

The choice and definition of summary measure for meta-analysis of clinical studies with binary outcomes: Effect on clinical interpretation

Authors and institutions:

Andrew A Plumb¹

Steve Halligan¹

Susan Mallett²

¹Centre for Medical Imaging, University College London, 43-45 Foley St, London, W1W 7TS

²Institute of Applied Health Research, University of Birmingham, B15 2TT

Correspondence to:

Andrew Plumb, Associate Professor of Medical Imaging, Centre for Medical Imaging, University College London, 43-45 Foley St, London, W1W 7TS.

andrew.plumb@ucl.ac.uk

Article type:

Commentary

Word Count (excluding figures and legends):

2,287

Figures:

4

Tables:

1

1 **The choice and definition of summary measure for meta-analysis of clinical**
2
3 **studies with binary outcomes: Effect on clinical interpretation**
4
5
6
7

8 **Article type:**
9

10 Commentary
11
12
13

14 **Word Count (excluding figures and legends):**
15

16 2,287
17
18
19
20

21 **Figures:**
22

23 4
24
25
26
27

28 **Tables:**
29

30 1
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

BJR UNCORRECTED PROOFS

Abstract:

1
2 Many systematic reviews and meta-analyses concern the effect of a healthcare intervention
3
4 on a binary outcome i.e. occurrence (or not) of a particular event. Usually, the overall effect,
5
6 pooled across all studies included in the meta-analysis, is summarised using the odds ratio
7
8 (OR) or the relative risk (RR). Under most circumstances, it is obvious how to identify what
9
10 should be considered as the event of interest – for example, death or a clinically-important
11
12 side effect. However, on occasion it may not be clear in which “direction” the event should
13
14 be specified – such as attendance (vs. non-attendance) at cancer screening. Usually, this
15
16 choice is not critical to the overall conclusion of the meta-analysis, but occasionally it can
17
18 lead to differences in how the included studies are pooled, ultimately affecting the overall
19
20 meta-analytic result, particularly when using relative risks rather than odds ratios. In this
21
22 commentary, we will explain this phenomenon in more detail using examples from the
23
24 literature, and explore how analysts and readers can avoid some potential pitfalls.
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Commentary:

1
2
3
4 There are a number of interesting examples in the literature whereby the choice of summary
5 measure for a meta-analysis can substantially affect its clinical interpretation¹, sometimes
6 leading the reader to draw different conclusions regarding the statistical significance of the
7 results. A recently published article by Zhu et al² compared the attendance rate at screening
8 CT colonography with that at colonoscopy, correctly identifying via meta-analysis that the
9 risk of *attendance* at screening was not significantly different between colonoscopy than
10 CTC; but the directly opposite scenario, risk of *non-attendance* at screening, was
11 significantly worse for colonoscopy than CTC. This highly counter-intuitive phenomenon,
12 whereby simple reversal of the outcome of interest can lead to apparently different results,
13 has been long-known^{1,3,4}, but is perhaps under-recognised. In this commentary, we aim to
14 explain how this occurs, using examples in the literature where necessary, before making
15 suggestions for how analysts and readers can mitigate the problem.
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30

31
32
33 When considering a dichotomous outcome, researchers have to decide what represents the
34 outcome of interest. Usually, this is obvious – death, myocardial infarction, or diagnosis of
35 cancer are all unambiguous, clearly-defined, clinically-relevant endpoints, and are suitable
36 outcomes for both component primary studies and subsequent meta-analysis. However, in
37 some circumstances, it may be more difficult to define the relevant outcome. For example,
38 when investigating screening, it is arguable whether attendance or non-attendance should
39 be chosen as the outcome. Similarly, studies examining fertility treatment (for example)
40 could choose to define successful conception as the outcome, or non-conception as an
41 adverse outcome and thus the “risk event”. Why does this seemingly arbitrary choice
42 matter? Because, for meta-analysis under certain circumstances, it can be extremely
43 important depending on the analysis method chosen³. Analysing relative risks requires
44 caution when the prevalence of the outcome varies across the component studies.
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 A summary statistic for meta-analysis is generated by pooling the individual estimates of the
2 effects observed in the component primary studies. For binary outcomes, these are usually
3 expressed as a relative risk (RR, also called the risk ratio) or an odds ratio (OR). Although
4 exact methods for pooling component studies individual RRs or ORs vary, in essence meta-
5 analysis assigns a “weight” to each component study based on how precisely the outcome
6 measure can be estimated; and, when random effects meta-analysis is used, an adjustment
7 for variation between studies⁵. The weighting given to an individual study determines its
8 influence over the final meta-analysis pooled summary statistic; larger weightings exert
9 greater effect. For fixed effect meta-analysis, individual study weights are affected primarily
10 by the 95% confidence interval around the point estimate of the RR (or OR) – larger
11 weightings are given to studies with narrower confidence intervals than the other studies in
12 the meta-analysis.

13 This process seems logical – studies whose individual results are more precise should exert
14 more effect on the final meta-analysis outcome measure. In random effects meta-analysis,
15 smaller studies have larger relative weight, as the meta-analysis aims to estimate the
16 average effect across all studies (rather than assuming there is an underlying “standard”
17 effect that should be the same across all studies). Nonetheless, the width of the 95%
18 confidence intervals around the risk estimate from each individual study still influences the
19 final weight assigned even in random effects meta-analyses. Usually, this is not problematic;
20 however, using relative risks can introduce unpredictable behaviour in meta-analysis when
21 the prevalence of the risk event varies across the studies. Specifically, it can result in very
22 different weights being assigned to studies that are otherwise similar. For example, consider
23 the first forest plot presented by Zhu et al² (Figure 2a of their article; redrawn here for
24 convenience).

1 We can see here that the two large studies by Stoop et al⁶ (8844 patients) and Sali et al⁷
2 (5861 patients) receive the largest weightings, of 23.9% and 22.3% respectively. However,
3
4 You et al⁸, which randomised only 131 patients, receives almost the same weighting, at
5
6 21.7%; and more than Scott et al⁹ (weighting 17.4% for a sample of 709 patients) and the
7
8 MACS group¹⁰ (weighting 14.7% for a sample of 429 patients). This would not matter if
9
10 individual study results were identical, but they are not (which, after all, is why we perform
11
12 meta-analysis). For example, the RR for You et al is less than 1.0, but greater than 1.0 for
13
14 Stoop et al. Accordingly, although the two largest studies show a clearly significant effect in
15
16 favour of CTC, the overall meta-analytic point estimate is not significant at the 5% level
17
18 (p=0.07) because it is “dragged down” by the weighting ascribed to smaller studies with
19
20 conflicting findings. How has this apparently counterintuitive situation occurred, whereby a
21
22 small study of 131 patients receives weighting virtually equivalent to a 6000 patient RCT?
23
24
25
26
27
28

29 Relative risks for randomised trials are a simple concept; the probability (or risk) of the
30
31 outcome in one trial arm is compared to the probability of the same outcome in the other
32
33 arm, expressed as a ratio. If 200 of every 1000 patients die with placebo (200/1000 = 0.2)
34
35 and 150 die with treatment (150/1000 = 0.15), the relative risk of death is 0.15 / 0.20 = 0.75
36
37 (95%CI 0.62 to 0.91), strong evidence supporting treatment. The alternative is to calculate
38
39 the odds ratio¹¹, which is, exactly as the name suggests, the ratio of the *odds* of the outcome
40
41 in each study arm (rather than the probability). In this example, the odds of death with
42
43 placebo are 200:800 i.e. 1:4 or 0.25, versus odds of death with treatment of 150:850 i.e.
44
45 0.176. The odds ratio is therefore 0.176 / 0.25 = 0.71 (95%CI 0.56 to 0.89). Little has
46
47 changed to our conclusion – but the absolute value of the outcome metric and its confidence
48
49 intervals are different.
50
51
52
53
54
55

56 Relative risks are often preferred to odds ratios because they are simpler to understand¹²⁻¹⁴.
57
58 However, they have an unfortunate statistical property – their 95% confidence interval
59
60 depends greatly on how frequently the outcome occurs. The commoner the outcome (i.e. the
61
62
63
64
65

1 greater its prevalence), the narrower the 95% confidence interval. Consider a trial in which
2 the outcome (e.g. death) occurs in 10 of 100 untreated patients, and 9 of 100 treated
3 patients; the RR is 0.90 and the 95% confidence interval is very wide, at 0.38 to 2.12.
4
5 However, now consider a trial in which death is common, occurring in 50 of 100 untreated
6 patients and 45 of 100 treated patients; the RR remains 0.90, but the 95%CI is far narrower,
7
8 at 0.67 to 1.20. As the outcome becomes increasingly common, the 95%CI for the RR
9
10 narrows progressively; if 90 of 100 untreated patients were to die versus 81 of 100 treated
11 patients, the result is RR = 0.90, 95%CI 0.80 to 1.01. This explains the apparent discrepancy
12
13 in our example – the prevalence of the outcome (attendance at screening) was far higher for
14
15 You et al (80.3% for colonoscopy, 76.9% for CTC) than for all other studies (which ranged
16
17 from 14.8% to 33.6%). This likely happened because You et al pre-selected their
18
19 participants via an expression of interest in screening. The 95% confidence interval for the
20
21 relative risk is narrower than expected and the study is weighted accordingly, despite
22
23 randomising far fewer patients, fewer attendances at screening, and fewer non-attendances
24
25 than Scott et al (i.e. all event categories were smaller).
26
27
28
29
30
31
32
33
34
35

36 Moreover, this effect can cause bizarre results when it is unclear in which “direction” we
37
38 should define the outcome. For example, consider what happens if we choose to reverse the
39
40 outcome categories, so that *non-attendance* at screening becomes the “risk event” (rather
41
42 than attendance). This generates the forest plot shown in Figure 2.
43
44
45

46 Now we can see You et al has wide confidence intervals (0.61 to 2.26) and a low weighting
47
48 in the meta-analysis (1.05% for RR of missed appointments compared to 21.07% for RR of
49
50 attendance). Meta-analysis of RR for missed appointments is 0.92 with 95%CI 0.85 to 0.98,
51
52 $p = 0.01$, conventionally significant at the 5% level. Therefore, different conclusions can be
53
54 drawn from the same data, depending on whether attendance or non-attendance is defined
55
56 as the “risk event”. Both analyses have been performed correctly, but their results vary
57
58
59
60
61
62
63
64
65

1 because 95% CIs for relative risks are not symmetric with respect to what is defined as the
2 outcome⁴. This is clearly highly counterintuitive for clinical decision-making.
3
4
5

6 Although we are not aware of this phenomenon being described for previous radiological
7 meta-analyses, it has been reported in other scenarios. For example, describing a meta-
8 analysis comparing eradication of *Helicobacter pylori* for non-ulcer dyspepsia versus
9 placebo for trials conducted before 2000⁴, Deeks observed the same phenomenon. When
10 considering the outcome to be “ongoing dyspepsia”, the RR for treatment was 0.92 (95%CI
11 0.85 to 0.99), significant at the 5% level. However, if the outcome definition were to be
12 reversed (i.e. the outcome is “no dyspepsia”), the RR for treatment was 1.28 (95%CI 0.92 to
13 1.77), no longer significant at the 5% level (as more trials have been published, the benefit
14 of *H. pylori* eradication has become clear). It is clear that, for meta-analyses using relative
15 risks, there are in fact two possible summary risks – one for benefit (sometimes called RR_B),
16 and one for harm, RR_H ⁴; and there is no easily-predictable mathematical relationship
17 between the two.
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

35 So, how can the problem be resolved? Firstly, we can avoid over-reliance on the arbitrary
36 5% threshold to define significance, and inspect the data itself, which should prevent
37 spurious conclusions. A second option is to use odds ratios, the 95% confidence intervals for
38 which are unaffected by prevalence, or by which event category is defined as the outcome.
39 For example, Figure 3 shows the same data as in figures 1 and 2, but now presented as
40 odds ratios. Irrespective of whether we define attendance or non-attendance as the “risk
41 event”, study weights are assigned consistently and each summary estimate is simply the
42 reciprocal of the other, which seems intuitive if we have reversed the outcome.
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57

58 It is important to note that, under most circumstances, it is clear in which direction the “risk
59 event” should be specified (i.e. what constitutes the outcome of interest); and empirical data
60
61
62
63
64
65

shows that defining an adverse event as the outcome typically gives more consistent results and is preferred, certainly for preventative (rather than therapeutic) interventions⁴. Generally speaking, whichever is the less common state is usually preferable as being the “risk event”. Moreover, relative risks are much easier to interpret than odds ratios¹²⁻¹⁴, meaning they are often preferred for this reason alone. Fortunately, the situation that we have outlined above is rare, since it depends on both (a) a wide range of prevalences of the outcome of interest occurring in the component studies and (b) effect sizes varying between these different studies, which may not always be the case. Nonetheless, using relative risks may introduce a bias towards larger weights being assigned to smaller, early-phase RCTs (which typically are targeted to high prevalence scenarios in order to maximise event rates) than to larger, pragmatic RCTs (which aim to recruit more representative patient populations, sometimes including many without the outcome of interest). It seems fundamentally wrong that the outcome prevalence can influence weighting within meta-analysis, so that smaller studies, paradoxically recruiting higher-risk, less-representative patients, can outweigh larger studies. For example, imagine six studies comparing death rates in a series of placebo-controlled trials, all of which have the same average effect size when expressed as a relative risk (RR = 0.9), but recruiting different numbers of participants and with two differing death rates, 26% and 70%. Table 1 summarises these hypothetical data:

Table 1: Hypothetical data for six randomised trials of varying sizes and with varying event rates.

Study name	Treatment		Placebo		Total number of participants	Average event rate (death rate)
	Died	Survived	Died	Survived		
A	600	300	600	210	1710	70%
B	600	1800	600	1560	4560	26%
C	300	150	300	105	855	70%
D	300	900	300	780	2280	26%
E	120	60	120	42	342	70%

F	120	360	120	312	912	26%
---	-----	-----	-----	-----	-----	-----

1
2
3
4
5
6 Intuitively, we might expect studies B and D to contribute the greatest weight to a meta-
7 analysis, since they are the largest; however, meta-analysis using relative risks, whether
8 with fixed or random effects, shows this is not the case (figure 4).
9

10
11
12
13
14
15
16
17 Despite identical relative risks for all studies, the largest have wider 95% confidence
18 intervals and thus contribute smaller weights to meta-analysis, because the prevalence of
19 the outcome (death) is lower than in the smaller studies. The largest study (Study B) has
20 less weight within meta-analysis than the much smaller Study C, despite its raw data
21 contributing more patients in all categories (i.e. died or survived, for both treatment and
22 placebo). Therefore, as outlined by the Cochrane handbook¹⁵, it “may be wise to plan to
23 undertake a sensitivity analysis to investigate whether choice of summary statistic (and
24 selection of the event category) is critical to the conclusions of the meta-analysis” where
25 component study prevalence is variable; and consider using odds ratios rather than relative
26 risks for the primary analysis (bearing in mind the difficulties with their subsequent
27 interpretation). A further option, adopted by Zhu et al², is to report both relative risks i.e. for
28 both definitions of the outcome. This is probably only appropriate when it is arguable in
29 which direction the outcome should be specified (e.g. neither is clearly a negative or harmful
30 event, or much less common – in which case, the rarer and/or negative event should be
31 specified as the outcome).
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

51
52
53 In summary, we urge readers of meta-analyses, and researchers themselves, to consider
54 carefully their choice of summary statistical measure when analysing dichotomous
55 outcomes. Although relative risks are commonly chosen for simplicity, where outcome
56 prevalence varies greatly between component studies, and particularly where it is not clear
57
58
59
60
61
62
63
64
65

which category of outcome should be regarded as the event, researchers should exercise
caution and follow the Cochrane handbook guidance outlined above.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

BJR UNCORRECTED PROOFS

Figure Legends:

1
2
3
4 Figure 1: Forest plot similar to that generated by Zhu et al via random effects meta-analysis
5 using relative risks for attendance at screening; larger values imply greater attendance at
6
7
8
9 CTC when compared to colonoscopy. The summary estimate is 1.26 (95%CI 0.98 to 1.63), p
10 = 0.07, not statistically significant at a 5% level (despite the two largest trials finding a
11
12
13 significant result in favour of CTC).
14
15
16

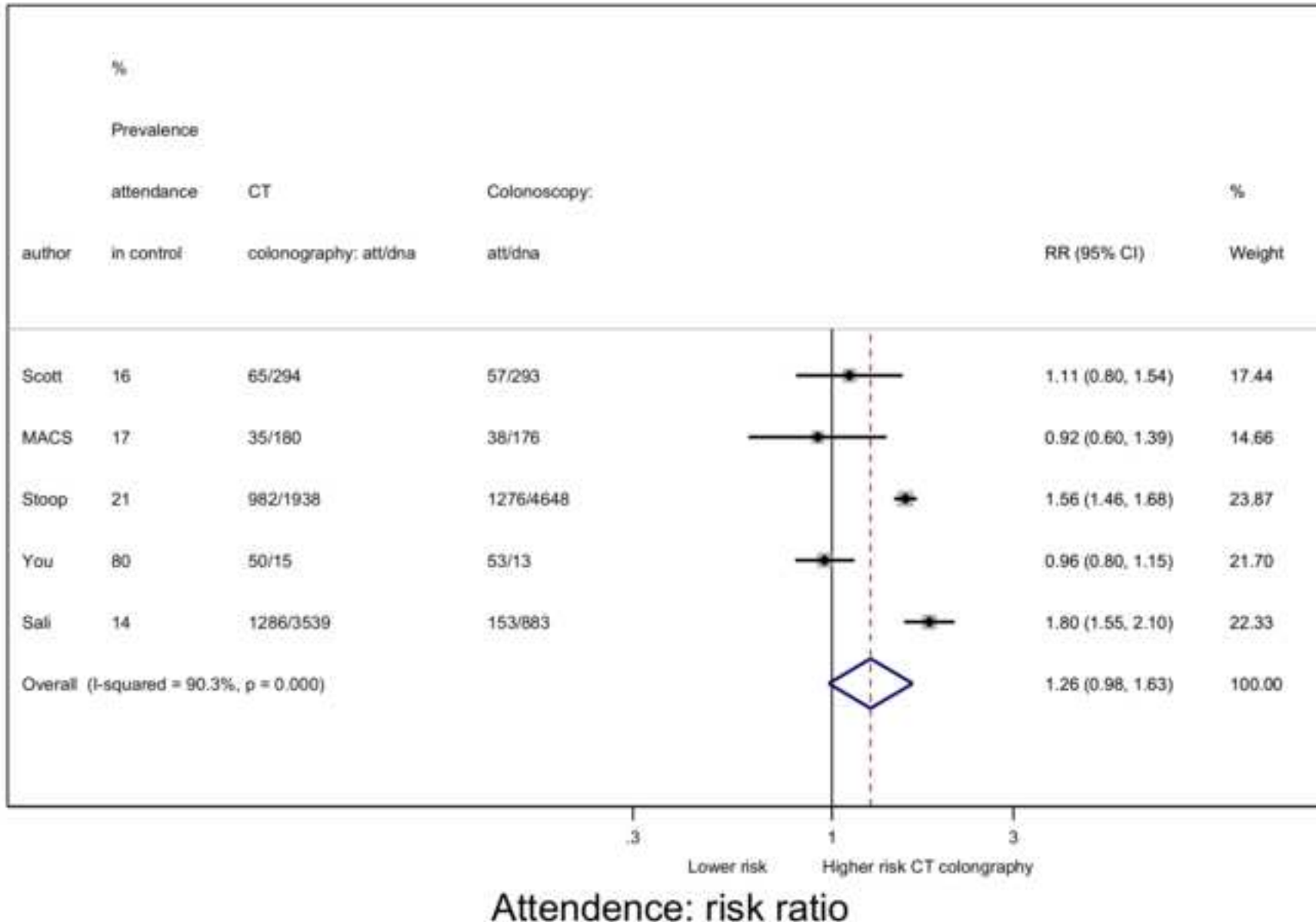
17
18 Figure 2: Forest plot generated using identical data for that in Figure 1, again via random
19 effects meta-analysis using relative risks, but reversing the category of the outcome, such
20 that relative risks indicate the risk of non-attendance at screening CTC vs colonoscopy. The
21
22
23
24 summary estimate is 0.92 (95%CI 0.85 to 0.98), $p = 0.01$.
25
26
27

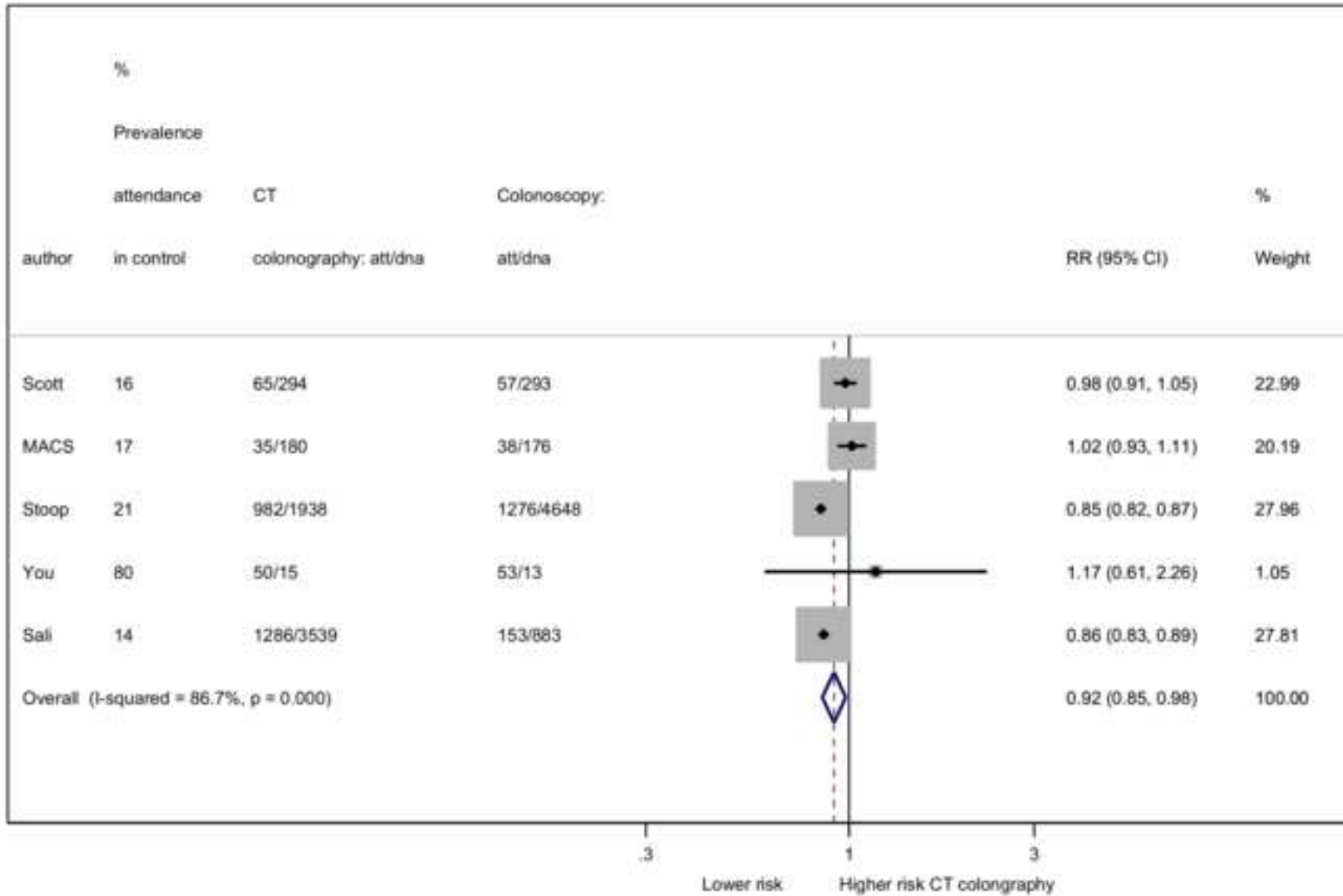
28
29 Figure 3: Forest plots generated using identical data for that in Figures 1 and 2, again via
30 random effects meta-analysis but now using odds ratios. Figure 3a shows attendance at
31
32
33 screening being the “risk event” (or outcome), whereas figure 3b uses non-attendance as
34
35
36 the outcome. The study weights are the same in each analysis, the forest plots are simple
37
38 “mirror images” and the summary odds ratios are reciprocals of each other.
39
40
41

42
43 Figure 4: Forest plot generated using the hypothetical data from Table 1 and for random
44 effects meta-analysis using relative risks. The same would be seen for a fixed effect meta-
45
46
47 analysis (since these hypothetical studies all have the same effect size, there is no between-
48
49
50 study variance and so the weights are the same for random and fixed effects meta-
51
52 analyses).
53
54
55
56
57
58
59
60
61
62
63
64
65

References

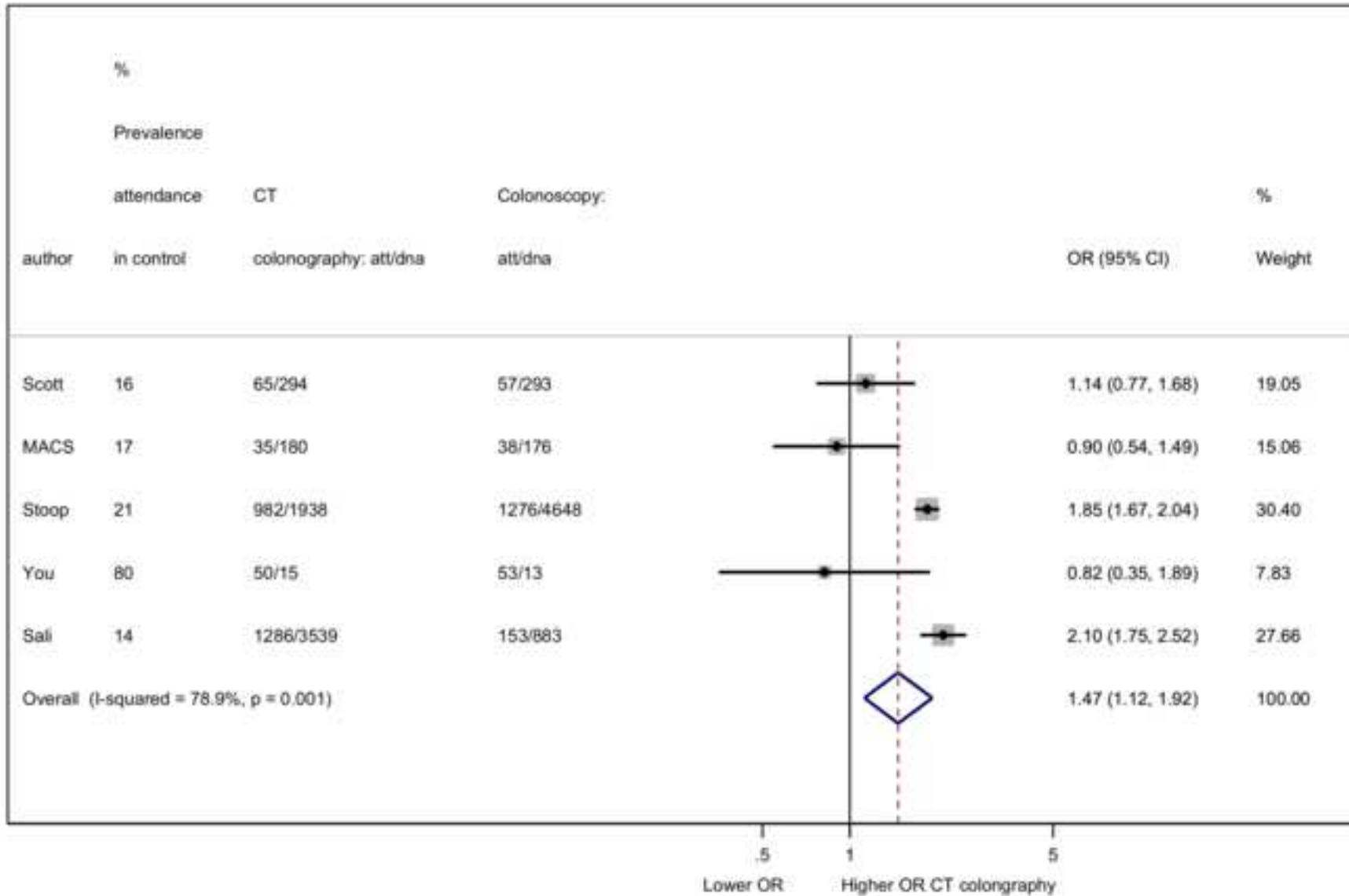
1. Cox DR, Snell EJ. Analysis of binary data. 2nd ed. ed: Chapman and Hall; 1989.
2. Zhu H, Li F, Tao K, et al. Comparison of the participation rate between CT colonography and colonoscopy in screening populations: A systematic review and meta-analysis of randomized controlled trials. *British Journal of Radiology* 2019; **92**: 20190240.
3. Engels EA, Schmid CH, Terrin N, Olkin I, Lau J. Heterogeneity and statistical significance in meta-analysis: an empirical study of 125 meta-analyses. *Stat Med* 2000; **19**(13): 1707-28.
4. Deeks JJ. Issues in the selection of a summary statistic for meta-analysis of clinical trials with binary outcomes. *Stat Med* 2002; **21**(11): 1575-600.
5. Borenstein M, Hedges LV, Higgins JP, Rothstein HR. A basic introduction to fixed-effect and random-effects models for meta-analysis. *Res Synth Methods* 2010; **1**(2): 97-111.
6. Stoop EM, de Haan MC, de Wijkerslooth TR, et al. Participation and yield of colonoscopy versus non-cathartic CT colonography in population-based screening for colorectal cancer: a randomised controlled trial. *Lancet Oncology* 2012; **13**(1): 55-64.
7. Sali L, Mascacchi M, Falchini M, et al. Reduced and Full-Preparation CT Colonography, Fecal Immunochemical Test, and Colonoscopy for Population Screening of Colorectal Cancer: A Randomized Trial. *JNCI J Natl Cancer Inst* 2016; **108** (2) (no pagination)(d3v319).
8. You JJ, Liu Y, Kirby J, Vora P, Moayyedi P. Virtual colonoscopy, optical colonoscopy, or fecal occult blood testing for colorectal cancer screening: results of a pilot randomized controlled trial. *Trials [Electronic Resource]* 2015; **16**: 296.
9. Scott RG, Edwards JT, Fritschi L, Foster NM, Mendelson RM, Forbes GM. Community-based screening by colonoscopy or computed tomographic colonography in asymptomatic average-risk subjects. *American Journal of Gastroenterology* 2004; **99**(6): 1145-51.
10. Multicentre Australian Colorectal-neoplasia Screening G. A comparison of colorectal neoplasia screening tests: a multicentre community-based study of the impact of consumer choice. *Med J Aust* 2006; **184**(11): 546-50.
11. Bland JM, Altman DG. Statistics notes. The odds ratio. *BMJ* 2000; **320**(7247): 1468.
12. Davies HT, Crombie IK, Tavakoli M. When can odds ratios mislead? *BMJ* 1998; **316**(7136): 989-91.
13. Cummings P. The relative merits of risk ratios and odds ratios. *Arch Pediatr Adolesc Med* 2009; **163**(5): 438-45.
14. Sackett DL, Deeks JJ, Altman DG. Down with odds ratios! *Evidence Based Medicine* 1996; **1**(6): 164-6.
15. Higgins JPT, Green S. Cochrane Handbook for Systematic Reviews of Interventions. 2011.





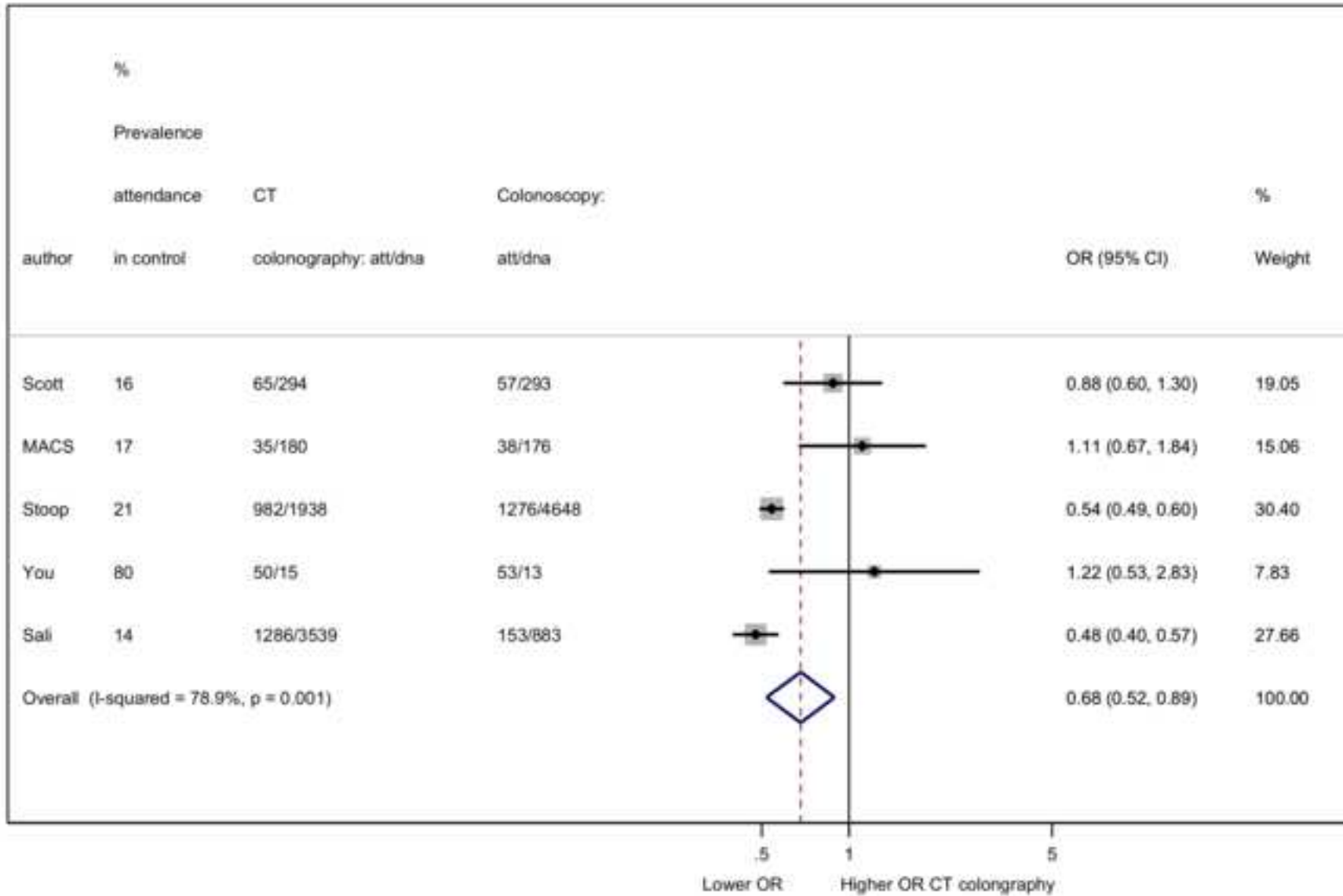
Missed appointments: risk ratio





Attendance: odds ratio





Missed appointments:odds ratio



