

# BMJ Open Plea for routinely presenting prediction intervals in meta-analysis

Joanna IntHout,<sup>1</sup> John P A Ioannidis,<sup>2,3,4,5</sup> Maroeska M Rovers,<sup>1</sup> Jelle J Goeman<sup>1</sup>

**To cite:** IntHout J, Ioannidis JPA, Rovers MM, *et al.* Plea for routinely presenting prediction intervals in meta-analysis. *BMJ Open* 2016;**6**:e010247. doi:10.1136/bmjopen-2015-010247

► Prepublication history and additional material is available. To view please visit the journal (<http://dx.doi.org/10.1136/bmjopen-2015-010247>).

Received 12 October 2015  
Revised 7 March 2016  
Accepted 14 April 2016



CrossMark

<sup>1</sup>Radboud University Medical Center, Radboud Institute for Health Sciences (RIHS), Nijmegen, The Netherlands

<sup>2</sup>Department of Medicine, Stanford Prevention Research Center, Stanford University School of Humanities and Sciences, Stanford, California, USA

<sup>3</sup>Department of Health Research and Policy, Stanford University School of Medicine, Stanford, California, USA

<sup>4</sup>Department of Statistics, Stanford University School of Humanities and Sciences, Stanford, California, USA

<sup>5</sup>Meta-Research Innovation Center at Stanford (METRICS), Stanford University, Stanford, California, USA

## Correspondence to

Dr Joanna IntHout; Joanna.IntHout@radboudumc.nl

## ABSTRACT

**Objectives:** Evaluating the variation in the strength of the effect across studies is a key feature of meta-analyses. This variability is reflected by measures like  $\tau^2$  or  $I^2$ , but their clinical interpretation is not straightforward. A prediction interval is less complicated: it presents the expected range of true effects in similar studies. We aimed to show the advantages of having the prediction interval routinely reported in meta-analyses.

**Design:** We show how the prediction interval can help understand the uncertainty about whether an intervention works or not. To evaluate the implications of using this interval to interpret the results, we selected the first meta-analysis per intervention review of the Cochrane Database of Systematic Reviews Issues 2009–2013 with a dichotomous (n=2009) or continuous (n=1254) outcome, and generated 95% prediction intervals for them.

**Results:** In 72.4% of 479 statistically significant (random-effects  $p < 0.05$ ) meta-analyses in the Cochrane Database 2009–2013 with heterogeneity ( $I^2 > 0$ ), the 95% prediction interval suggested that the intervention effect could be null or even be in the opposite direction. In 20.3% of those 479 meta-analyses, the prediction interval showed that the effect could be completely opposite to the point estimate of the meta-analysis. We demonstrate also how the prediction interval can be used to calculate the probability that a new trial will show a negative effect and to improve the calculations of the power of a new trial.

**Conclusions:** The prediction interval reflects the variation in treatment effects over different settings, including what effect is to be expected in future patients, such as the patients that a clinician is interested to treat. Prediction intervals should be routinely reported to allow more informative inferences in meta-analyses.

## INTRODUCTION

Interventions may have heterogeneous effects across studies because of differences in study populations, interventions, follow-up length or other factors like publication bias.<sup>1</sup> Nevertheless, the usual reporting of a meta-analysis is focused on the summary effect size combined with a CI and p value. Typically also some measure of the between-study heterogeneity is presented such as  $\tau^2$  or the

## Strengths and limitations of this study

- In many meta-analyses, there is large variation in the strength of the effect.
- The prediction interval helps in the clinical interpretation of the heterogeneity by estimating what true treatment effects can be expected in future settings.
- In case of heterogeneity, prediction intervals will show a wider range of expected treatment effects than CIs, and thus may lead to different conclusions. This occurred in over 70% of statistically significant meta-analyses with heterogeneity of the Cochrane Database of Systematic Reviews. Completely opposite effects were not excluded in over 20% of those meta-analyses.
- Prediction intervals should be routinely reported to allow more informative inferences in meta-analyses.
- Limitations are that the calculations and inferences for the prediction interval are based on the normality assumption, which is difficult to ensure. Further, the interval will be imprecise if the estimates of the summary effect and the between-study heterogeneity are imprecise, for example, if they are based on only a few, small studies. Inferences based on the prediction interval are only valid for settings that are similar (exchangeable) to those on which the meta-analysis is based.

inconsistency measure  $I^2$ .<sup>2,3</sup> However, neither of these two metrics can readily point to the clinical implications of the observed heterogeneity. Our objective in the current article is to show the potential advantages of obtaining and reporting the prediction interval routinely in meta-analyses because its clinical meaning is much more straightforward. The prediction interval presents the heterogeneity in the same metric as the original effect size measure, in contrast to  $\tau^2$  or  $I^2$ . Reporting a prediction interval in addition to the summary estimate and CI will illustrate which range of true effects can be expected in future settings. We describe its merits and provide working examples to show how it can be calculated.

## METHODS

### Interpretation of heterogeneity

Between-study variation in the magnitude of treatment effects cannot be neglected. One of the main merits of a meta-analysis may even be that it reveals the variation of effects in different studies.<sup>4</sup> Therefore, summarising the findings of a meta-analysis in a single summary value sacrifices potentially informative variation.<sup>5</sup> However, the information that can be directly retrieved from  $\tau^2$  and  $I^2$  with respect to the variation in the effects is limited. The clinical interpretation of  $I^2$  is ambiguous: a high  $I^2$  does not necessarily imply that the study effects are dispersed over a wide range<sup>6</sup> and a low  $I^2$  might correspond to high dispersion,<sup>7</sup> because  $I^2$  depends on sample size of the included studies.<sup>8</sup> With very large (highly precise) studies, even tiny differences in effect size may result in a high  $I^2$ , while with small (imprecise) studies, very different treatment effects can yield an  $I^2$  of 0. Dispersion in treatment effects is better reflected by  $\tau$  because  $\tau$  is the SD of the between-study effects. One could, for example, estimate the ratio of the effect size over  $\tau$ , which can convey how many times larger the treatment effect is compared with the SD of the effect across studies.<sup>9</sup> But this may still be not very intuitive to a clinical reader. Another popular way to express variation in effect sizes is the CI, for example, the 95% CI. The CI in a random-effects model contains highly probable values for the summary treatment effect. However, it does not convey what range of treatment effects are likely to be seen in other patients, for example, in the next study or in the patients a clinician wants to treat in her clinic.

### Prediction intervals

Not so often reported but much more insightful is the prediction interval.<sup>10</sup> A prediction interval always presents the heterogeneity on the same scale as the original outcomes, in contrast to  $\tau$  (eg, in case of ORs),  $\tau^2$  or  $I^2$ . A 95% prediction interval estimates where the true effects are to be expected for 95% of similar (exchangeable) studies that might be conducted in the future.<sup>4</sup> Therefore, it is well suited to evaluate the variability of the effect of an intervention over different settings. For example, in a meta-analysis on sedentary time in adults and the association with diabetes, cardiovascular disease and death, CIs were thought to represent insufficiently the different study populations. Therefore, also prediction intervals were reported.<sup>11</sup> In the absence of between-study heterogeneity, the prediction interval coincides with the respective CI. However, in case of heterogeneity, a prediction interval covers a wider range than a CI. Consequently, in case of a statistically significant effect (where all values of the 95% CI are on the same side of the null), the corresponding 95% prediction interval may indicate that values are possible on both sides of the null. This means that there will be settings where conclusions based on CIs will not hold. In the same framework, one can also calculate the probability that the true effect will be harmful (on the other

side of the null) in a next study. Table 1 presents an overview of measures of between-study heterogeneity.

### Example: topical steroids for nasal polyps

A 2012 review on the use of topical steroids for treatment of chronic rhinosinusitis with nasal polyps, based on seven randomised studies, resulted in a larger decrease in overall symptom scores in favour of steroids compared with placebo.<sup>12</sup> This is reflected by a standardised mean difference (SMD) of  $-0.51$ , with a 95% CI  $-0.96$  to  $-0.07$  (figure 1). The  $I^2$  is 73.9% (95% CI 44.2% to 87.8%), which can be considered substantial heterogeneity,<sup>13</sup> and the estimated  $\tau^2$  is 0.148. Notwithstanding these numbers, it is difficult to evaluate what the clinical consequences of this heterogeneity may be for future settings.

In order to estimate the prediction interval for the SMD, we need the point estimate of the SMD, its SE and the estimated  $\tau^2$ . We derive the SE from the 95% CI of the SMD (see online supplementary appendix formula 1), which results in an SE of 0.227. We can calculate the SD of the prediction interval  $SD_{PI}$  as  $\sqrt{(0.148+0.227^2)}$  and the lower and upper limit of the 95% prediction interval as  $-0.51 \pm 2.45 \times SD_{PI}$ . The value 2.45 results from the  $t_{1-0.05/2,6}$  distribution. Prediction intervals with a different coverage could be calculated by using a different t-value, for example,  $t_{1-0.20/2,6}$  for an 80% prediction interval (see online supplementary appendix formula 1).

The resulting prediction interval, ranging from  $-1.60$  to  $0.58$ , can be interpreted as the 95% range of true SMDs to be expected in similar studies. We present it in figure 1 as a rectangle below the diamond for the 95% CI.<sup>14</sup> The prediction interval contains values below zero, which correspond to a decrease in symptom scores of at best  $\sim 1.6$  SD after steroid use compared with placebo. But it also contains values above zero which means that the steroids may exhibit no or even a harmful effect (SMD $>0$ ) in some settings, with a (95%) worst case increase in SMD of 0.58. Consequently, the effect in a new study may be even the exact opposite to the summary point estimate of the meta-analysis, that is, an increase of 0.51 instead of a decrease of  $-0.51$  may occur. The estimated probability that the true effect of the steroids will be null or higher in a new study is equal to 14.7%, based on the t-distribution with 6 degrees of freedom (see online supplementary appendix formula 2).

### Cochrane database

In order to investigate how often there is a discrepancy in conclusions based on prediction intervals and CIs, we evaluated this in statistically significant meta-analyses ( $p<0.05$  by random-effects calculations) of the Cochrane Database of Systematic Reviews Issues 2009–2013, kindly provided by the UK Cochrane Editorial Unit. To avoid subjectivity in the selection, we used the first meta-analysis with a dichotomous or continuous outcome and based on at least two studies in the data and analyses section when these studies were also combined in the original

**Table 1** Some frequently used measures for heterogeneity

Measure	Advantages	Disadvantages
$\tau^2$	<ul style="list-style-type: none"> <li>▶ The <math>\tau</math> (the square root of <math>\tau^2</math>) is the SD of the between-study variation on the scale of the original outcome.</li> <li>▶ The <math>\tau^2</math> is the direct estimate of the between-study variation and therefore useful in calculations, for example, for the prediction interval.</li> </ul>	<ul style="list-style-type: none"> <li>▶ A direct clinical interpretation based on <math>\tau^2</math> is difficult, especially when <math>\tau^2</math> belongs to outcomes that were analysed on log scale, for example, ORs.</li> <li>▶ When the <math>\tau^2</math> estimate is based on only a few studies, it will be imprecise.</li> </ul>
$I^2$	<ul style="list-style-type: none"> <li>▶ <math>I^2</math> presents the inconsistency between the study results and quantifies the proportion of observed dispersion that is real, that is, due to between-study differences and not due to random error.<sup>2,3</sup></li> <li>▶ <math>I^2</math> reflects the extent of overlap of the CIs of the study effects.</li> <li>▶ <math>I^2</math> represents the inconsistency always on a scale between 0 and 100, therefore it can be compared with suggested limits for low or high inconsistency.<sup>13</sup></li> </ul>	<ul style="list-style-type: none"> <li>▶ A direct clinical interpretation of <math>I^2</math> is difficult.</li> <li>▶ <math>I^2</math> is also ambiguous because its size depends on sample size: <ul style="list-style-type: none"> <li>– With very large studies, even tiny between-study differences in effect size may result in a high <math>I^2</math>;</li> <li>– With small (imprecise) studies, very different treatment effects can yield an <math>I^2</math> of 0.</li> </ul> </li> </ul>
CI	<ul style="list-style-type: none"> <li>▶ The CI in a random-effects model contains highly probable values for the summary (mean) treatment effect.</li> </ul>	<ul style="list-style-type: none"> <li>▶ The CI gives no information on the range of true treatment effects that are likely to be seen in other settings, for example, in the next study or in the patients a clinician wants to treat in her clinic.</li> </ul>
Prediction interval	<ul style="list-style-type: none"> <li>▶ The prediction interval in a random-effects model contains highly probable values for the true treatment effects in future settings, if those settings are similar to the settings in the meta-analysis.</li> <li>▶ The values in the interval can be compared with clinically relevant thresholds to see whether they correspond to benefit, null effects or harm.</li> <li>▶ The prediction interval can be used to estimate the probability that the treatment in a future setting will have a true-positive or true-negative effect, and to perform better power calculations.</li> </ul>	<ul style="list-style-type: none"> <li>▶ Conclusions drawn from the prediction interval are based on the assumption that <math>\tau^2</math> and the study effects are normally distributed.</li> <li>▶ The estimate of the prediction interval will be imprecise if the estimates of the summary effect and the <math>\tau^2</math> are imprecise, for example, if they are based on only a few studies and if these studies are small.</li> </ul>

review, as we wanted to reflect the status quo as precise as possible. Details can be found in another paper.<sup>15</sup> In brief, of a total of 3263 meta-analyses, 920 were statistically significant: 479 with an estimated  $I^2 > 0$  and 441 with an estimated  $I^2 = 0$ .

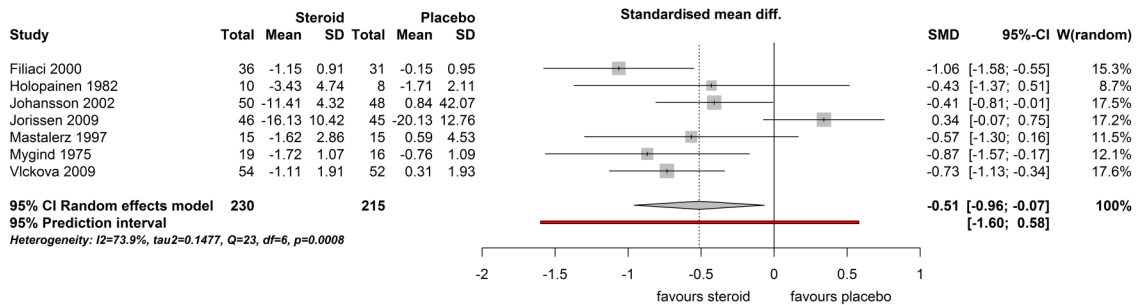
### Calculations

We used the Hartung-Knapp/Sidik-Jonkman<sup>16</sup> (HKSJ) random-effects meta-analysis approach combined with the empirical Bayes estimator for  $\tau^2$ . We estimated  $\tau^2$  for all meta-analyses, even when the authors originally performed a fixed-effects analysis. Prediction intervals were calculated according to online supplementary appendix formula 1). We categorised the statistically significant meta-analyses with heterogeneity ( $\tau^2 > 0$ ) by number of studies (2–6 studies or  $> 6$ ) and heterogeneity ( $I^2 < 30\%$ , 30% to 60% or  $> 60\%$ , based on the Cochrane Handbook<sup>13</sup> stating that an  $I^2$  between 30% and 60% corresponds to moderate heterogeneity). For significant meta-analyses where the heterogeneity estimate was zero, we assessed the impact of possibly low but non-zero heterogeneity by assuming an  $I^2$  of 20%, calculating prediction intervals using online supplementary appendix formula 3). Categorical outcomes were compared

between groups by means of the  $\chi^2$  test. We used R software (R: A language and environment for statistical computing. Retrieved from <http://www.R-project.org/>. [program]. Vienna, Austria: R Foundation for Statistical Computing, 2014) V.3.1.2 and the R packages metafor<sup>17</sup> V.1.9-5 and meta (meta: General Package for Meta-Analysis. R package version 4.1-0. <http://CRAN.R-project.org/package=meta> [program], 2015) V.4.1-0.

### RESULTS

Overall, 132 (27.6%) of the 479 statistically significant meta-analyses with an  $I^2 > 0$  had both the 95% CI and the 95% prediction interval excluding the null effect (table 2). Consequently, almost three-quarter (347, 72.4%) had a prediction interval that contained the null effect. This means that it is likely that for these comparisons, some patient populations might experience null effects or effects in the opposite direction, that is, a treatment might be more harmful than the comparator even though the point estimate suggests benefit (or vice versa). Not surprisingly, significant meta-analyses with low heterogeneity more often had prediction intervals that excluded the null than meta-analyses with high



**Figure 1** Forest plot of the standardised mean difference (SMD) in symptom scores in nasal polyps. Steroids versus placebo, analysis 1.1 in Cochrane Review CD006549.<sup>12</sup> Note that our results differ from the original analysis, as we used a random-effects analysis with the Hartung-Knapp/Sidik-Jonkman adjustment<sup>16</sup> and the empirical Bayes estimator for  $\tau^2$ .

heterogeneity. The percentage of prediction intervals containing the null effect was slightly higher for meta-analyses with a continuous outcome (80.4%) than for those with a dichotomous outcome (65.8%;  $p<0.001$ ), but not significantly different for meta-analyses based on more than six studies (74.1%) than for those with at most six studies (69.1%;  $p=0.25$ ; web table W1).

Of the 347 meta-analyses with a prediction interval that contained the null or opposite effect, 199 (57.3%) had also at least one study with an opposite effect. This happened more often in meta-analyses with more than six studies (181/235, 77.0%) than in those based on at most six studies (18/102, 17.6%). Especially in meta-analyses with few studies and substantial heterogeneity, the prediction interval was wider than the range of study outcomes. The opposite (ie, a smaller prediction interval) occurred in meta-analyses based on many studies and with low estimated heterogeneity. Results for meta-analyses with dichotomous and continuous outcomes were not notably different.

### Prediction intervals containing the opposite effect

If the prediction interval just includes the null effect, this may be less worrying than when it contains the exact opposite effect of the pooled summary effect, for example, if it contains an OR of 0.5 when the meta-analysis summary estimate is an OR of 2, or if it contains an SMD of  $-0.7$  when the summary estimate was 0.7. Of the 479 significant meta-analyses with an  $I^2>0.97$  (20.3%) had a prediction interval that contained the opposite effect. This percentage was higher for the meta-analyses with a continuous outcome (65/219, 29.7%) than for those with a dichotomous outcome (32/260, 12.3%;  $p<0.001$ ). It occurred also more frequently in

meta-analyses with more than six primary studies (57/139, 41.0% and 30/178, 20.3% for meta-analyses with a continuous or dichotomous outcome, respectively) than for those based on at most six studies (8/80, 10.0% and 2/82, 2.4%;  $p<0.001$  and  $p=0.001$ , respectively).

### Meta-analyses with estimated $I^2=0$

A substantial part of meta-analyses have an estimated  $I^2$  of 0. However, there is typically very large uncertainty about the exact amount of heterogeneity, and this is demonstrated by very large 95% CIs for the values of  $I^2$ .<sup>18</sup> The same applies to  $\tau$ : an estimate of 0 is often accompanied by large uncertainty. The true  $I^2$  and  $\tau$  are unlikely to ever be exactly 0, although low values are possible. To assess the impact of possibly low but non-zero heterogeneity among the 441 Cochrane meta-analyses with estimated  $I^2=0$  and statistically significant results, we imputed an  $I^2=20%$  (suggestive of low between-study heterogeneity). Under this assumption, in 329 (74.6%) of these 441 meta-analyses the 95% prediction interval would span both sides of the null (table 2), similar for meta-analyses with a dichotomous (74.7%) or continuous (74.4%) outcome (web table W1). This is a sensitivity analysis that is useful to perform to see whether the inferences of a meta-analysis that seemingly does not have detectable heterogeneity may be influenced by even a small amount of heterogeneity.

## DISCUSSION AND OUTLOOK

In meta-analyses, a CI is inadequate for clinical decision-making because it only summarises the average effect for the average study. The prediction interval is more informative as it shows the range of possible effects in

**Table 2** Proportion of statistically significant meta-analyses where both the 95% CIs and PIs excluded the null

Statistically significant meta-analyses	Estimated heterogeneity $I^2$				
	$I^2=0^*$	$I^2>0$	$>0$ and $<30%$	30% to 60%	$>60%$
N	441	479	123	150	206
Both 95% CI and 95% PI excluded null, n (%)	112 (25.4)	132 (27.6)	88 (71.5)	39 (26.0)	5 (2.4)

\*When the estimated heterogeneity  $I^2$  was equal to 0,  $I^2=20%$  was imputed for the calculation of the PI. PI, prediction interval.

relation to harm and clinical benefit thresholds. While we have focused on the situation where the separating threshold is the null, a different threshold may be considered. For example, in the prediction interval framework, one can calculate the probability that an effect is larger than  $B$ , where  $B$  may be a clinically meaningful effect (if the treatment benefit is less than  $B$ , then it is felt not to be worth it). A narrow prediction interval that lies completely on the beneficial side of a clinically relevant threshold increases confidence in an intervention. A broad prediction interval may indicate the existence of settings where the treatment has a suboptimal and possibly even harmful effect. In more than 70% of statistically significant meta-analyses of the Cochrane Database with some estimated or assumed between-study heterogeneity, the prediction intervals crossed the no-effect threshold, indicating that there are settings where those treatments will have no effect or even an effect in the opposite direction. In 20.3% of those meta-analyses, the prediction interval even contained the opposite effect of the summary estimate, for example, an OR of 0.5 when the summary point estimate was an OR of 2. This occurred most frequently for meta-analyses with a continuous outcome, probably because heterogeneity can be more prominent in many topics where outcomes are assessed on continuous scales; higher heterogeneity for the continuous outcomes was also observed in the full set of 3263 meta-analyses.<sup>15</sup> It was also slightly more common for meta-analyses based on more than six studies, probably because such meta-analyses have more power to detect smaller effects, which means that also the opposite effects will be smaller.

Graham and Moran<sup>19</sup> evaluated prediction intervals in 72 meta-analyses with a dichotomous outcome in critical care published between 2002 and 2010. They found a higher percentage of significant meta-analyses (50/72, 69.4%), compared with 28.5% (572/2009) in our set of meta-analyses with an OR outcome. The difference may be caused by publication bias, the higher number of primary studies in their sample (median 9 vs 4 in our set<sup>15</sup>) and by their use of the DerSimonian-Laird approach which can result in too many statistically significant findings, whereas we used the HKSJ approach.<sup>16</sup> However, results with respect to the prediction interval were remarkably similar. In 32 (64.0%) of their 50 significant meta-analyses, the 95% prediction interval included the null, similar to 65.8% in our data set. Seven (14.0%) of their 50 meta-analyses suggested a high probability of exact reversal of the efficacy or harm, similar to 12.3% of our meta-analyses where the prediction interval contained the opposite effect, despite the fact that they used a different definition for possible 'harm' and that they did not mention whether there was positive between-study heterogeneity in their significant meta-analyses.

It is straightforward to calculate a prediction interval if we can assume that the effects are normally distributed

and that  $\tau^2$  is known and stable across studies. However, one should realise that the prediction interval is dependent on this assumption and on the precisions of the estimated  $\tau^2$  and study effect, and will be imprecise if the number of studies in the meta-analysis is small. If the number of studies is large, estimates will be more precise and the normality of the distribution of  $\tau^2$  can be empirically evaluated. A final caveat is that the uncertainty conveyed by the prediction interval pertains to the uncertainty about the extent to which future studies are similar (exchangeable) to those that have already been done, but this applies to all inferences from a meta-analysis. If the future studies evaluate patients and settings that are entirely different from what was evaluated in past studies, this exchangeability is questionable and uncertainty may be even more prominent than what the prediction interval conveys. In practical terms, if the patients treated by a physician are considered to be very different from the patients seen in all studies that have been done in the past, even the prediction interval cannot tell us what we might expect for these patients.

#### Power calculations for a future study

Meta-analysis results can also be used for power calculations for a new study. However, the expected true effect in a new study is not necessarily equal to the point estimate of the meta-analysis: it can be any of the values in the prediction interval. In case of heterogeneity, the probability of a statistically significant result in a new study may differ substantially from an apparent power of 80% based on the point estimate. The latter will be overly optimistic because the power function is asymmetric. If the true study effect is larger than the point estimate, the real probability of a significant study will be higher, up to a maximum of 100%, but if the effect is smaller, the probability may decrease substantially, even to 5% or less in case of a null effect. Consequently the expected probability of a significant new study in case of heterogeneity will be lower than 80% (online supplementary appendix formula 4). For example, if the prediction interval shows that 30% of future studies may have a true null or negative effect, the probability of a significant new study can never be much larger than 70%. The sample size should be increased to compensate for this loss, see also Roloff *et al.*<sup>20</sup>

Summarising, the prediction interval reflects the variation in true treatment effects over different settings, including what effect is to be expected in future patients, such as the patients that a clinician is interested to treat. Therefore, it should be routinely reported in addition to the summary effect and its CI, and used as a main tool for interpreting evidence, to enable more informed clinical decision-making.

**Contributors** JI originated the idea for this study together with JJG. JI drafted the manuscript and conducted the data analysis. All authors read and critically revised the manuscript for important intellectual content and approved the final manuscript.

**Funding** This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors.

**Competing interests** None declared.

**Provenance and peer review** Not commissioned; externally peer reviewed.

**Data sharing statement** Data sets are available on request from the corresponding author.

**Open Access** This is an Open Access article distributed in accordance with the Creative Commons Attribution Non Commercial (CC BY-NC 4.0) license, which permits others to distribute, remix, adapt, build upon this work non-commercially, and license their derivative works on different terms, provided the original work is properly cited and the use is non-commercial. See: <http://creativecommons.org/licenses/by-nc/4.0/>

## REFERENCES

- Riley RD, Higgins JPT, Deeks JJ. Interpretation of random effects meta-analyses. *BMJ* 2011;342:d549.
- Thompson SG, Higgins JPT. How should meta-regression analyses be undertaken and interpreted? *Stat Med* 2002;21:1559–73.
- Higgins JP, Thompson SG, Deeks JJ, et al. Measuring inconsistency in meta-analyses. *BMJ* 2003;327:557–60.
- Higgins JP, Thompson SG, Spiegelhalter DJ. A re-evaluation of random-effects meta-analysis. *J R Stat Soc Ser A Stat Soc* 2009;172:137–59.
- Saha S, Chant D, McGrath J. Meta-analyses of the incidence and prevalence of schizophrenia: conceptual and methodological issues. *Int J Methods Psychiatr Res* 2008;17:55–61.
- Borenstein M, Hedges LV, Higgins JPT, et al. *Introduction to meta-analysis*. Chichester, UK: Wiley, 2009.
- Melsen WG, Bootsma MCJ, Rovers MM, et al. The effects of clinical and statistical heterogeneity on the predictive values of results from meta-analyses. *Clin Microbiol Infect* 2014;20:123–9.
- Rücker G, Schwarzer G, Carpenter JR, et al. Undue reliance on  $I^2$  in assessing heterogeneity may mislead. *BMC Med Res Methodol* 2008;8:79.
- Moonesinghe R, Khoury MJ, Liu T, et al. Required sample size and nonreplicability thresholds for heterogeneous genetic associations. *Proc Natl Acad Sci USA* 2008;105:617–22.
- Chiolero A, Santschi V, Burnand B, et al. Meta-analyses: with confidence or prediction intervals? *Eur J Epidemiol* 2012;27:823–5.
- Wilmot EG, Edwardson CL, Achana FA, et al. Sedentary time in adults and the association with diabetes, cardiovascular disease and death: systematic review and meta-analysis. *Diabetologia* 2012;55:2895–905.
- Kalish L, Snidvongs K, Sivasubramaniam R, et al. Topical steroids for nasal polyps. *Cochrane Database Syst Rev* 2012;12:CD006549.
- Higgins JPT, Green S, Collaboration C. *Cochrane handbook for systematic reviews of interventions*. Wiley Online Library, 2008.
- Guddat C, Grouven U, Bender R, et al. A note on the graphical presentation of prediction intervals in random-effects meta-analyses. *Syst Rev* 2012;1:34.
- IntHout J, Ioannidis JPA, Borm GF, et al. Small studies are more heterogeneous than large ones: a meta-meta-analysis. *J Clin Epidemiol* 2015;68:860–9.
- IntHout J, Ioannidis JP, Borm GF. The Hartung-Knapp-Sidik-Jonkman method for random effects meta-analysis is straightforward and considerably outperforms the standard DerSimonian-Laird method. *BMC Med Res Methodol* 2014;14:25.
- Viechtbauer W. Conducting meta-analyses in R with the metafor package. *J Stat Software* 2010;36:1–48.
- Ioannidis JPA, Patsopoulos NA, Evangelou E. Uncertainty in heterogeneity estimates in meta-analyses. *BMJ* 2007;335:914–16.
- Graham PL, Moran JL. Robust meta-analytic conclusions mandate the provision of prediction intervals in meta-analysis summaries. *J Clin Epidemiol* 2012;65:503–10.
- Roloff V, Higgins JP, Sutton AJ. Planning future studies based on the conditional power of a meta-analysis. *Stat Med* 2013;32:11–24.