a red center. In this example, it is obvious that $\lambda = 0.5$ (square root) is certainly a reasonable choice. The large, shaded area in the middle of the plot highlights the values of λ that are within the confidence interval for the optimum.

Note that this plot was created using the Scatter Plot procedure. The shading effects and different plot symbols were made by making several groups of data.

Box-Cox Transformation for Two or More Groups (T-Test and One-Way ANOVA)

Plots for Assessing Normality at Various λ 's



These plots let you see the improvement towards normality achieved by the power transformation. The top row shows the histogram and probability plot of the original data. The lack of normality is evident in the histogram, although the non-normality does not appear to be very severe.

The bottom row of plots shows the same two plots applied to the data that has been transformed by the optimum λ . The histogram is now much closer to being bell shaped. The normal probability plots don't appear that different in this case.