

# Standardised mean differences: A tutorial

Daniel Gallardo Gómez<sup>1</sup>, Rachel Richardson<sup>2</sup>, and Kerry Dwan<sup>3</sup>

<sup>1</sup>University of Seville Faculty of Education Sciences

<sup>2</sup>Cochrane

<sup>3</sup>Liverpool School of Tropical Medicine

March 15, 2024

## Abstract

This tutorial focuses on standardised mean differences (SMD) as effect measures in meta-analyses. We will explain what they are, when they should be used, how to correctly compute and interpret them, and some of the most common error made within evidence synthesis

## Standardised mean differences in meta-analysis: A tutorial

Daniel Gallardo-Gómez<sup>1</sup>, Kerry Dwan<sup>2</sup>, Rachel Richardson<sup>3</sup>

<sup>1</sup> Department of Physical Education and Sport, University of Seville, Seville, Spain. Epidemiology of Physical Activity and Fitness Across Lifespan (EPAFit) Research Group

<sup>2</sup>Centre for Evidence Based Medicine, Department of Clinical Sciences, Liverpool School of Tropical Medicine, Liverpool, United Kingdom of Great Britain and Northern Ireland

<sup>3</sup>Cochrane Methods Support Unit, Evidence Production and Methods Department, Cochrane, London, United Kingdom of Great Britain and Northern Ireland

Correspondence: Daniel Gallardo Gómez; Department of Physical Education and Sport, University of Seville; [daniel.gallardogomez200@gmail.com](mailto:daniel.gallardogomez200@gmail.com)

## Abstract

This tutorial focuses on standardised mean differences (SMD) as effect measures in meta-analyses. We will explain what they are, when they should be used, how to correctly compute and interpret them, and some of the most common error made within evidence synthesis.

## Graphical abstract

COCHRANE EVIDENCE SYNTHESIS AND METHODS

Click here to read the article

This tutorial is accompanying an article published in *Cochrane Evidence Synthesis and Methods* journal. If you want to read the article before taking the tutorial, click on the banner above.

## Standardised mean differences (SMD) in meta-analysis

SMD can be used as an effect measure when you have several studies measuring a concept, e.g. depression, but they are not all using the same measurement scale.

In some cases, such as studies measuring weight in pounds vs kilograms, it is straightforward to convert one unit into the other. In other cases, such as two different scales for measuring depression, the relative 'weight' of points on each scale is more qualitative, and it's not valid to mathematically convert measures this way.

In these cases, you can standardize the results of each study using its standard deviations, so that you can compare scores across studies in the review.



This tutorial includes two practicals, each consisting of a video, followed by a brief exercise.

Click on the buttons to take the practicals.

Use the SMD calculator

Interpret an SMD

This tutorial focuses on standardised mean differences (SMD) as effect measures in meta-analyses. We will explain what they are, when they should be used, how to correctly compute and interpret them, and some of the most common error made within evidence synthesis. Additionally, we have created a micro-learning module to accompany this article which includes a real-life example of meta-analysis of SMDs. Moreover, there is the opportunity to test out how to use a tool to correctly calculate SMDs, and how to interpret them. Standardised mean differences micro-learning module.

## INTRODUCTION

In this tutorial, we focus on **standardised meandifferences** (SMD); what they are, when they should be used, how to correctly compute and interpret them, and some of the common errors made within systematic reviews.

Review authors use the SMD as a summary statistic in meta-analyses of continuous outcomes when the studies all assess the same outcome but measure it in a variety of ways [1]. For example, we can look at a meta-analysis in which included studies have measured the depressive symptoms of their participants using different scales/questionnaires [2] (e.g., the Beck Depression Inventory, the Geriatric Depression Scale, the Hamilton Rating Scale, or the Montgomery-Asberg Depression Scale). From a statistical perspective, it is virtually impossible to quantitatively synthesise the results directly. Instead, a solution would be standardising all these data to a common effect size measure: the SMD (Figure 1).

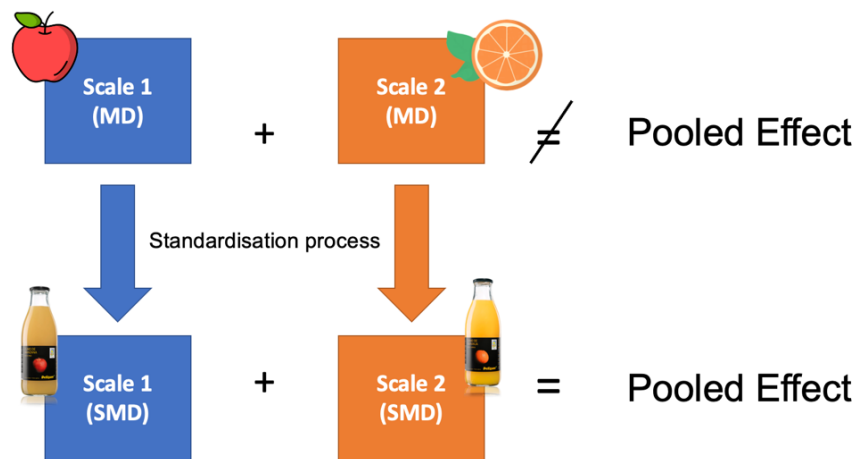


Figure 1. How SMDs enable us to combine different scale-specific units into one pooled effect. The illustration shows how data standardisation process allows us to mix ‘apples’ and ‘oranges’ making them ‘juice’. MD: Mean Difference.

## HOW TO COMPUTE STANDARDISED MEAN DIFFERENCES

When we standardise data, we divide the mean difference (MD) between the treatment and control groups (i.e., the effect size of the treatment) by the pooled sample standard deviation (SD) in each study (i.e., the between-participant variability in outcome measurements observed in each study) at one specific follow-up time point [3].

$$SMD = \frac{\text{MD between groups}}{\text{SD of outcome among participants at follow-up time point}}$$

Equation 1. SMD calculation using pooled sample SD at a specific follow-up time point.

For a better understanding of this terminology, we are going to apply different standardisation methods on data extracted from a published meta-analysis [4]. Therefore, we select the Ortiz-Alonso et al. (2020) study [5] included in this meta-analysis, which reported results using the overall score of the Short Physical Performance Battery (an instrument to assess the physical function; SPPB), and extract the ‘raw’ data (i.e., data directly extracted from the study without any transformation) (Table 1).

Table 1. ‘Raw’ data extracted from Ortiz-Alonso et al. (2020) study.

Study	Scale	Study-arm	Baseline	Baseline	Baseline	Post-treatment	Post-treatment	P
			N	Mean	SD	N	Mean	S
Ortiz-Alonso, 2020	SPPB	Treatment	143	4	2.5	125	3.2	2.
		Control	125	4.2	3.1	125	3.8	2.

Notes. N: Sample size; SD: Standard Deviation.

Next, using Equation 1, we standardise the data using the pooled SD of the outcome among participants at

the post-treatment time point. First, we calculate both pooled sample SDs (baseline and post-treatment)

$$SD_{\text{pooled}} = \sqrt{\frac{SD_t^2 * (n_t - 1) + SD_c^2 * (n_c - 1)}{n_t + n_c - 2}}$$

Equation 2. Calculation of a pooled sample SD. The ‘t’ suffix indicates treatment, and the ‘c’ suffix refers to control arm.

Pooled sample SD at baseline time point:

$$SD_{\text{pooled}} = \sqrt{\frac{2.5^2 * (143 - 1) + 3.1^2 * (125 - 1)}{143 + 125 - 2}} = 2.796$$

Pooled sample SD at post-treatment time point:

$$SD_{\text{pooled}} = \sqrt{\frac{2.5^2 * (125 - 1) + 2.9^2 * (125 - 1)}{125 + 125 - 2}} = 2.707$$

Next, we convert arm-based data into contrast-based data (i.e., a single effect measure that summarises the MD between the two study-arms) using Equations 3 and 4.

$$MD = \text{Mean score}_{\text{treatment}} - \text{Mean score}_{\text{control}}$$

Equation 3. MD computation at post-treatment time point.

$$SE_{\text{MD}} = \sqrt{\frac{n_t + n_c}{n_t * n_c} * \frac{SD_t^2 * (n_t - 1) + SD_c^2 * (n_c - 1)}{n_t + n_c - 2}}$$

Equation 4. Computation of the SE of the MD at post-treatment time point.

$$MD = 3.2 - 3.8 = -0.6$$

$$SE_{\text{MD}} = \sqrt{\frac{125 + 125}{125 * 125} * \frac{2.5^2 * (125 - 1) + 2.9^2 * (125 - 1)}{125 + 125 - 2}} = 0.342$$

Then, we standardise our MD and SE dividing them by the corresponding pooled sample SDs. Methodologists support the use of pooled SDs at baseline over follow-up SDs, but it is common that studies only report follow-up data. Therefore, we are going to standardise data in both cases: (1) supposing that we have baseline data; and (2) supposing that we only have follow-up data.

Standardised data (SMD and SE) using the pooled sample SD at baseline time point.

$$SMD = \frac{MD}{SD_{\text{pooled}}} = \frac{-0.6}{2.796} = -0.215$$

$$SE_{\text{SMD}} = \frac{SE}{SD_{\text{pooled}}} = \frac{0.342}{2.796} = 0.122$$

Standardised data (SMD and SE) using the pooled sample SD at post-treatment time point.

$$SMD = \frac{MD}{SD_{pooled}} = \frac{-0.6}{2.707} = -0.222$$

$$SE_{SMD} = \frac{SE}{SD_{pooled}} = \frac{0.342}{2.707} = 0.126$$

Although this method is the most common applied in meta-analyses, the use of a fixed scale-specific SD reference is recommended [6,7]. A more-in-depth explanation of this method can be found in Gallardo-Gómez et al. (2023) [3] and the online content.

## HOW TO INTERPRET STANDARDISED MEAN DIFFERENCES

The SMDs express the size of the treatment effect in each study relative to the variability observed in that study. However, the overall treatment effect could be difficult to interpret as it is reported in units of standard deviation rather than the original units of measurement. Without guidance, clinicians and patients may have little idea how to interpret results presented as SMDs. There are two possibilities for **re-expressing** such results in more helpful ways:

Re-expressing SMDs using **rules of thumb** for effect sizes. One example based on Cohen (1998) [8] is as follows: 0.2 represents a small effect; 0.5 a moderate effect; and 0.8 a large effect. Nevertheless, some methodologists believe that such interpretations are problematic because the importance of a finding is context-dependent and not amenable to generic statements [7].

Re-expressing SMDs using a **familiar instrument**. The second (and recommended) option is to re-express the SMD in the units of one or more of the specific measurement instruments. This method could be performed by multiplying the SMD by a typical among-person SD for a particular scale (e.g., an external SD reference from a large cohort or cross-sectional study that matches the target population, an internal SD reference, or a pooled sample SD), preferably, the same used for data standardisation [3]. In this way, using the original scale-specific units, the clinical relevance and impact of a pooled treatment effect can be interpreted more easily. In our example, when authors pooled all effect sizes, they obtained a pooled treatment effect of  $SMD = 0.40$  (95% Confidence Interval 0.02 to 0.77). We then re-express this effect size into SPPB units multiplying by the external SD reference for the study population (external SD reference = 3.14), obtaining a scale-specific pooled effect of  $MD = 0.97$  (95% CI 0.06 to 2.42). Considering a predefined minimally clinically important difference of 1 point in the SPPB [9], we could support the use of an intervention (physical activity in this case [4]) in a specific population due to its *clinically* meaningful benefit in the outcome of interest.

## COMMON PITFALLS USING STANDARDISED MEAN DIFFERENCES

1. **Unnecessary data standardisation**. Reviewers do not need to standardise their data when there are not different scales assessing the outcome of interest. The belief that the term ‘effect size’ is a synonym of ‘SMD’ can lead to authors reporting the treatment effect in SMD units when it is not needed. One example of this is when only one study is reported in a forest plot; an SMD is not needed, and this should be reported as MD.
2. **Use of SEs rather than SDs to calculate SMDs**. As we have seen in the Equation 1, we use the post-treatment pooled sample SD to calculate SMDs. Nonetheless, primary studies could wrongly report the SE of an assessment as the SD or not specify whether they are reporting SD or SE. A red flag for this could be a quite low  $SD$  (i.e.,  $<1$ ), though it is highly dependent on the score range of the specific scale. This mistake could lead to ‘effect size inflation’ because when you use SEs to calculate SMDs, you are dividing the MD by a lower value of the truly corresponding one, obtaining a higher value. Therefore, if you obtain SMDs greater than one, you should check whether the SD or SE has been used.
3. **Combination of change from baseline and post-treatment effect measures**. Although mixing change from baseline and post-treatment outcomes is not a problem when it comes to meta-analysis

of MDs [7], they should not in principle be combined using SMDs. This is because the SDs used in standardising post-treatment values reflect between-person variability at a single time point, where SDs used in change scores standardization reflect variation in between-person changes over time, so will depend on both within-person (dependent on the length of time between measurements) and between-person variability [7].

4. **Effect size direction** . There are scales where an improvement in the outcome is reflected by a reduction in the score (e.g., in our illustrative example, the less time spent in walking a distance, the better functional capacity). In addition, to interpret the magnitude of an effect, we must consider the specific outcome (e.g., a more negative effect could be positive if the review is investigating depressive symptoms, meaning a reduction in these symptoms). To correct an effect that is not in line with the direction of our meta-analysis, we should multiply the effect size value by  $-1$ , (no modifications are needed for the SD), ensuring that all effects are in the same direction.
5. **No interpretation of SMDs** . A huge number of meta-analyses often leave their effect estimates as SMDs, which can make interpretation difficult. We have talked about different available options to re-express SMDs to more-interpretable estimates above.

---

**RECOMMENDATIONS** To promote replicability and transparency in evidence synthesis, provide details about the dire

---

COCHRANE EVIDENCE SYNTHESIS AND METHODS

Open Access

Click here to read the article


This tutorial is accompanying an article published in *Cochrane Evidence Synthesis and Methods* journal. If you want to read the article before taking the tutorial, click on the banner above.

## Standardised mean differences (SMD) in meta-analysis

SMD can be used as an effect measure when you have several studies measuring a concept, e.g. depression, but they are not all using the same measurement scale.

In some cases, such as studies measuring weight in pounds vs kilograms, it is straightforward to convert one unit into the other. In other cases, such as two different scales for measuring depression, the relative 'weight' of points on each scale is more qualitative, and it's not valid to mathematically convert measures this way.

In these cases, you can standardize the results of each study using its standard deviations, so that you can compare scores across studies in the review.



This tutorial includes two practicals, each consisting of a video, followed by a brief exercise.

Click on the buttons to take the practicals.

Use the SMD calculator ●

Interpret an SMD ●

Figure 2. Screenshot of the micro-learning module

## ADDITIONAL INFORMATION

The particular definition used in Cochrane Reviews and in this tutorial is the effect size known in social sciences as **Hedges'  $g$** , which is similar to the effect size so-called **Cohen's  $d$**  with a small-sample correction. Hedges'  $g$  uses a pooled SD in the denominator, which is an estimate of the SD based on outcome data from

both intervention groups, assuming that the SDs in the two groups are similar [7]. In contrast, **Glass’ delta** ( $\Delta$ ) uses only the SD from the comparator group, on the basis that if the experimental treatment affects between-person variation, then such an impact of the treatment should not influence the effect estimate.

All these effect measures referred to as SMDs can be calculated by hand or in any statistical package. Statistical packages in R software include *metafor* [10], *esc* [11], or *compute.es* [12]. A hands-on useful resource is *thebookdown* of Harrer et al. (2021) [13], which serves an accessible introduction into how meta-analyses are covered, including different SMD calculation, and pooling methods with examples.

## FURTHER READING AND ONLINE CONTENT

More information on SMDs, can be found in Chapter 6.5 of The Cochrane Handbook for Systematic Reviews of Interventions [1].

Cochrane Training have produced a micro-learning module on how to calculate SMDs to accompany this article [link: <https://share.gomolearning.com/sharelink/b17d5bf8ee76fd1056d6a2505eb81375793889d0773d0141fd/>]. [Figure 2].

## AUTHOR CONTRIBUTIONS

**Daniel Gallardo Gomez** : Conceptualization; writing—original draft; writing—review and editing. **Kerry Dwan** : Conceptualization; supervision; writing—review and editing. **Rachel Richardson**: Supervision; writing-review and editing.

## ACKNOWLEDGMENTS

The authors would like to thank Dario Sambunjak, Development Directorate, Cochrane Central Executive, who designed and built the micro-learning module that accompanies this paper.

## CONFLICT OF INTEREST STATEMENT

Rachel Richardson is employed by Cochrane. Kerry Dwan is a former employee of Cochrane.

## DATA AVAILABILITY STATEMENT

Data sharing not applicable to this article as no data sets were generated or analysed during the current study.

## ORCID

Daniel Gallardo Gómez <http://orcid.org/0000-0002-3029-026X>

Kerry Dwan <http://orcid.org/0000-0001-6918-1215>

Rachel Richardson <http://orcid.org/0000-0003-2848-6260>

## REFERENCES

- <sup>1</sup>Higgins JPT, Li T, Deeks JJ (editors). Chapter 6: Choosing effect measures and computing estimates of effect. In: Higgins JPT, Thomas J, Chandler J, Cumpston M, Li T, Page MJ, Welch VA (editors). *Cochrane Handbook for Systematic Reviews of Interventions* version 6.4 (updated August 2023). Cochrane, 2023. Available from [www.training.cochrane.org/handbook](http://www.training.cochrane.org/handbook).
- <sup>2</sup>Noetel M, Sanders T, Gallardo-Gómez D, Taylor P, del Pozo B, van den Hoek D, et al. Running from Depression: A Systematic Review and Network Meta-Analysis of Exercise Dose and Modality in the Treatment for Depression. Available at SSRN: <https://ssrn.com/abstract=4388153> or <http://dx.doi.org/10.2139/ssrn.4388153>
- <sup>3</sup>Gallardo-Gómez D, Pedder H, Welton NJ, et al. Variability in meta-analysis estimates of continuous outcomes using different standardization and scale-specific re-expression methods. *J Clin Epidemiol.* 2023 Nov. doi: 10.1016/j.jclinepi.2023.11.003.

- <sup>4</sup>Gallardo-Gómez D, del Pozo-Cruz J, Pedder H, et al. Optimal dose and type of physical activity to improve functional capacity and minimize adverse events in acutely hospitalised older adults: a systematic review with dose-response network meta-analysis of randomised controlled trials. *British Journal of Sports Medicine*.2023;57: 1272-1278.
- <sup>5</sup>Ortiz-Alonso J, Bustamante-Ara N, Valenzuela PL, et al. Effect of a simple exercise programme on hospitalization-associated disability in older patients: a randomised controlled trial. *Geriatric Medicine*. 2019. doi:10.1101/19008151.
- <sup>6</sup>Daly C, Welton NJ, Dias S, Anwer S, Ades AE. Meta-Analyses of Continuous Outcomes: Guideline Methodology Document 2. NICE Guidelines Technical Support Unit, 2021. 49 p.
- <sup>7</sup>Higgins JPT, Thomas J, Chandler J, Cumpston M, Li T, Page MJ, et al. Cochrane Handbook for Systematic Reviews of Interventions. John Wiley & Sons; 2019.
- <sup>8</sup>Cohen, J. Statistical power analysis for the behavioral sciences (2<sup>nd</sup> ed.) 1998. Hillsdale, NJ: Erlbaum.
- <sup>9</sup>Perera S, Mody SH, Woodman RC, et al. Meaningful change and responsiveness in common physical performance measures in older adults. *Journal of the American Geriatrics Society*.2006;54(5): 743-749.
- <sup>10</sup>Viechtbauer W. Conducting Meta-Analyses in R with the metafor Package. *J Stat Softw* 2010;36: 1–48.
- <sup>11</sup>Lüdtke D, Lüdtke MD, Calculator' from David BW. Package "esc" 2017.
- <sup>12</sup>AC Del Re. compute.es: Compute Effect Sizes. 2013. R package version 0.2-2- URL <https://cran.r-project.org/package=compute.es>.
- <sup>13</sup>Harrer M, Cuijpers P, Furukawa TA, Ebert DD. *Doing Meta-Analysis with R: A Hands-On Guide* . 2021. Boca Raton, FL and London: Chapman & Hall/CRC Press. ISBN 978-0-367-61007-4.