

# Model Diagnostics and Remedial Measures

## One-Way ANOVA – Stat 230

In this class, you will work with your group to explore the use of normal Q-Q plots to assess the normality of residuals in ANOVA models. You will also learn how to apply transformations to the response variable in order to remedy violations to the model assumptions.

The [R Manual](#) has useful R code for today’s activities.

### Reading Normal Q-Q Plots

In the reading you learned how normal Q-Q plots are constructed and why they are useful. Today, we’ll start off by creating a number of these plots to practice reading them.

**Activity A.** Go to <https://loya.shinyapps.io/qplots/>. This is a web app where you can generate a histogram and normal Q-Q plot for a variety of sample sizes and distributional shapes. Looking at multiple examples will help you develop your intuition about what constitute a substantial deviation from a straight line on a normal Q-Q plot. Do the following:

- The default settings are a sample size of 15 observations drawn from a normal population distribution. Generate a few of these plots by pressing “Plot it!”. Do the points always fall on the line?
- Now, increase the sample size. How does the Q-Q plot change as the sample size increases?
- Now, change the shape of the population distribution to “right skewed”. Generate a few Q-Q plots at a few sample sizes. Describe the form/shape of the Q-Q plot for right-skewed data.
- Now, change the shape of the population distribution to “left skewed”. Generate a few Q-Q plots at a few sample sizes. Describe the form/shape of the Q-Q plot for left-skewed data.
- Now, change the shape of the population distribution to “heavy tailed”. Generate a few Q-Q plots at a few sample sizes. Describe the form/shape of the Q-Q plot for heavy-tailed data.

Now that you have calibrated your intuition, let’s create a normal Q-Q plot in R.

**Activity 29.** Load the `Normal` data set using the code (which can also be copied from the R Manual):

```
normal <- read.csv("https://aloy.rbind.io/kuiper_data/Normal.csv")
```

In this activity you will practice creating normal Q-Q plots in R, and also how outliers impact the plot. Work through parts (a)-(c). All of the data is available in the `Normal` data set, so there is no “extra” R work needed.

## Diagnosing Normality

Now that you have had more practice reading normal Q-Q plots, let's use the understanding you've gained to diagnose models!

In Activity 26 you calculated the residuals from the one-way ANOVA model for the `Games1` data set using the following code:

```
games1 <- read.csv("https://aloy.rbind.io/kuiper_data/Games1.csv")
games_anova <- aov(Time ~ Type, data = games1)
games_resid <- resid(games_anova)
```

### Exercise E.10

- Create a normal Q-Q plot and histogram of the residuals. Comment on the normality assumption for the random error terms.
- Create a normal Q-Q plot and histogram of the observed response variable, `Time`.
- Explain why residuals should be used instead of the observed responses to test the normality assumption.

## Applying Transformations

Transforming the response variable can often help remedy non-normal error terms or subgroups with dramatically different variances. The textbook has a nice discussion on when the square root, log, reciprocal, and logit transformations are often useful, so be sure to review this.

**Extended Activity 31.** To practice working with transformations work through activity 31 c-d in the textbook. You can load the data via

```
emissions <- read.csv("https://aloy.rbind.io/kuiper_data/Emissions.csv")
```

In activity 31, you applied a transformation to remedy violations so that you could conduct an appropriate F-test for a one-way ANOVA model. Now, let's consider transforming the response for the simple linear regression formulation of a two-sample t-test.

The below code extracts the rows of the data set for the pre-63 and 70-71 cases.

```
sub_emissions <- emissions |>
  filter(Year == "1963-1967" | Year == "1970-1971")
```

Now, let's fit a simple linear regression model with an indicator variable that flags 1970-1971 (i.e., the indicator is 1 when for cases from 1970-1971). In this model, a natural log transformation was applied to the response variable.

```
emissions_lm <- lm(log(Emissions) ~ Year, data = sub_emissions)
```

As seen in the book, a 95% confidence interval for the difference in mean emissions between the years is given by

```
confint(emissions_lm, parm = 2) # CI only for the second parameter
```

	2.5 %	97.5 %
Year1970-1971	-1.295767	-0.3771826

However, this interval is on the natural log scale, which isn't as meaningful/interpretable as the original emissions scale. To get back to the original scale, we need to undue the transformation, which can be done in R (no need for a calculator!)

```
exp(confint(emissions_lm, parm = 2))
```

	2.5 %	97.5 %
Year1970-1971	0.2736878	0.6857909

### Danger

Back transforming a logarithmic transformation does not make the results have the “same meaning” as the original data. Caution is needed when interpreting parameter estimates and confidence intervals with log-transformed data.

**Activity B.** To discover how this back transformation can be interpreted, let's work with the two samples directly:

```
y1 <- sub_emissions |>
  filter(Year == "1963-1967") |>
  pull(Emissions)
y2 <- sub_emissions |>
  filter(Year == "1970-1971") |>
  pull(Emissions)
```

- Compute the sample means and medians for the untransformed samples, `y1` and `y2`.
- Compute the sample means and medians for the transformed samples, `log(y1)` and `log(y2)`.
- Is the natural log of the difference in sample means approximately equal to the difference in means of the logged data?
- Is the natural log of the ratio of sample means from the untransformed samples approximately equal to the difference in means of the logged data?
- Is the natural log of the ratio of sample medians from the untransformed samples approximately equal to the difference in medians of the logged data?
- Based on (c)-(e), how would you suggest interpreting the back transformation given by:

```
exp(mean(log(y1)) - mean(log(y2)))
```

```
[1] 2.308216
```

Recall that if a distribution is normal, then the mean and median will be equal.

Based on your suggestion, write a one-sentence interpretation of the back-transformed confidence interval given above.

**Activity C.** Load the `AnovaTrans` data set using the below code chunk, then do the following for each of the two response variables.

```
anova_trans <- read.csv("https://aloy.rbind.io/data/AnovaTrans.csv")
```

- (a) Create an individual values plot of the response (`y1`, `y2`, or `y3`) against `x`.
- (b) Fit the one-way ANOVA model and extract the residuals.
- (c) Create a plot of the residuals against `x` and a normal Q-Q plot of the residuals.
- (d) Try various transformations of the response variable to create a better one-way ANOVA model (as validated by graphical analysis of the residuals).