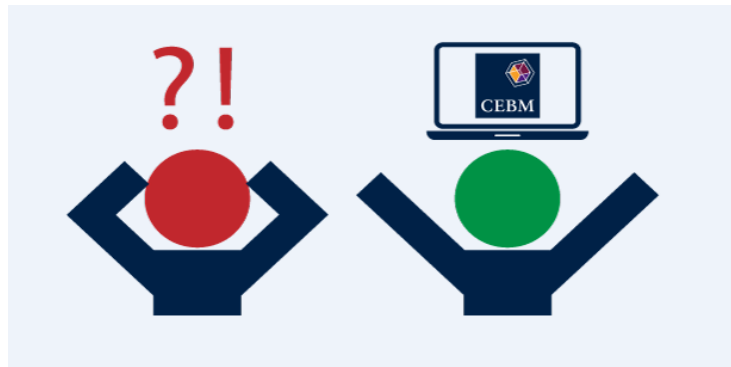


## Tip for data extraction for meta-analysis - 25



### What if you're missing a mean and only a similar statistical statistic is given?

Kathy Taylor

[Previously](#), I highlighted a list of ways where, when extracting data for meta-analysis of continuous outcomes, you might find that a summary statistic that you want is missing. In this post I'll give some examples of the 3<sup>rd</sup> way - **a similar summary statistic is reported, but it's not the statistical measure that you want** - when you have missing means.

### Finding a median reported

You may find that instead of a mean, a median is reported. A median is a different type of average. The reporting of medians indicates that the distribution of outcome data is skewed. The median and mean are equal if the distribution of the data is perfectly symmetrical (Figure 1). When the distribution is skewed, the mean and median will differ, and the difference between them will depend on the degree of skewness.

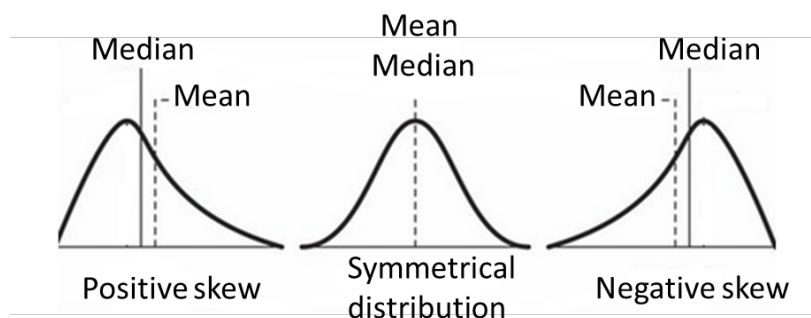


Figure 1. The relative position of the mean and median depending on the data's distribution

[The Cochrane Handbook \(section 6.5.2.9\)](#) highlights three papers which provide equations for estimating means from other summary statistics.

One is the paper by [Hozo et al](#) who concluded that even with skewed data, the sample mean can be estimated by the median. They provided estimates for the mean based on the sample size and range (min and max are the bounds of the range):

$$\text{mean} \approx \text{median} \quad \text{if } n \geq 25$$

$$\text{mean} \approx \frac{\text{min} + 2\text{median} + \text{max}}{4} \quad \text{if } n < 25$$

The curved equal sign means ‘approximately equal’. The estimates were tested using simulation, and drawing samples from normal and skewed distributions.

The second paper is by [Bland](#) who uses more information by providing estimates of the sample mean based on the median, range, sample size and interquartile range ( $q_1$  to  $q_2$ ):

$$\frac{\text{min}(n + 3) + 2(n - 1)(q_1 + \text{median} + q_3) + \text{max}(n + 3)}{8n}$$

When  $n$  is large, this equation simplifies to

$$\frac{\text{min} + 2(q_1 + \text{median} + q_3) + \text{max}}{8}$$

He tested his estimates on three real data sets with simulated data drawing samples from normal and skewed distributions.

The formulae of Bland and Hozo et al both work better with small samples.

The third paper is by [Wan et al](#) who provide estimates of the sample mean based on the median and interquartile range. This has the advantage of not being influenced by extreme values.

$$\frac{q_1 + \text{median} + q_3}{3}$$

Using simulation, they also tested their estimates, drawing samples from normal and skewed distributions and found smaller relative errors compared to Bland’s approach. Wan et al also provide a very useful spreadsheet which you can use to calculate and compare their estimated means with those of Bland and Hozo et al.

Another paper by [Luo et al](#) provide improved estimates of the sample mean based on the sample size, median, range and interquartile range. Their estimates use a weighted formulation. They consider three scenarios.

Firstly, when the sample size, median and range are reported.

$$\left(\frac{4}{4 + n^{0.75}}\right)\left(\frac{\min + \max}{2}\right) + \left(\frac{n^{0.75}}{4 + n^{0.75}}\right) \text{median}$$

Secondly, when the sample size, median and interquartile range are reported.

$$\left(0.7 + \frac{0.39}{n}\right)\left(\frac{q_1 + q_3}{2}\right) + \left(0.3 - \frac{0.39}{n}\right) \text{median}$$

Thirdly, when the sample size, median, range and interquartile range are reported.

$$\left(\frac{2.2}{2.2 + n^{0.75}}\right)\left(\frac{\min + \max}{2}\right) + \left(0.7 - \frac{0.72}{n^{0.55}}\right)\left(\frac{q_1 + q_3}{2}\right) + \left(0.3 + \frac{0.72}{n^{0.55}} - \frac{2.2}{2.2 + n^{0.75}}\right) \text{median}$$

These are approximations of their more complicated formulae that are reported in their paper. They demonstrate the accuracy of their estimates using simulation and they provide an excel spreadsheet.

### **Finding a geometric mean reported**

Sometimes a geometric mean is reported. This is another type of average, which arises from the analysis of skewed data which have been [log-transformed](#) and then back-transformed (using the exponential function) when presenting results. With small samples, skewed data is often log-transformed, before analysis, because standard inferences on the means of skewed data is only acceptable for large samples. With large samples we assume that the means of outcome measurements are approximately normally distributed due to the [central limit theorem](#).

So, instead of means (which are more formally known as arithmetic means) and standard deviations (SD), geometric means are reported, either with confidence intervals (CIs), the exponential of the SD of the log-transformed values (often referred to as the tolerance factor or the inappropriately named as the 'SD of the geometric mean'), or the exponential of the standard error (SE) of the log-transformed values. Geometric means and arithmetic means should [not be pooled](#). If most of your studies report arithmetic means, you will want to convert geometric mean summary data to arithmetic mean summary data. Pooling is possible by using the conversion equations of [Higgins et al](#). It's a two-stage process (Figure 2) as the geometric mean data has to be log-transformed first.

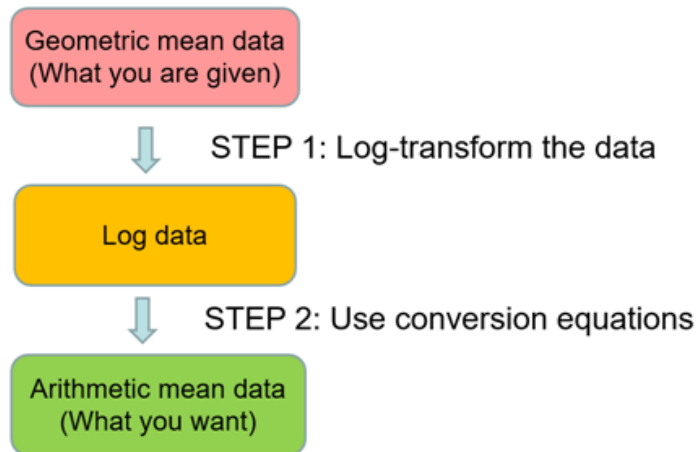


Figure 2. From geometric means to arithmetic means

Using the following notation:

$g$  is the geometric mean of a treatment arm

$(g_{lower} \text{ to } g_{upper})$  is the confidence interval of  $g$

$e^{S_z}$  is the incorrectly named SD of  $g$

$e^{SE_z}$  is the incorrectly named SE of  $g$

### STEP 1

Calculate the log-transformed measurements ( $\bar{z}$  and  $s_z$ ) from the geometric mean data.

$$\bar{z} = \ln(g) \quad \text{and}$$

$$s_z = \frac{(\ln(g_{upper}) - \ln(g_{lower}))\sqrt{n}}{2t} \quad \text{OR}$$

$$s_z = \ln(e^{S_z}) \text{ in cases where } e^{S_z} \text{ has been reported OR}$$

$$s_z = \ln(e^{\sqrt{n} \times SE_z})$$

where  $t$  is the 97.5 percentage point of the t-distribution with  $(n-1)$  degrees of freedom.

### STEP 2

Apply the conversion equations to the log-transformed data to calculate the arithmetic mean summary data ( $\bar{x}$  and  $s_x$ ). There are two sets of equations depending on the similarities between the SDs of the two treatment arms. Higgins et al recommend comparing the SDs on the log scale as it's more plausible. If the SDs are different, use Method 1. If the SDs are similar, use Method 2.

## Method 1

For each treatment arm, calculate

$$\bar{x} = \exp\left(\bar{z} + \frac{s_z^2}{2}\right)$$
$$s_x = \sqrt{(\exp(s_z^2) - 1)\exp(2\bar{z} + s_z^2)}$$

## Method 2

First calculate

$$s_{z,pooled} = \sqrt{\frac{(n_1 - 1)s_{z,1}^2 + (n_2 - 1)s_{z,2}^2}{n_1 + n_2 - 2}}$$

Then, for each treatment arm, calculate

$$\bar{x} = \exp\left(\bar{z} + \frac{s_{z,pooled}}{2}\right)$$
$$s_x = \sqrt{(\exp(s_{z,pooled}^2) - 1)\exp(2\bar{z} + s_{z,pooled}^2)}$$

Higgins et al also provide equations to convert the other way, from arithmetic means to geometric means (Figure 3). You might want this if the majority of your included studies report geometric means.

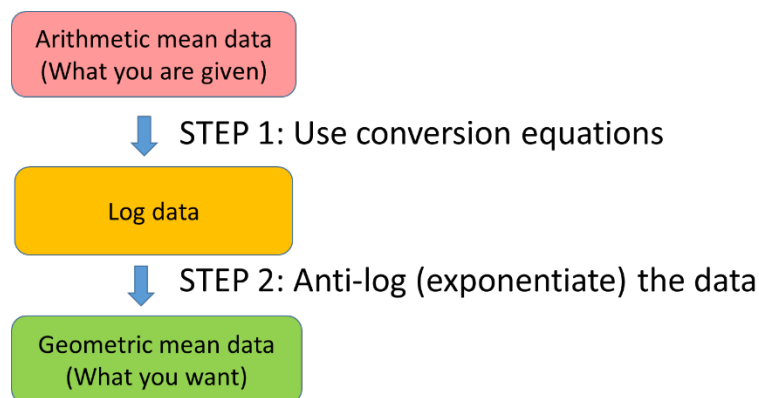


Figure 3. From arithmetic means to geometric means

### STEP 1

Convert the arithmetic mean summary data to log summary data.

$\bar{x}$  and  $s_x$  are the arithmetic mean and SD

$\bar{z}$  and  $s_z$  are the mean and SD of the log data

If the SDs are different, use Method 1. If the SDs are similar, use Method 2.

### Method 1

For each treatment arm, calculate

$$\bar{z} = \ln(\bar{x}) - \frac{1}{2} \ln\left(\frac{s_x^2}{\bar{x}^2} + 1\right)$$

$$s_z = \sqrt{\ln\left(\frac{s_x^2}{\bar{x}^2} + 1\right)}$$

### Method 2

For each treatment arm, calculate

$$s_z = \sqrt{\ln\left(\frac{s_x^2}{\bar{x}^2} + 1\right)}$$

Then calculate

$$s_{z,pooled} = \sqrt{\frac{(n_1 - 1)s_{z,1}^2 + (n_2 - 1)s_{z,2}^2}{n_1 + n_2 - 2}}$$

Then for each treatment arm, calculate

$$\bar{z} = \ln(\bar{x}) - \frac{1}{2} s_{z,pooled}^2$$

## STEP 2

Back-transforming (exponentiating) the log data calculates the geometric mean data.

$\bar{z}$  and  $s_z$  are the log data.

$g = e^{\bar{z}}$  and  $e^{s_z}$  are the geometric mean data.

Meta-analysis can be carried out on the log scale and SD for the log values can be calculated using the following equation:

$$SD = \frac{(\text{upper CI} - \text{lower CI})}{3.92} \sqrt{n}$$

Let me show you an example of converting geometric mean data to arithmetic mean data.

For the [review](#) that I worked on we used these conversion equations when we extracted data from the study by [Romero et al.](#) They reported microalbuminuria (albumin excretion rate) at 6 months with geometric mean (95% CI). Data for the intervention group, treated with

Captopril, was converted from 60 (35 to 104) mg/24hr to mean (SD) of 90 (101) mg/24hr, and data for the untreated group was converted from 91 (58 to 141) mg/24hr was converted to 119 (101) mg/24hr (Table). I applied the equations from Method 1.

Table. Calculating arithmetic mean data

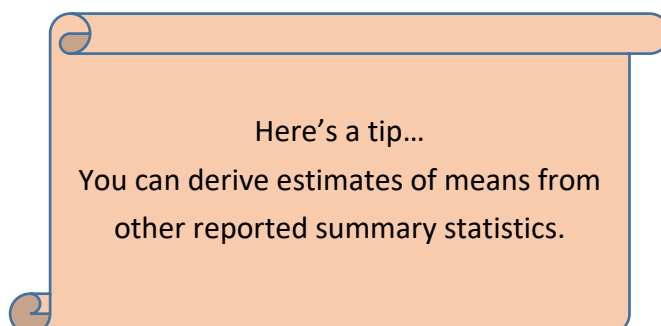
	Captopril	Untreated
n	13	13
gm	60	91
g lower	35	58
g upper	104	141
$z=\ln(g)$	4.09	4.51
dof	12	12
t-value	2.18	2.18
Sz	0.90	0.74
x	90.05	119.22
Sx	100.77	100.91

See below to find out where the t-value came from. Higgins et al highlight that their estimates are likely to be biased in small sample studies. As this study had small samples, it was removed as part of a sensitivity analysis.

### Other methods

Other approaches of dealing with missing means highlighted by the review of [Weir et al](#) include the simulation-based approximate Bayesian computation (ABC) approach of [Kwon and Reis](#).

If a large proportion of studies have missing means, pooling is not recommended.



In my next blog post I'll give some more examples of a similar summary statistic is reported, but it's not the statistical measure that you want when you have missing SDs.

Where did the equations and t-value come from?

*Converting geometric means to arithmetic means and the reverse*

[Higgins et al](#) derive their equations in their paper.

*Calculating means from medians, range and interquartile range*

[Hozo et al](#), [Bland](#), [Wan et al](#) and [Luo et al](#) also derive their equations in their respective papers.

*Calculating an SD from a 95% confidence interval:*

This was [derived earlier](#) and in my next blog post I will give more details.

*What about the t-value?*

This came from a t-distribution table (Figure 4).

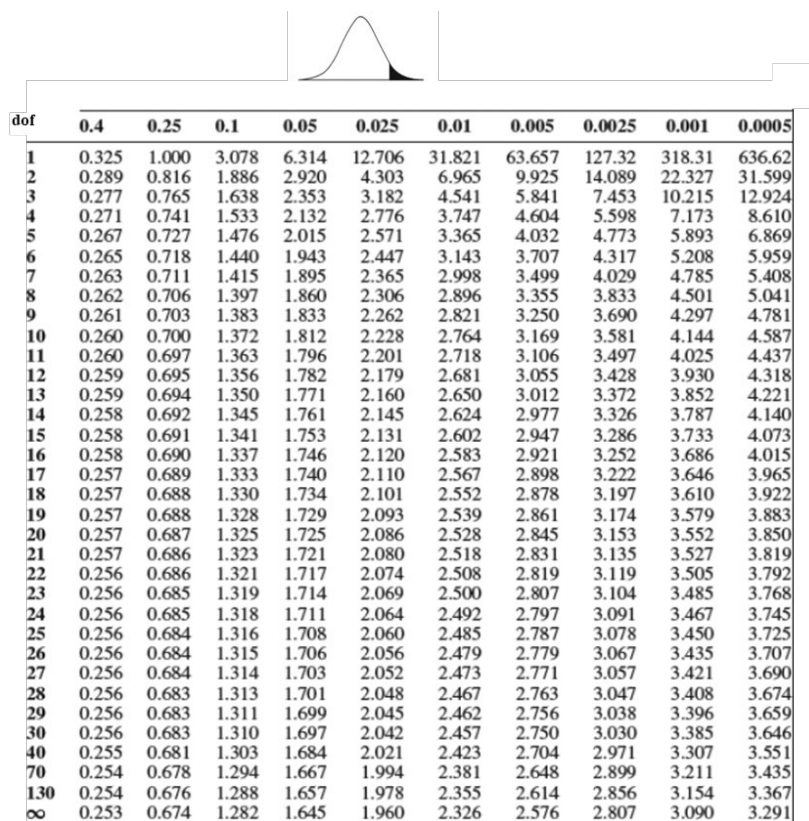
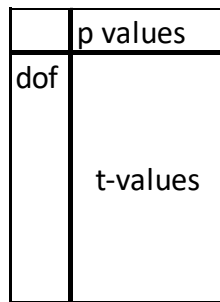


Figure 4. t-distribution table

The first column of the table shows the degrees of freedom (dof) and the area probabilities (also known as percentages or p-values) are shown in the first row. As indicated in the t-



distribution curve above the table, the p-values represent the area under the t-distribution curve in the tail, from the t-value to infinity (shaded black) for different dofs.



In the example I gave, the Captopril group had  $n=13$  so  $dof=n-1=12$ .

The area probability for the 97.5 percentage point of the t-distribution using the above table  $=1-0.975=0.025$ . The corresponding t-value is 2.179 (Figure 5).

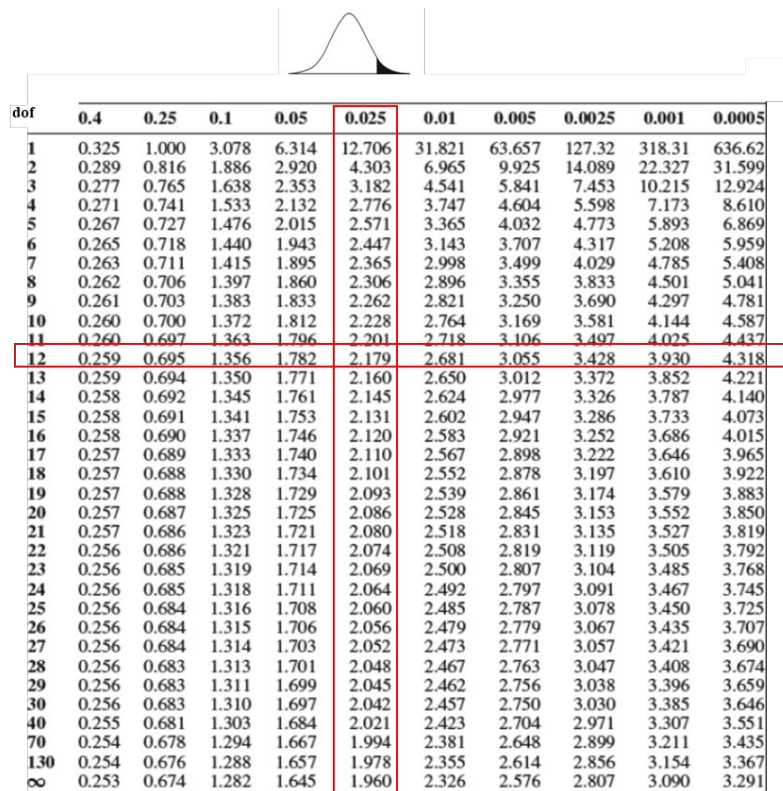


Figure 5. t-distribution table for  $dof=12$  and  $p\text{-value}=0.025$

The same t-value can also be calculated using EXCEL by typing in an EXCEL cell  $=ABS(T.INV(\text{one sided p-value}, dof))$ . With my example this is  $=ABS(T.INV(0.025,12))$ .

Note that I'm using EXCEL 2016. Earlier versions use the term  $TINV(\text{two-sided p-value}, dof)$  i.e.  $TINV(0.05,12)=ABS(T.INV(0.025,12))=2.178813$ .

Dr Kathy Taylor teaches data extraction in Meta-analysis,  
<https://www.conted.ox.ac.uk/courses/meta-analysis> This is a short course that is  
also available as part of our MSc in Evidence-Based Health Care  
<https://www.conted.ox.ac.uk/about/msc-in-evidence-based-health-care>, MSc in Medical  
Statistics  
<https://www.conted.ox.ac.uk/about/msc-in-ebhc-medical-statistics>, and MSc in  
Systematic Reviews  
<https://www.conted.ox.ac.uk/about/msc-in-ebhc-systematic-reviews>

Follow me on Twitter @dataextips for updates on my blog, related news, and to find out  
about further examples where others, like me, are trying to make statistics more broadly  
accessible.