

Log Transformations

This document provides details about the model interpretation when the predictor and/or response variables are log-transformed. For simplicity, we will discuss transformations for the simple linear regression model:

$$y = \beta_0 + \beta_1 x \quad (1)$$

All results and interpretations can be easily extended to transformations in multiple regression models.

Note: \log refers to the natural logarithm.

Log-transformation on the response variable

Suppose we fit a linear regression model with $\log(y)$, the log-transformed y , as the response variable. Under this model, we assume a linear relationship exists between x and $\log(y)$, such that $\log(y) \sim N(\beta_0 + \beta_1 x, \sigma^2)$ for some β_0 , β_1 and σ^2 . In other words, we can model the relationship between x and $\log(y)$ using the model in (2).

$$\log(y) = \beta_0 + \beta_1 x \quad (2)$$

If we interpret the model in terms of $\log(y)$, then we can use the usual interpretations for slope and intercept. When reporting results, however, it is best to give all interpretations in terms of the original response variable y , since interpretations using log-transformed variables are often more difficult to truly understand.

In order to get back on the original scale, we need to use the exponential function (also known as the anti-log), $\exp\{x\} = e^x$. Therefore, we use the model in (2) for interpretations and predictions, we will use (3) to state our conclusions in terms of y .

$$\begin{aligned} \exp\{\log(y)\} &= \exp\{\beta_0 + \beta_1 x\} \\ \Rightarrow y &= \exp\{\beta_0 + \beta_1 x\} \\ \Rightarrow y &= \exp\{\beta_0\} \exp\{\beta_1 x\} \end{aligned} \quad (3)$$

In order to interpret the slope and intercept, we need to first understand the relationship between the mean, median and log transformations.

Mean, Median, and Log Transformations

Suppose we have a dataset y that contains the following observations:

```
y <- c(3,5,6,7,8)
y
```

```
## [1] 3 5 6 7 8
```

If we log-transform the values of y then calculate the mean and median, we have

```
log_y <- tibble(log_y = log(y))
summary <- log_y %>%
  summarise(mean_log_y = mean(log_y), median_log_y = median(log_y))
kable(summary,digits=5)
```

mean_log_y	median_log_y
1.70503	1.79176

If we calculate the mean and median of y , then log-transform the mean and median, we have

```
centers <- tibble(y) %>% summarise(mean_y = mean(y), median_y = median(y))
summary2 <- centers %>%
  summarise(log_mean = log(mean_y), log_median = log(median_y))
kable(summary2,digits=5)
```

log_mean	log_median
1.75786	1.79176

This is a simple illustration to show

1. $\text{Mean}[\log(y)] \neq \log[\text{Mean}(y)]$ - the mean and log are not commutable
2. $\text{Median}[\log(y)] = \log[\text{Median}(y)]$ - the median and log are commutable

Interpretation of model coefficients

Using (2), the mean $\log(y)$ for any given value of x is $\beta_0 + \beta_1 x$; however, this does **not** indicate that the mean of $y = \exp\{\beta_0 + \beta_1 x\}$ (see previous section). From the assumptions of linear regression, we assume that for any given value of x , the distribution of $\log(y)$ is Normal, and therefore symmetric. Thus the median of $\log(y)$ is equal to the mean of $\log(y)$, i.e. $\text{Median}(\log(y)) = \beta_0 + \beta_1 x$.

Since the log and the median are commutable, $\text{Median}(\log(y)) = \beta_0 + \beta_1 x \Rightarrow \text{Median}(y) = \exp\{\beta_0 + \beta_1 x\}$. Thus, when we log-transform the response variable, the interpretation of the intercept and slope are in terms of the effect on the **median** of y .

Intercept: The intercept is expected median of y when the predictor variable equals 0. Therefore, when $x = 0$,

$$\begin{aligned}\log(y) &= \beta_0 + \beta_1 \times 0 = \beta_0 \\ \Rightarrow y &= \exp\{\beta_0\}\end{aligned}\tag{4}$$

Interpretation: When $x = 0$, the median of y is expected to be $\exp\{\beta_0\}$.

Slope: The slope is the expected change in the median of y when x increases by 1 unit. The change in the median of y is

$$\exp\{[\beta_0 + \beta_1(x+1)] - [\beta_0 + \beta_1 x]\} = \frac{\exp\{\beta_0 + \beta_1(x+1)\}}{\exp\{\beta_0 + \beta_1 x\}} = \frac{\exp\{\beta_0\} \exp\{\beta_1 x\} \exp\{\beta_1\}}{\exp\{\beta_0\} \exp\{\beta_1 x\}} = \exp\{\beta_1\} \quad (5)$$

Thus, the median of y for $x+1$ is $\exp\{\beta_1\}$ times the median of y for x .

Interpretation: When x increases by one unit, the median of y is expected to multiply by a factor of $\exp\{\beta_1\}$.

Log-transformation on the predictor variable

Suppose we fit a linear regression model with $\log(x)$, the log-transformed x , as the predictor variable. Under this model, we assume a linear relationship exists between $\log(x)$ and y , such that $y \sim N(\beta_0 + \beta_1 \log(x), \sigma^2)$ for some β_0 , β_1 and σ^2 . In other words, we can model the relationship between $\log(x)$ and y using the model in (6).

$$y = \beta_0 + \beta_1 \log(x) \quad (6)$$

Intercept: The intercept is the mean of y when $\log(x) = 0$, i.e. $x = 1$.

Interpretation: When $x = 1$ ($\log(x) = 0$), the mean of y is expected to be β_0 .

Slope: The slope is interpreted in terms of the change in the mean of y when x is multiplied by a factor of C , since $\log(Cx) = \log(x) + \log(C)$. Thus, when x is multiplied by a factor of C , the change in the mean of y is

$$\begin{aligned} -[\beta_0 + \beta_1 \log(x)] &= \beta_1 [\log(Cx) - \log(x)] \\ &= \beta_1 [\log(C) + \log(x) - \log(x)] \\ &= \beta_1 \log(C) \end{aligned} \quad (7)$$

Thus the mean of y changes by $\beta_1 \log(C)$ units.

Interpretation: When x is multiplied by a factor of C , the mean of y is expected to change by $\beta_1 \log(C)$ units. For example, if x is doubled, then the mean of y is expected to change by $\beta_1 \log(2)$ units.

Log-transformation on the the response and predictor variable

Suppose we fit a linear regression model with $\log(x)$, the log-transformed x , as the predictor variable and $\log(y)$, the log-transformed y , as the response variable. Under this model, we assume a linear relationship exists between $\log(x)$ and $\log(y)$, such that $\log(y) \sim N(\beta_0 + \beta_1 \log(x), \sigma^2)$ for some β_0 , β_1 and σ^2 . In other words, we can model the relationship between $\log(x)$ and $\log(y)$ using the model in (8).

$$\log(y) = \beta_0 + \beta_1 \log(x) \quad (8)$$

Because the response variable is log-transformed, the interpretations on the original scale will be in terms of the median of y (see the section on the log-transformed response variable for more detail).

Intercept: The intercept is the mean of y when $\log(x) = 0$, i.e. $x = 1$. Therefore, when $\log(x) = 0$,

$$\begin{aligned} \log(y) &= \beta_0 + \beta_1 \times 0 = \beta_0 \\ \Rightarrow y &= \exp\{\beta_0\} \end{aligned} \quad (9)$$

Interpretation: When $x = 1$ ($\log(x) = 0$), the median of y is expected to be $\exp\{\beta_0\}$.

Slope: The slope is interpreted in terms of the change in the median y when x is multiplied by a factor of C , since $\log(Cx) = \log(x) + \log(C)$. Thus, when x is multiplied by a factor of C , the change in the median of y is

$$\begin{aligned}
\exp\{[\beta_0 + \beta_1 \log(Cx)] - [\beta_0 + \beta_1 \log(x)]\} &= \exp\{\beta_1[\log(Cx) - \log(x)]\} \\
&= \exp\{\beta_1[\log(C) + \log(x) - \log(x)]\} \\
&= \exp\{\beta_1 \log(C)\} = C^{\beta_1}
\end{aligned} \tag{10}$$

Thus, the median of y for Cx is C^{β_1} times the median of y for x .

Interpretation: When x is multiplied by a factor of C , the median of y is expected to multiple by a factor of C^{β_1} . For example, if x is doubled, then the median of y is expected to multiply by 2^{β_1} .