

# Statistical methods for meta-analyses including information from studies without any events—add nothing to nothing and succeed nevertheless

O. Kuss\*<sup>†</sup>

Meta-analyses with rare events, especially those that include studies with no event in one ('single-zero') or even both ('double-zero') treatment arms, are still a statistical challenge. In the case of double-zero studies, researchers in general delete these studies or use continuity corrections to avoid them. A number of arguments against both options has been given, and statistical methods that use the information from double-zero studies without using continuity corrections have been proposed. In this paper, we collect them and compare them by simulation. This simulation study tries to mirror real-life situations as completely as possible by deriving true underlying parameters from empirical data on actually performed meta-analyses. It is shown that for each of the commonly encountered effect estimators valid statistical methods are available that use the information from double-zero studies without using continuity corrections. Interestingly, all of them are truly random effects models, and so also the current standard method for very sparse data as recommended from the Cochrane collaboration, the Yusuf-Peto odds ratio, can be improved on. For actual analysis, we recommend to use beta-binomial regression methods to arrive at summary estimates for the odds ratio, the relative risk, or the risk difference. Methods that ignore information from double-zero studies or use continuity corrections should no longer be used. We illustrate the situation with an example where the original analysis ignores 35 double-zero studies, and a superior analysis discovers a clinically relevant advantage of off-pump surgery in coronary artery bypass grafting. Copyright © 2014 John Wiley & Sons, Ltd.

**Keywords:** meta-analysis; safety; sparse data; rare events; continuity correction

## 1. Introduction

Meta-analyses with rare events, for example, safety or adverse event outcomes, are methodically more challenging as compared with those with frequently observed efficacy outcomes [1, 2]. The main problem from a statistical viewpoint is the handling of studies that observe no event in one ('single-zero') or even both ('double-zero') treatment arms. Occurrence of meta-analyses with such extreme studies might seem to be a rare event itself, but empirical research suggests otherwise. For example, Vandermeer *et al.* [3] found that 30% of meta-analyses from a random sample of 500 Cochrane reviews contained at least one single-zero study. Using the same sample, we [4] found that in 34% of these reviews there was at least one meta-analysis that included a double-zero study. Moreover, we found a number of systematic reviews that have double-zero studies in their primary meta-analysis. An extreme example can be found in a systematic review that investigates the occurrence of fatal and nonfatal lactic acidosis with metformin use in type 2 diabetes mellitus [5]. There the primary analysis has 148 randomized trials (with more than 62,000 patients), which were all double-zero studies; that is, not a single event of lactic acidosis was observed in any of the 148 trials! Another systematic review reports on three meta-analyses, all without a

Institute for Biometry and Epidemiology, German Diabetes Center, Leibniz Institute for Diabetes Research at Heinrich Heine University Düsseldorf, Düsseldorf, Germany

\*Correspondence to: O. Kuss, Institut für Biometrie und Epidemiologie, Deutsches Diabetes-Zentrum, Leibniz-Zentrum für Diabetes-Forschung an der Heinrich-Heine-Universität Düsseldorf, Auf'm Hennekamp 65, 40225 Düsseldorf, Germany.

<sup>†</sup>E-mail: oliver.kuss@ddz.uni-duesseldorf.de

single event; that is, here the whole systematic review (and not only one meta-analysis within the review) constitutes double-zero studies [6].

Methods for meta-analyses with rare events are of considerable interest also in clinical research. For example, Nietert *et al.* [7] assessed 989 statistical journal articles from 45 statistical journals between 2000 and 2009 that were most frequently cited within general/internal medicine research articles. The article of Sweeting *et al.* [8] on meta-analytic methods for sparse data scored third in the cohort of 2000–2006 articles, and the article of Bradburn *et al.* [9] on the same topic scored second in the 2007–2009 cohort.

In the case of double-zero studies, most researchers pick one of three options: First, they explicitly delete these studies; second, they implicitly delete them by using statistical methods that do not use the information from double-zero studies; or third, they use artificial continuity corrections to avoid double-zero studies. For the following reasons, we and others consider all of these three options to be wrong or at least questionable. As double-zero studies point to no differences in treatment effects, at least in balanced trials, deleting them might bias the treatment effect away from the null [10]. Moreover, double-zero studies contain information through their mere sample size: a double-zero study with 2,000 treated patients provides stronger evidence for a null effect than a double-zero study with only 20 treated patients [11]. From an ethical viewpoint, patients who have been recruited in double-zero studies have a right to their data being also included in meta-analyses [12]. The use of continuity corrections, for example, adding 0.5 to each cell in the fourfold table of a double-zero study, adds an unpleasant element of arbitrariness into the analysis [13, 14], because results can depend on the respective value for correction, and a number of clinical examples has been given where different procedures to deal with double-zero studies resulted in different results [12, 15, 16]. With unbalanced data, even obvious paradoxes can be observed [17] if continuity corrections are used.

The problems of the current standard statistical methods for meta-analysis with rare events have been known for long. As early as in 1981, Breslow [18] reported on the failure of the still-standard-today inverse-variance fixed effects approach in the sparse data situation if the odds ratio is used to describe the treatment effect. In the sequel, a number of simulation studies [8, 9, 19–22] have investigated other procedures in the sparse data situation, and it is maybe agreed on today that the Mantel–Haenszel (MH) approach outperforms the standard inverse-variance methods, that is, the standard fixed and the DerSimonian–Laird random effects approach [23, Chapter 9.4.4.1]. In the situation of extreme sparseness (event rates below 1%), the Cochrane collaboration [23, Chapter 16.9.5], recommends to use the Yusuf–Peto odds ratio, relying on probably the most comprehensive simulation studies of Sweeting *et al.* [8] and Bradburn *et al.* [9] for rare events data. However, this recommendation ignores that the MaHa approach (at least for the odds ratio and the relative risk) and also the Yusuf–Peto odds ratio ignore studies with double-zero studies unless continuity corrections are applied.

In the following, we report on a simulation study that compares statistical methods for meta-analysis that explicitly include information from double-zero studies without using continuity corrections. In choosing the methods, we restricted on methods that can be computed within a manageable number of code lines and computing time within SAS (SAS Inc., Cary, NC, USA). This facilitates later usage of the code for other users and restricted computing time for our simulation study to reasonable time spans. To be concrete, we omitted especially the exact method of Tian *et al.* [24]. In terms of Bayesian analyses, Stijnen *et al.* [25] argued that these methods are not always well suited for meta-analysis with rare events. This is because with rare events there is only little information in the data, and this information is easily dominated by the information from the prior distribution, even if this is thought to be noninformative. However, we nevertheless included an analysis using Markov chain Monte Carlo (MCMC) sampling with noninformative priors for parameter estimation in generalized linear mixed models.

The paper is organized as follows. In Chapter 2 we give an example data set from a Cochrane review to motivate the study, Chapter 3 introduces the various methods, and in Chapter 4 we report on the simulation study. In Chapter 5 we come back to the example and compare the original results with those from a superior method that was identified in the simulation study. In Chapter 6 we conclude and give recommendations for practical work.

## 2. Example

We illustrate the methods under study with data from a Cochrane review on postoperative stroke occurrence when comparing off-pump and on-pump coronary artery bypass grafting for ischemic heart disease [26]. The full data set is given in Table I. This meta-analysis included 60 trials with 9,044 participants,

**Table I.** Example data set from a Cochrane review on postoperative stroke occurrence when comparing off-pump and on-pump coronary artery bypass grafting for ischemic heart disease [26].

Study	Off-pump		On-pump	
	Events	Total	Events	Total
OCTOPUS 2001	2	142	5	139
BHACAS I+II 2002	3	200	6	201
SMART 2003	2	100	2	100
Al-Ruzzeh 2006	2	84	1	84
DOORS 2009	10	450	18	450
MASS III 2009	3	156	5	155
ROOBY 2009	14	1104	8	1099
PROMISS 2010	0	73	0	74
BBS 2011	16	176	11	163
Matata 2000	0	10	0	10
Penttila 2001	0	11	0	11
Caputo 2002	0	20	0	20
Zamvar 2002	0	30	0	30
Carrier 2003	0	32	1	33
Raja 2003	3	150	4	150
Gerola 2004	0	80	0	80
PRAGUE-4 2004	0	208	2	192
Legare 2004	2	150	0	150
Lingaas 2004	0	60	2	60
Gasz 2005	0	10	0	20
JOCRI 2005	0	81	1	86
Ascione 2005	0	10	0	10
Niranjan 2006	1	40	1	40
Michaux 2006	0	25	0	25
Ascione 2006	0	20	0	20
Motallebzadeh 2006	1	108	5	104
Tatoulis 2006	0	50	0	50
Ozkara 2007	0	22	0	22
Hernandez 2007	0	102	3	102
Rasmussen 2007	0	18	0	17
Mandak 2008	0	20	0	20
Vural 1995	0	25	0	25
Gulielmos 1999	0	20	0	20
Czerny 2000	0	15	0	15
Diegeler 2000	0	20	1	20
Kochamba 2000	0	29	0	29
Wandschneider 2000	0	52	0	67
Czerny 2001	0	40	0	40
Sahlman 2003	1	24	1	26
Muneretto 2003	0	88	2	88
Lee 2003	0	30	1	30
Vedin 2003	0	33	0	37
Velissaris 2003	0	27	0	27
Parolari 2003	0	11	0	14
Motallebzadeh 2004	0	15	1	20
Selvanayagam 2004	0	30	1	30
Gasz 2004	0	10	0	10
Synnergren 2004	0	26	0	26
Blacher 2005	0	13	0	15
Rachwalik 2006	0	21	0	21
Malik 2006	0	25	0	25
Cavalca 2006	0	25	0	25
Gnenc 2006	0	30	0	12
Rainio 2007	0	10	0	10
Kunes 2007	0	17	0	17
Parolari 2007	0	14	0	15
Sajja 2007	0	60	1	60
Jares 2007	0	10	1	10
Formica 2009	0	30	0	30
Modine 2010	0	35	0	36

where a stroke was observed in 60 (from 4,527 participants) in the off-pump group and in 84 (from 4,517 participants) in the on-pump group. The authors used a standard inverse-variance random effects meta-analysis model for the relative risk and reported an effect of 0.76 with a 95% confidence interval (CI) of 0.54 to 1.06 and a  $p$ -value of  $p = 0.10$  in favor of the off-pump method. However, by using this method, the analysis misses the data of 35 double-zero studies (from 1,832 participants).

### 3. Statistical methods for meta-analysis including information from all studies without using continuity corrections

We defined a statistical method to include information from double-zero studies if the respective estimate or its standard error changed when double-zero studies were included in a meta-analysis without using any continuity correction. To this task, we used the toy example of six studies as given in Table II. This example included two double-zero studies (and also two single-zero studies), and statistical methods had to give different results (in terms of the effect estimate or its standard error) when the two double-zero studies were included or excluded.

The methods given in the following met the described test. Each subsection heading gives the name of the method and, in parentheses, its abbreviation in the results figures; the treatment effects (odds ratio [OR], relative risk [RR], risk difference [RD], or arcsine difference [AD]) that can be estimated by this method; and the SAS procedure that was used for parameter estimation.

#### 3.1. Collapsed table (COLL); OR, RR, RD; FREQ

These methods, sometimes also called ‘crude’ methods, ‘marginal’ methods, ‘pooled table’ [9], or ‘treat-as-one-trial’ [27], refer to aggregating all studies from the meta-analysis in a single fourfold table and computing effect estimates from this single table by standard methods. This explicitly ignores that data were collected from several studies, thus assuming that the underlying risk of an event is constant across trials. This results in a certain vulnerability of these methods to Simpson’s paradox [27,28].

#### 3.2. Generalized linear mixed model (GHQ, PQL, MCMC); OR, RR, RD; GLIMMIX, MCMC

As proposed by Platt *et al.* [29], a meta-analysis with binary outcomes can be perceived as a generalized linear mixed model (GLMM) with a binary outcome, the treatment as a fixed effect, and a normally distributed random intercept term for the study. By using a random study effect, we simultaneously allow that outcomes of patients within a single study might be correlated and that baseline event probabilities between studies might be different (but follow a normal distribution). As such, all the theory (see, e.g., [30]) and software for this mighty class of statistical models can be used for meta-analysis. Changing the link function from the canonical logit link (which gives an estimated log odds ratio for the treatment effect), we can also obtain estimates for the log relative risk (by using the log link), or the risk difference (by using the identity link). In terms of parameter estimation, approximate maximum likelihood methods (penalized quasi-likelihood, PQL), numerical integration (by Gaussian quadrature, GHQ), or stochastic integration (by MCMC methods, MCMC) can be used. Gao [31] found in a simulation study that the PQL approach is superior to GHQ in meta-analysis.

**Table II.** Example (‘toy’) meta-analysis that was used to assess if a statistical method includes information from double-zero-studies.

Study	Treatment		Control	
	Events	Total	Events	Total
1	0	100	0	100
2	0	100	0	100
3	0	100	2	100
4	0	100	2	100
5	2	100	4	100
6	2	100	4	100

### 3.3. Marginal Generalized linear model for correlated responses (GEE); OR, RR, RD; GENMOD

Similar to GLMM models (see 3.2), we can use all available methods for logistic regression models for correlated responses to estimate treatment effects in meta-analysis. Thus, also the generalized estimation equation (GEE) method can be applied with three different link functions to arrive at estimates for the odds ratio, the relative risk, and the risk difference.

### 3.4. Beta-binomial regression (BBIN); OR, RR, RD; NLMIXED

The beta-binomial model is another method for logistic regression with correlated responses and is comprehensively described in [32, 13.3] or [33, 13.4.2]. In general, assume we observe proportions  $p_i = y_i/n_i, i = 1, \dots, I$  from binomial distributions  $\text{bin}(n_i, \pi)$  where  $\pi$  has a beta distribution with parameters  $\alpha$  and  $\beta$ . The mean  $\mu$  of this beta distribution is  $E(\pi) = \mu = \alpha/(\alpha + \beta)$ , and the variance is  $\text{Var}(\pi) = \mu(1 - \mu)\theta/(1 + \theta)$  with  $\theta = 1/(\alpha + \beta)$ .

Marginally, that is when averaging with respect to the beta distribution for  $\pi$ , the  $y_i$  are beta-binomially distributed with mean  $E(y_i) = n_i\mu$ , variance  $\text{Var}(y_i) = n_i\mu(1 - \mu)[1 + (n_i - 1)\theta/(1 + \theta)]$ , and the correlation between two different individual observations  $y_{ij}$  and  $y_{ik}$  from  $y_i$  equal to  $\text{corr}(y_{ij}, y_{ik}) = \rho = 1/(\alpha + \beta + 1)$ .

The log-likelihood is defined conveniently in terms of  $y_i, n_i, \alpha$  and  $\beta$  by using  $\alpha = \mu(1 - \rho)/\rho$  and  $\beta = (1 - \mu)(1 - \rho)/\rho$  as follows:

$$\begin{aligned} \ell(\alpha, \beta) = \sum_{i=1}^I \ell_i(\alpha, \beta) = \sum_{i=1}^I & \lgamm(n_i + 1) + \lgamm(y_i + \alpha) + \lgamm(n_i - y_i + \beta) + \lgamm(\alpha + \beta) \\ & - \lgamm(y_i + 1) - \lgamm(n_i - y_i + 1) - \lgamm(n_i + \alpha + \beta) - \lgamm(\alpha) - \lgamm(\beta), \end{aligned}$$

where  $\lgamm$  denotes the natural logarithm of the gamma function.

In the meta-analytic context, each single study contributes two proportions, one from the control and one from the treatment group. To assess the treatment effect,  $\mu$  is modeled via

$$g(\mu) = b_0 + b_t x_t,$$

with  $g$  as one of the common link functions from the generalized linear model family,  $x_t = 0$  for the control,  $x_t = 1$  for the treatment group, and parameter  $b_t$ . Using the logit link for  $g$  gives a log odds ratio, and using the log link for  $g$  a log relative risk for the treatment effect. We avoid using the identity link and estimate the risk difference instead by using the estimated event probabilities  $\hat{p}_0 = g^{-1}(\hat{b}_0)$  and  $\hat{p}_1 = g^{-1}(\hat{b}_0 + \hat{b}_t)$  from the logit model for the control and treatment groups, respectively, and subtracting them. The respective standard error for the difference is then given by the delta method. Note that although the beta-binomial model is a true random effects model, it has a closed-form log-likelihood function, which facilitates parameter estimation. As such, virtually every tool that can maximize a self-written function with two parameters can be used. We use SAS PROC NLMIXED to this task although the 'MIXED' part of this procedure (that is, the RANDOM statement) that allows to include normally distributed random effects in the log-likelihood is not needed here.

### 3.5. Conditional logistic regression (CLRW); OR; PHREG

We saw in the previous paragraphs that all methods for logistic regression with correlated responses can be used for meta-analysis with a binary response. As such, we can also use the theory of conditional logistic regression. Here, the correlation of patients within studies is accounted for by conditioning on the study effect. This removes the study effect completely from the likelihood function and avoids assuming a distribution for this random effect as in GLMMs. Conveniently, the resulting conditional likelihood function is equivalent to the partial likelihood function from a Cox proportional hazard model, and any statistical software that allows fitting a stratified Cox model (with a flexible enough option for ties handling) can be used for meta-analysis [34, p.422]. It is important to note that in its original form the model does not use the information from double-null studies; however, by using the robust sandwich estimate of Lin and Wei [35] for the covariance matrix (issued via the COVSANDWICH statement in PROC PHREG), this information is included.

### 3.6. Logistic regression with fixed study effect and Firth likelihood (FIRT); OR; LOGISTIC

Firth [36] proposed a general bias correction term for maximum likelihood estimation and showed that this correction can be easily applied to the logistic regression model where it defines a penalized likelihood. Heinze [37] showed that Firth's method is especially useful in sparse data situation in the logistic model. As such, we use a logistic regression model with the Firth likelihood and the treatment effect and a fixed categorical study effect as covariates in our simulation.

### 3.7. Stijnen's bivariate Binomial-Normal model (StBN); OR; NLMIXED

Stijnen *et al.* [25, Chapter 3.5.1] borrowed from the current literature on meta-analysis of diagnostic studies and introduced a GLMM with a bivariate response for the meta-analysis of intervention studies. In this model, the log odds in the treatment and control groups are modeled separately, allowing different degrees of heterogeneity in both groups by including a random intercept term for the studies. The association between treatment and control groups within the same study is modeled by a covariance term in the covariance matrix of the random effects. Stijnen *et al.* point to the advantages of the bivariate approach, in that it allows different heterogeneity in both groups, allows missing arms in single studies, and can be used to assess the association between the treatment effect and the baseline risk.

### 3.8. Cai's Poisson-Gamma models (RBFR, RBRR); RR, NLMIXED

Using the same idea of conjugacy as in the beta-binomial model (3.4), Cai *et al.* [38] proposed a number of models that consider the response to be Poisson and the underlying event probabilities in the control group to follow a gamma distribution. Assuming the treatment to be fixed, the resulting Poisson-gamma likelihood [38, Chapter 2.1, (3), RBFR] has a closed form. Allowing the treatment effect to be random results in an extended model [38, Chapter 2.2, (4), RBRR]. It should be noted that this extension to a true GLMM does not assume the intercept but the treatment effect to be random. Cai *et al.* further argued to use a conditional likelihood, but their models do not use information from double-zero-studies ([38, Chapter 2.2.2, (6)]) or proved to be very unstable ([38, Chapter 2.2.2, (7)]).

### 3.9. Stijnen's bivariate Poisson-normal model (StPN); RR, NLMIXED

Similar to the model in Section 3.7, Stijnen *et al.* [25, Chapter 3.5.1] also gave a bivariate model for the log relative risk by using a Poisson assumption for the event proportions in the treatment and control groups.

### 3.10. Mantel-Haenszel method (MaHa); RD, DATA STEP

In contrast to the MaHa estimators for the odds ratio and the relative risk, the MaHa for the risk difference [39] explicitly uses information from double-zero studies. The common estimator with the variance estimate of Sato *et al.* [40] was improved on by Kuhnert and Böhning [41], but as their method includes a continuity correction for the weights (which would not meet our requirement of not using continuity corrections), we stay with the original method as given, for example, by Newman [42, Chapter 7.3].

### 3.11. Arcsine difference (CONS, STRD); AD, DATA STEP

Rücker *et al.* [43] expanded the collection of effect measures in meta-analysis for binary outcomes and introduced the arcsine difference. Writing  $p_0$  and  $p_1$  for the event probabilities in the control and treatment groups, respectively, the arcsine difference for a single study is defined via  $\arcsin \sqrt{p_1} - \arcsin \sqrt{p_0}$ , where for actual computation the observed proportions are plugged in. Study-specific estimates of the arcsine differences are then combined by the standard inverse-variance fixed effects approach to arrive at a meta-analytic estimate. Rücker *et al.* [43] gave two different estimators of the variance of a study-specific arcsine difference. The standard estimator (later abbreviated as STRD) is given by  $(0.25/n_0 + 0.25/n_1)$  with  $n_0, n_1$  denoting the overall sample sizes in the control and treatment groups, respectively. The conservative (later abbreviated as CONS) variance estimator will always give larger variances than the standard estimator and is expected to work better in the very sparse data case. It is motivated by a maximum variance argument and calculated by  $(0.42/n_0 + 0.42/n_1)$ . As the variance of the study-specific estimate does in both cases depend only on the number of total observations in treatment groups (and not on the number of events), the overall arcsine difference also includes information from double-zero studies. Because of this construction principle, the arcsine difference gives also sensible estimates and



CIs in the most extreme situation of meta-analyses that only include double-zero studies, an advantage that is not shared by any of the methods described here. Although Rücker *et al.* give a nice graph to compare the arcsine difference with the risk difference, the arcsine difference suffers from its limited interpretability and is very rarely used in practice.

### 3.12. Excluded methods

The following methods were excluded because they did not meet the described test of differing estimates or standard errors when including double-zero studies in the toy example from Table II or could not be evaluated because they did not converge with the toy example: The standard inverse-variance meta-analysis model, both fixed effect and random effect (OR, RR, RD; MIXED); the exact MaHa test (OR; LOGISTIC with EXACT-statement) as described in Kuss and Gromann [22]; a logistic regression model with a fixed, categorical study effect (OR; LOGISTIC); the hypergeometric-normal model from Stijnen *et al.* [25, Chapter 3.2.1] (OR; NLMIXED); the binomial-normal model from Stijnen *et al.* [25, Chapter 3.2., equation (11)] (OR; GLIMMIX); the binomial-normal model from Stijnen *et al.* [25, Chapter 3.4., equation (15)] (RR; NLMIXED); and, finally, bivariate models with beta-binomial margins and the association parameter modeled with a Gauss or a Plackett copula (OR, RR, RD; NLMIXED) as described in Kuss *et al.* [44].

## 4. Simulation

To compare the statistical properties of the different methods when estimating treatment effects, we conducted a simulation study. The simulation program was written in SAS and is available from the author on request, so all definitions and calculations might be checked in the original code. In an effort to mirror reality as complete as possible, true values for the design factors in the simulation study were gathered, where possible, from empirical data on performed meta-analyses. The most valuable source for this was the review of Turner *et al.* [45], which collected 14,886 meta-analyses with binary outcomes (required to include at least two studies) from 1,991 Cochrane Reviews. Rebecca Turner was kind enough to deliver some additional information on request.

### 4.1. Design

The following five design factors were varied in the simulation study.

- **Proportion of double-zero studies** (Reference, 25%, 50%, and 75%)

In a recent paper in this journal, Stijnen *et al.* [25] introduced a new method for meta-analysis of sparse data and at the same time stated that ‘future research is needed to study what amount of sparseness is allowed for the reliable application of these methods.’ Following this call, our main motivation was to check the robustness of the different methods against increasing proportions of double-zero studies. To achieve this, we defined four scenarios with increasing proportions. The reference scenario was intended to represent a typical not-too-sparse situation; as such the event probability was set to 10% in the control group, resulting in an overall proportion of double-zero studies < 5% throughout. In the remaining scenarios, the event probabilities were chosen to result in proportions of double-zero studies of 25%, 50%, and 75%. It should be noted that by generating double-zero studies, also a number of single-zero studies with zero events in the control or in the treatment arms are generated. In Table III, we report on this and some other characteristics (percentages of overall events or percentages of meta-analyses without any events) of the simulation scenarios.

- **Size of treatment effect** (no effect, medium effect)

In all simulation settings, we imagined the situation of a meta-analysis with the outcome being a rare adverse (or safety) event where the treatment is aiming for a further lowering of events as compared with the control. As such, we consider an odds ratio of 1 as the true treatment effect in the ‘no effect’ situation. In the medium effect situation, we use an odds ratio of 0.684, which corresponds to the median odds ratio from 7,887 fixed effects meta-analyses with an odds ratio below 1 from the Turner *et al.* [45] study (R. Turner, personal communication).

- **Number of studies per meta-analysis** (Cochrane, non-Cochrane)

The number of studies in the ‘Cochrane’ scenarios was generated from a log-normal distribution with mean 0.65 and standard deviation 1.2 with a subsequent ceiling step (a function returning the smallest integer that is greater than or equal to the generated number) and finally adding 1. This

**Table III.** Description of the simulation scenarios in terms of percentages of overall events, double-zero studies, single-zero studies with zeros in the control or the treatment arm, and meta-analyses without any events.

Simulation scenario				Results				
Treatment effect	Random effects variance	Number of studies	Proportion of double-zero studies	% Events overall	% Double-zero studies	% Single-zero studies control	% Single-zero studies treatment	% Meta-analyses without any event
No	FEM	NC	R	10.0	2.3	5.0	5.1	0.0
			25	1.7	25.0	16.7	16.7	0.0
			50	0.6	50.0	16.5	16.3	0.3
			75	0.2	75.1	10.3	10.3	3.3
		CO	R	10.0	2.4	5.0	5.0	0.0
			25	1.7	25.0	16.4	16.4	2.2
			50	0.6	50.0	16.3	16.6	11.0
			75	0.2	75.0	10.4	10.4	35.0
	REM	NC	R	11.2	2.4	4.8	7.7	0.0
			25	2.3	25.0	17.4	17.7	0.1
			50	0.9	50.1	17.8	15.9	0.4
			75	0.3	75.2	11.6	9.3	3.7
		CO	R	11.1	2.4	4.9	7.5	0.0
			25	2.3	25.0	17.3	17.8	2.0
			50	0.9	50.0	17.8	16.0	11.2
			75	0.3	75.0	11.7	9.3	35.2
Medium	FEM	NC	R	8.5	3.1	4.3	8.6	0.0
			25	1.7	25.0	12.4	21.9	0.0
			50	0.6	50.0	12.6	20.8	0.3
			75	0.2	75.0	8.2	12.6	3.1
		CO	R	8.5	3.1	4.3	8.6	0.0
			25	1.7	25.0	12.2	21.6	2.3
			50	0.6	50.0	12.8	20.4	11.0
			75	0.2	75.0	8.3	12.5	34.0
	REM	NC	R	9.7	3.1	4.1	11.3	0.0
			25	2.2	25.0	13.2	22.6	0.0
			50	0.9	50.0	14.2	20.0	0.4
			75	0.3	75.0	9.6	11.6	3.9
		CO	R	9.6	3.1	4.3	11.2	0.0
			25	2.2	25.0	13.3	22.3	2.0
			50	0.9	50.0	14.0	19.9	11.2
			75	0.3	75.0	9.8	11.2	34.6

The abbreviations 'NC' and 'CO' refer to the non-Cochrane and the Cochrane scenarios, respectively.

procedure reproduces minimum, 25% percentile, median, and 75% percentile of the distribution of the number of studies in Turner *et al.* [45, Table 2]. As this procedure results in rather small meta-analyses, for example, only 5% of all meta-analysis would include 15 or more studies, we also defined a 'non-Cochrane' scenario that relies on a systematic review of Moher *et al.* [46] reporting on 88 systematic reviews evaluating therapeutic interventions that were published outside the Cochrane library. These reviews had a median number of 23 studies, the 25%/75% percentiles being 12/38. Using the procedure described previously with a log-normal mean of 3.05 and a log-normal standard deviation of 0.97 reproduces the median and the 25% percentile and gives a 75% percentile of 42. To avoid excessive simulation time, the maximum number of studies for one meta-analysis was set to 100 in all scenarios.

#### • Sample size of single study

Turner *et al.* [45, Table 2] report on the distribution of overall sample sizes from 77,237 single studies. Using the same procedure as described previously with a log-normal mean of 4.615 and a log-normal standard deviation of 1.1 reproduced minimum (two patients), the 25% percentile (50 patients), and the median (102 patients) of their distribution.

#### • Random effects variance (FEM, REM)

In the fixed effects (FEM) situation, events numbers were generated from the underlying true event probabilities in both groups. In the random effects (REM) situation, a random  $\tau^2$  was sampled



from the distribution in the Turner *et al.* [45] review, which is log-normal with mean  $-1.47$ , standard deviation  $1.65$ , and skewness  $-0.55$  (R. Turner, personal communication), equivalent to a distribution with 25% percentile/median/75% percentile of  $\tau^2$  of  $0.079/0.274/0.806$ . Fleishman's [47] power transformation method as given by Fan *et al.* [48, Chapter 4.2.2.2] was used to generate this skewed distribution.

Combination of the five design factors (proportion of double-zero studies, size of treatment effect, number of studies per meta-analysis, sample size of single study, and random effects variance) resulted in a total of 32 simulation scenarios. For each fixed scenario, 10,000 meta-analyses were generated. This number was determined mainly to keep computing time within a reasonable limit. A reviewer pointed out that, as an additional justification, by using 10,000 meta-analyses the standard error of an estimated percentage (e.g., for the empirical coverage) is guaranteed to be smaller than 0.5. To obey to the desired proportion of double-zero studies across all studies in these meta-analyses, true event probabilities for both treatment groups were determined, depending on the real treatment effect. Then, for each single meta-analysis the number of studies was generated, and for each single study the number of patients. According to a simple binomial experiment with probability 0.5 (thus mimicking randomization), patients were allocated to the treatment groups. In a FEM scenario, the true event probabilities in both groups were held constant for all 10,000 meta-analyses. In a REM scenario, a random  $\tau^2$  was simulated for each single meta-analysis. For each single study, the true event probability in the control group then was again held constant, while the true event probability in the treatment group was generated from a standard inverse-variance random effects model for the log odds ratio. Finally and for each single study, the number of events in both groups were generated by two binomial draws with respective event probabilities and sample sizes.

#### 4.2. Estimation

From each of the generated meta-analyses, we estimated the treatment effect, its standard error, and a 95% Wald CI. The treatment effect was measured on the log scale for the odds ratio and the relative risk and on the original scale for the risk and the arcsine difference. Standard SAS procedures and data step calculations (as given in the subsection headings in the previous chapter) were used for the different methods. To enable a fair comparison between models, all SAS procedures were run with the default options (with the exception of the methods of Stijnen *et al.* where adding the NOAD option in PROC NLMIXED improved convergence). In all procedures using PROC NLMIXED, the starting values for each single meta-analysis were computed from raw proportions, their variances, and correlations as requested (and possibly after transformations) for the respective procedure. For the MCMC methods and for each single meta-analysis, we used a single chain of length 10,000 with an additional burn-in of 1,000 draws. Prior distributions for the models parameters were normal distributions with zero mean and variance 10,000 for the fixed effects parameters and the inverse gamma distribution with scale and shape parameter 0.01 for the random effects variance.

We used bias and empirical coverage (to the 95% level) as outcomes to compare the different estimation methods under the null hypothesis of no treatment effect. In the medium treatment effect scenarios, we report empirical power and empirical coverage (to the 95% level). To assess the numerical robustness of the different methods, we also report the number of non-missing values for the estimated treatment effect. As convergence diagnosis is somewhat arbitrary and sometimes exceedingly large or small values were given as properly converged estimates from the different procedures, we always report median values for bias.

In the simulation study, we also included three reference methods, although all of them either ignore information from double-zero studies or use continuity corrections to arrive at estimates. First, the Yusuf-Peto method (abbreviated as YPET) for the odds ratio was included, because this method is recommended by the Cochrane collaboration in the case of extreme sparseness (below 1%). We also included the corrected Yusuf-Peto method of Sato [49] (abbreviated as YPSA), which uses an improved variance estimate. Second, we used the standard inverse-variance fixed effects methods for meta-analysis for odds ratios and relative risks within SAS PROC FREQ. These are abbreviated as 'FREQ' in the results, and 0.5 was added to each cell in a study table if no event was observed in at least one of the treatment groups. Third, we also included the results from the standard 'Mantel-Haenszel' method from PROC FREQ for odds ratios and relative risks. These Mantel-Haenszel estimates explicitly ignore information from double-zero studies.

## 4.3. Results

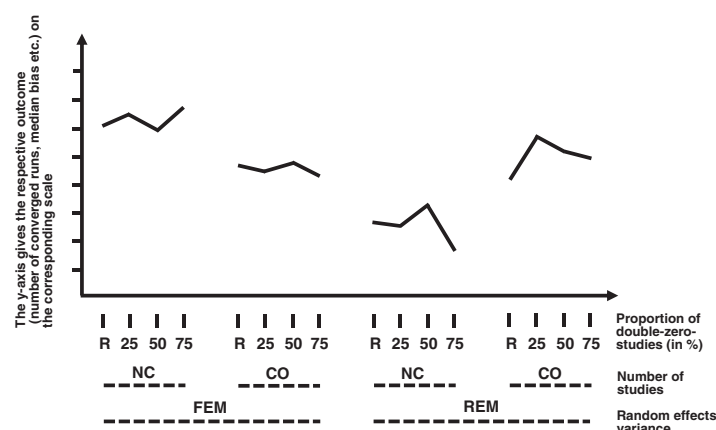
We report simulation results separately for each effect measure and for the two different treatment effects. To avoid lengthy tables, the results for the proportion of double-zero studies, the number of studies, and the random effects variance are given, separated by methods, as four lines in a graph. The legend in Figure 1 explains how the different factors are ordered in each single graph: The first two lines report on the FEM model and the second two lines on the REM model. The first and the third lines give the results for the ‘non-Cochrane’ (NC) setting with a larger number of studies and the second and fourth those for ‘Cochrane’ (CO) setting with a smaller number of studies. Finally, the four values within a single line are ordered by increasing proportion of double-zero studies: The first value reports on the reference scenario, the last on the 75% scenario. Roughly spoken, the more on the right side, the more challenging is the reported scenario under study in terms of proportion of double-zero studies, number of studies, or random effects variance. To refer to the exact numbers that underly the graphs, we give all the results in the same ordering of effect measure, response, and method in the Supporting information.

We first report the results for the zero treatment effect. As will be seen from these results, many of the methods do not keep the then requested empirical level of 95%. Comparing these anti-conservative methods under the alternative of an existing treatment effect to the methods that kept the empirical level would be unfair. As such, under the scenarios of the medium treatment effect, we only report results from the methods that kept at least roughly the empirical level under the null hypothesis of no treatment.

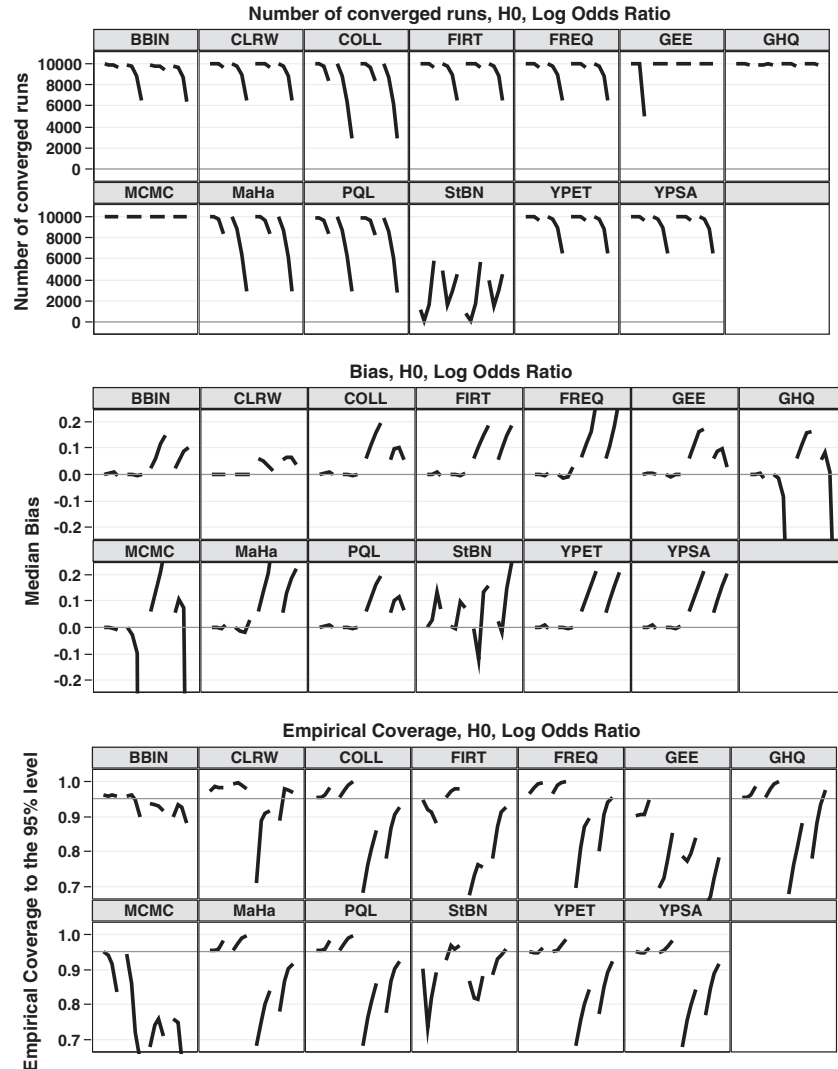
### 4.3.1. Null hypothesis of no treatment effect.

**General observations.** We start with some general observations across all scenarios. The more challenging the situation in terms of proportion of double-zero studies or random effects variance, the worse the convergence. As expected, methods that are only defined for a fixed treatment effect (COLL and MaHa for all effect measures, FIRT, FREQ, YPET, and YPSA for the log odds ratio, and both AD methods) perform well in the FEM situation, but have compromises in the REM situation. The reference methods that use continuity corrections (FREQ for the log odds ratio and the log relative risk) or ignore information from double-zero studies (FREQ, YPET, and YPSA for the log odds ratio and FREQ for the log relative risk) overall perform worse than the other methods (Figures 2–5).

In terms of the number of studies, results from ‘non-Cochrane’ and ‘Cochrane’ scenarios are rather similar, at least with respect to bias and coverage, indicating a small influence of the number of studies for these outcomes. With regard to convergence, ‘non-Cochrane’ scenarios show better convergence, where this holds independently from heterogeneity, that is, in FEM as well as in REM scenarios.

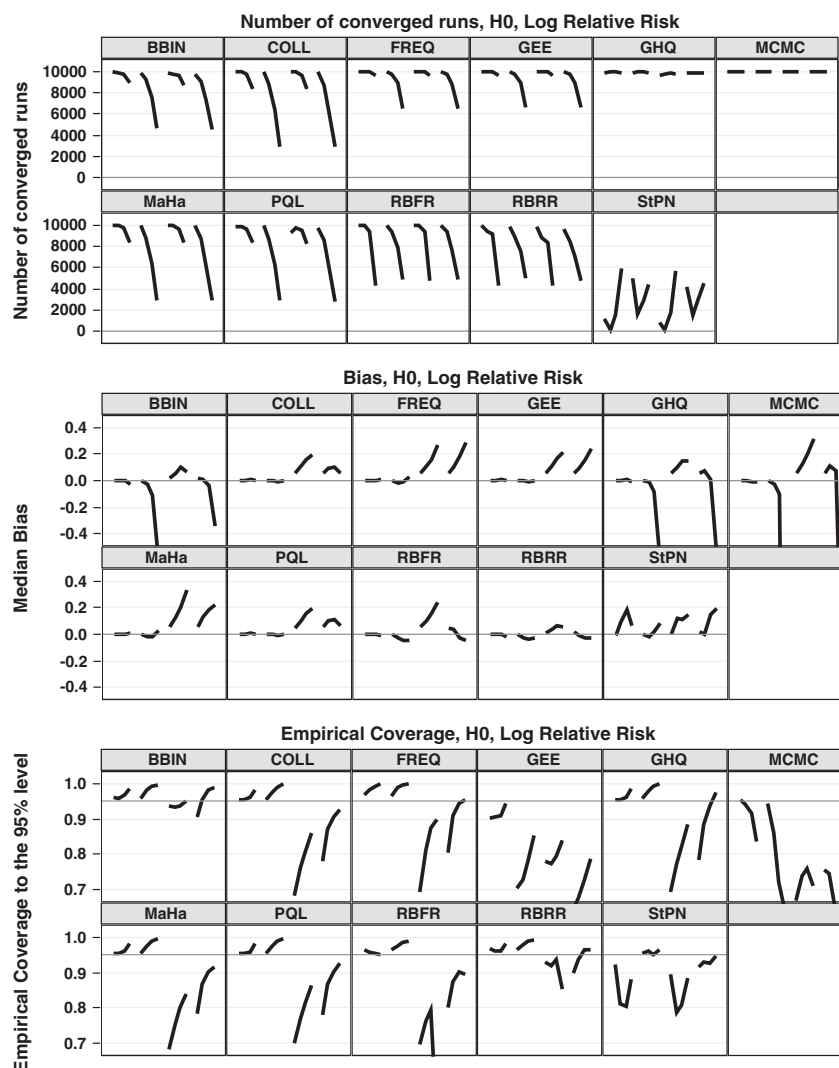


**Figure 1.** Legend for assessing simulation results in Figures 2 to 6: The first two lines in each single graph report on the fixed effects model (FEM) and the second two lines on the random effects model (REM). The first and the third lines give the results for the ‘Non-Cochrane’ (NC) setting with a larger number of studies, and the second and fourth those for ‘Cochrane’ (CO) setting with a smaller number of studies. Finally, the four values within a single line are ordered by increasing proportion of double-zero studies: The first value reports on the reference (R) scenario, the last on the 75% scenario. On the y-axis, the values for the respective outcome (number of converged runs, median bias, empirical coverage to the 95% level, or empirical power) on the respective scale are given. Please note that the black lines within the graph (actually the data of interest) are here just given for illustration.



**Figure 2.** Simulation results (number of converged runs, median bias, and empirical coverage to the 95% level) for the log odds ratio under the null hypothesis of no treatment effect. Lines within graphs are ordered as described in Figure 1. Methods are abbreviated as follows: BBIN, beta-binomial regression (3.4); CLRW, conditional logistic regression (3.5); COLL, collapsed table (3.1); FIRT, logistic regression with fixed study effect and Firth likelihood (3.6); FREQ, standard fixed effects meta-analysis with a 0.5 correction for empty cells; GEE, marginal generalized linear model for correlated responses (3.3); GHQ, generalized linear mixed model (estimation by Gauss-Hermite quadrature) (3.2); MCMC, generalized linear mixed model (estimation by Markov chain Monte Carlo) (3.2); MaHa, Mantel-Haenszel method; PQL, generalized linear mixed model (estimation by penalized quasi-likelihood) (3.2); StBN, Stijnen's bivariate binomial-normal model (3.7); YPET, Yusuf-Peto method; YPSA, Yusuf-Peto method with Sato's variance estimate.

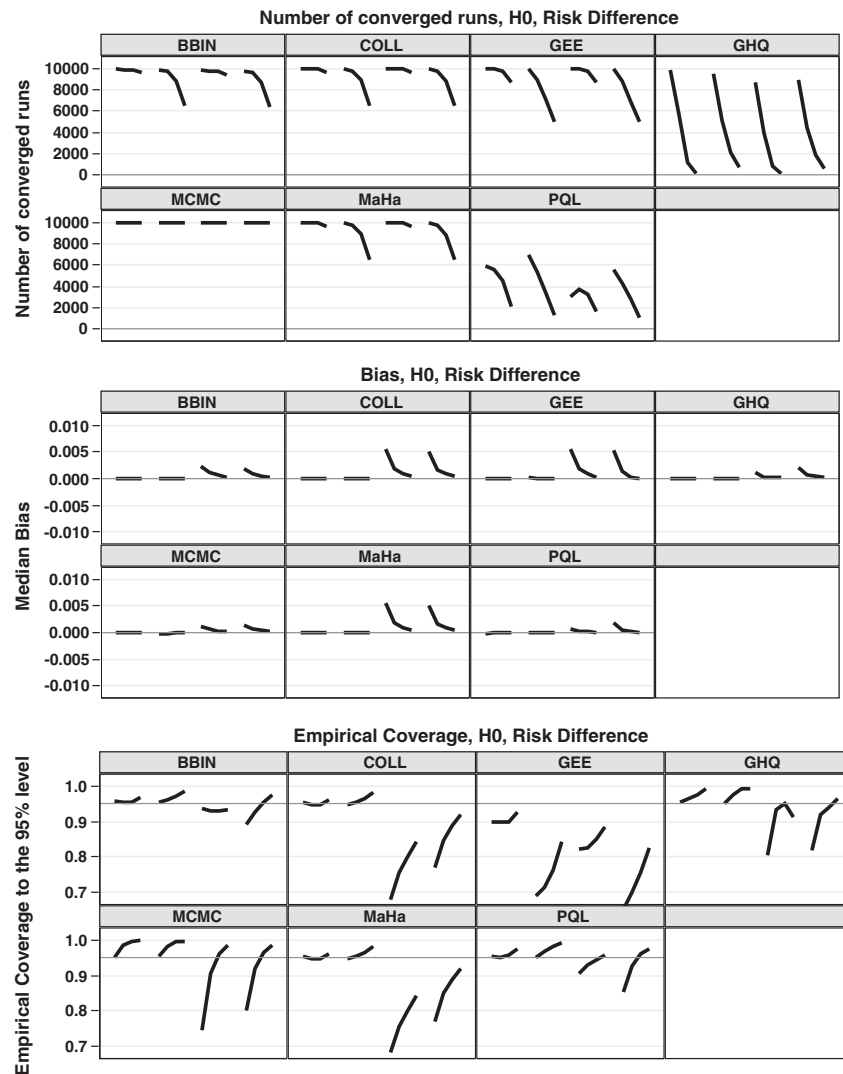
**Log odds ratio (Figure 2).** With respect to robustness, GEE, GHQ, and the MCMC method perform best, where the good performance of the GHQ can maybe explained by weaknesses in diagnosing convergence. As can be seen from the bias of the GHQ method, where many negative bias estimates point to many very negative parameter estimates, those parameter estimates should be rather diagnosed as non-converged. The perfect convergence of the MCMC method is expected, because here convergence means the error-free performance of the 10,000 posterior simulation runs. The low convergence rates of Stijnen's model (StBN) seem unacceptable but may not be that surprising, because with this method two random effects have to be estimated. All other methods perform similarly in terms of convergence. In terms of bias, the conditional logistic regression performs best; biases for GHQ, MCMC, and StBN method are unacceptable. All other methods perform similarly in terms of bias with satisfying values in the FEM and upward bias in the REM situation. Considering empirical coverage, the beta-binomial



**Figure 3.** Simulation results (number of converged runs, median bias, and empirical coverage to the 95% level) for the log relative risk under the null hypothesis of no treatment effect. Lines within graphs are ordered as described in Figure 1. Methods are abbreviated as follows: BBIN, beta-binomial regression (3.4); COLL, collapsed table (3.1); FREQ, standard fixed effects meta-analysis with a 0.5 correction for empty cells; GEE, marginal generalized linear model for correlated responses (3.3); GHQ, generalized linear mixed model (estimation by Gauss-Hermite quadrature) (3.2); MCMC, generalized linear mixed model (estimation by Markov chain Monte Carlo) (3.2); MaHa, Mantel-Haenszel method; PQL, generalized linear mixed model (estimation by penalized quasi-likelihood) (3.2); RBFR, Cai's Poisson-gamma model, fixed treatment effect (3.8); RBRR, Cai's Poisson-gamma model, random treatment effect (3.8); StPN, Stijnen's bivariate Poisson-normal model (3.9).

model (BBIN) clearly outperforms the other methods with coverages within the range of 96% to 88% in all scenarios. All other methods perform worse, with anti-conservativeness going down to below 70%. Sato's improved variance estimator (YPSA) for the Yusuf-Peto odds ratio does not lead to an improved behavior of the Yusuf-Peto method; results for YPSA and the standard method (YPET) are practically identical.

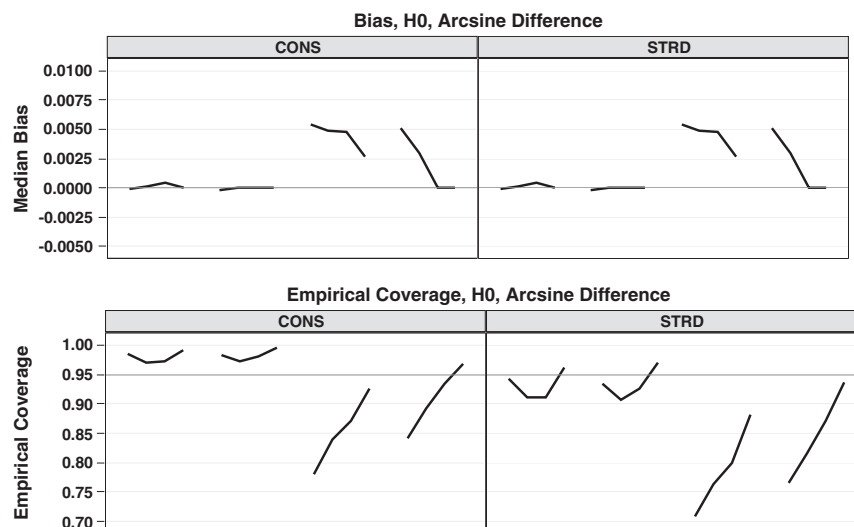
**Log relative risk (Figure 3.)** In terms of convergence, BBIN, COLL, GEE, and PQL perform rather similarly with satisfying convergence with a larger number of studies and a worse convergence with a smaller number of studies. In contrast, Cai's methods (RBFR, RBRR) also have converge problems in the less challenging situations with a large number of studies. The GHQ and the MCMC method show the same problems (very good convergence with unacceptable low bias in the more challenging situations) as seen for the log odds ratios. Also similar to the log odds ratio situation, Stijnen's method (StPN) has large convergence problems. Considering bias, Cai's method with random effects (RBRR) outperforms the other methods with all biases being smaller than 0.07 in absolute values. Unacceptable



**Figure 4.** Simulation results (number of converged runs, median bias, and empirical coverage to the 95% level) for the risk difference under the null hypothesis of no treatment effect. Lines within graphs are ordered as described in Figure 1. Methods are abbreviated as follows: BBIN, beta-binomial regression (3.4); COLL, collapsed table (3.1); GEE, marginal generalized linear model for correlated responses (3.3); GHQ, generalized linear mixed model (estimation by Gauss-Hermite quadrature) (3.2); MCMC, generalized linear mixed model (estimation by Markov chain Monte Carlo) (3.2); MaHa, Mantel-Haenszel method (3.10); PQL, generalized linear mixed model (estimation by penalized quasi-likelihood) (3.2).

biases are found for BBIN, GHQ, and MCMC; all other methods perform rather similarly. Two methods perform well in terms of empirical coverage, that is, the beta-binomial model (BBIN) with all observed coverages above 90% and Cai's method with random effects (RBRR) with only one observed coverage clearly below 90%.

**Risk difference (Figure 4).** With regard to convergence, MCMC again performs perfectly. BBIN, COLL, GEE, and MaHa perform similarly with good convergence for the larger number of studies and compromises for the smaller number of studies. The other methods for the GLMM (GHQ and PQL) have inferior convergence. In terms of bias, all methods work well in the fixed effect situation (first two lines in each graph); BBIN, GHQ, MCMC, and PQL (which are in fact random effects estimation methods) also perform well in the random effect situation, where COLL, GEE, and MaHa surprisingly show some upward bias in the situations with low proportions of zero-event trials. Considering empirical coverage, the BBIN is superior with only one value below 90%; the PQL also works well. The fixed effects methods (COLL and MaHa) again work well in the FEM situation but worse in the REM situation.



**Figure 5.** Simulation results (median bias and empirical coverage to the 95% level) for the arcsine difference under the null hypothesis of no treatment effect. The number of converged runs is omitted here because this was 10,000 in all of the scenarios. Lines within graphs are ordered as described in Figure 1. Methods are abbreviated as follows: CONS, arcsine difference with conservative variance (3.11); STRD, arcsine difference with standard variance (3.11).

**Arcsine difference (Figure 5).** The AD methods, as compared with the methods for the other effect estimates, perform perfectly in terms of robustness. Because of their construction principle, they give proper estimates and CIs in all of the 160,000 meta-analyses. Being fixed effect methods by definition, both AD methods work well in the FEM situation but show an upward bias and an anti-conservative coverage in the REM situation. Comparing the two different AD methods, they give, by definition, identical estimates for bias and the number of converged runs. In terms of empirical coverage, the conservative (CONS) method, as expected, gives more conservative CIs than the standard (STRD) method.

#### 4.3.2. Alternative hypothesis of a medium treatment effect (Figure 6).

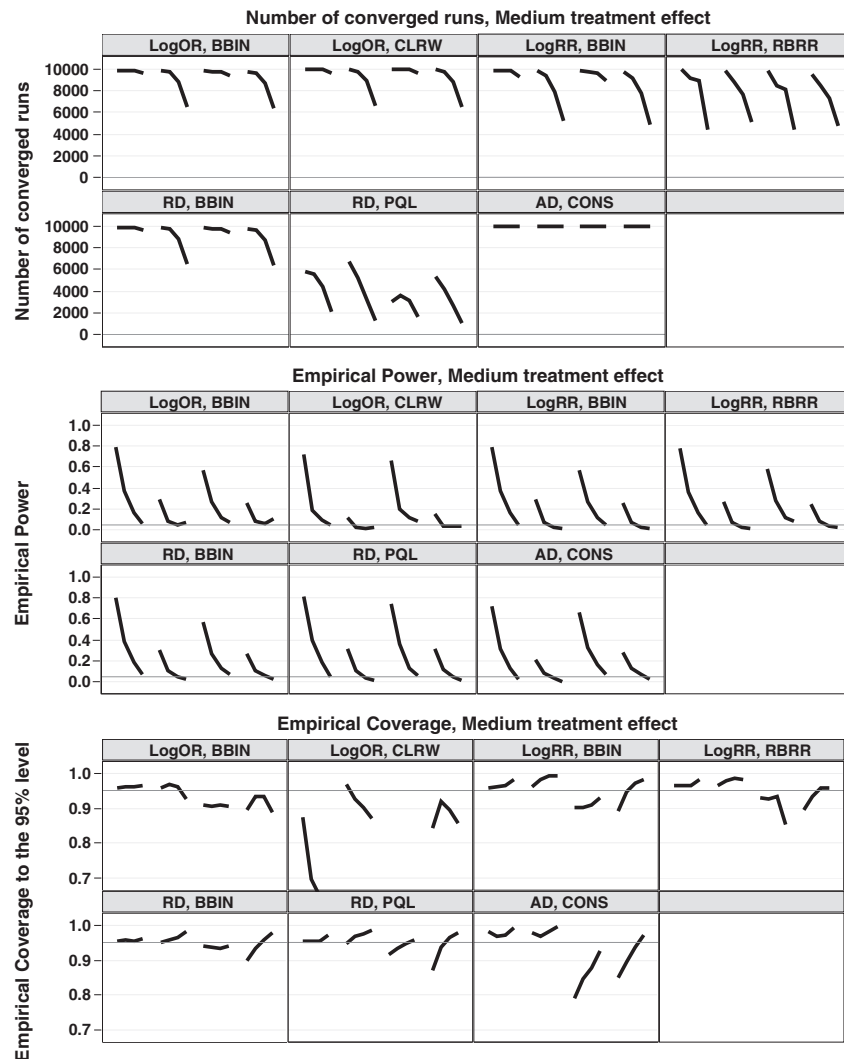
**General observations.** Four general points can be made. First, the convergence for the different methods is very similar to the respective convergence results from the no treatment effect scenarios. Second, the power depends heavily on the number of studies. Estimated powers are always smaller in the ‘Cochrane’ scenario which has smaller numbers of studies as compared with the ‘non-Cochrane’ scenario. Third, the empirical power decreases rapidly with the proportion of double-zero studies increasing. In the sparse scenario with 50 and 75% of double-zero studies, the power is sometimes even smaller than 5%, which is the expected power under the null hypothesis of no treatment. Fourth, in terms of empirical coverage, results are very similar to the results from the no treatment effect scenarios.

**Results for the different effect measures.** Comparing the methods separately for the different effect measures, results in terms of convergence and empirical power are similar for the two methods (BBIN and CLRW) for the log odds ratio. The CLRW is slightly superior with respect to convergence but has an inferior power (the median power across all scenarios is 11.1% for BBIN and 9.1% for CLRW). In terms of empirical coverage, the behavior of the CLRW method is unacceptable. For example, all coverages in the ‘non-Cochrane’ scenario with a random effects variance are below 70%, with results not even appearing in the graph. For the log relative risk, empirical power and coverage are similar for BBIN and RBRR; however, the convergence for BBIN is clearly better with the RBRR method having problems with increasing sparseness, independently of the random effects variance and the number of studies. For the risk difference again, the two methods under study (BBIN and PQL) are similar with respect to empirical power and coverage, whereas the convergence for the PQL method is rather unacceptable.

## 5. Updated analysis for the example data set

Using the beta-binomial model for computing the relative risk for our example data set, we find a value of 0.51 with a 95% CI of 0.28 to 0.92 ( $p = 0.026$ , see also Figure 7). That is, by using the standard

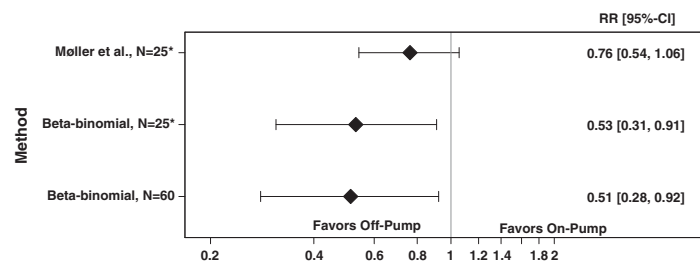




**Figure 6.** Simulation results (number of converged runs, empirical power, and empirical coverage to the 95% level) under the alternative hypothesis of a medium treatment effect for the seven methods that roughly kept the empirical level under the null hypothesis of a null treatment. Lines within graphs are ordered as described in Figure 1. Methods are abbreviated as follows: LogOR, BBIN: log odds ratio, beta-binomial regression (3.4); LogOR, CLRW: log odds ratio, conditional logistic regression (3.5); LogRR, BBIN: log relative risk, beta-binomial regression (3.4); LogRR, RBRR: log relative risk, Cai's Poisson-gamma model, random treatment effect (3.8); RD, BBIN: risk difference, beta-binomial regression (3.4); RD, PQL: risk difference, generalized linear mixed model (estimation by penalized quasi-likelihood) (3.2); AD, CONS: arcsine difference with conservative variance (3.11).

method which ignores 35 double-zero studies, the original analysis misses a clinically relevant (halving stroke occurrence!) and statistically significant advantage of the off-pump method to reduce postoperative stroke occurrence which is discovered by the method that was found superior in our simulation.

Remember that we postulated in Section 1 that ignoring information from double-zero studies would bias the estimate away from the null effect because double-zero studies point to equal treatments. As such, it seems paradoxical that including the information from double-zero studies in this example moves the relative risk *away* from 1 (from 0.76 to 0.51) and not *toward* 1. The solution to this problem is twofold. First, we are comparing results from different models here. The estimates from the beta-binomial model when only including the 25 trials that went into the original analysis are already different from those of the original analysis. To be concrete, the relative risk from the beta-binomial model when including only the 25 non-double-zero studies is 0.53 with a 95% CI of 0.31 to 0.91. However, this estimate is still closer to 1 than the beta-binomial relative risk from all 60 studies. This is because (and this is the second part of the solution to the seeming paradox) not only the treatment estimate itself is affected by including



**Figure 7.** Estimated relative risks with 95% confidence intervals from the original analysis by Møller *et al.* [26] and by the beta-binomial model. Results are given for the 25 studies that have at least one single event (marked with an \*) and for all 60 studies where these estimates also include the information from the double-zero studies.

double-zero studies but also the parameters of the random effects distribution, that is, the beta distribution that describes the event probabilities. And of course, all parameters interfere when maximizing their common likelihood. For example, in the beta-binomial model, we achieve a beta distribution for the event probabilities in the control group with mean 0.028 (95% CI: 0.019 to 0.036) and variance 0.00023 (95% CI: 0.00001 to 0.00044) when including only the 25 non-double-zero studies. When including all 60 studies, this beta distribution has mean 0.018 (95% CI: 0.012 to 0.023) and variance 0.00019 (95% CI: 0.00004 to 0.00034). To sum this up, the differences between the original analysis of Møller *et al.* [26] (standard inverse-variance random effects model ignoring double-zero studies) and the beta-binomial model arise from two sources: first, the model change itself, and second, the inclusion of double-zero studies, where in this example the model change makes a larger contribution than the inclusion of double-zero studies. Nevertheless, and as seen from the simulation, the beta-binomial model (which includes double-zero studies by default) should be used here.

One more point on the length of the CIs deserves mentioning. In the sample of the 25 studies without double-zero studies, the CI (on the log scale) for the original analysis (0.67) is shorter than that for the beta-binomial model (1.08). This reflects again the well-known anti-conservativeness of the standard inverse-variance methods that was seen in all previous simulation studies with sparse data.

Inspired by the original analysis, we focused on the relative risk for the example data set; however, results for the other effect measures are qualitatively similar. Because of the very low number of events, the odds ratio estimate (0.51 [0.28, 0.92]) is indistinguishable from the relative risk estimate; results for the risk difference and the arcsine difference would lead to the same clinical conclusions (RD:  $-0.009$  [ $-0.016, -0.001$ ],  $p = 0.022$ ; AD:  $-0.025$  [ $-0.052, 0.001$ ],  $p = 0.062$ ).

## 6. Discussion

We have shown that there are valid statistical methods for meta-analyses with binary responses that include the information from double-zero studies and that do not need any sort of continuity correction. Such a method is available for each of the commonly encountered effect estimators, that is, the odds ratio, the relative risk, and the risk difference, and also for the more exotic measure of the arcsine difference. Moreover, the current standard method for very sparse data as recommended from the Cochrane collaboration, the Yusuf–Peto odds ratio, can be improved on. As such, this study confirms previous recommendations that methods that do not use information from all studies or use continuity corrections need (or probably even should) no longer be used.

These results were found in a simulation study that tried to mirror real-life situations as completely as possible, where true values for the number of studies per meta-analysis, sample sizes of single studies, and the random effects variances were gathered from empirical data on a large number of actually performed meta-analyses. By including situations with large numbers of double-zero studies, this simulation also covered situations with single-zero studies and very sparse data situations in general. As such, the work reported here should not be seen as only considering the special case of double-zero studies but as a comprehensive study for sparse data in meta-analysis with binary endpoints.

For each of the three common treatment estimators (odds ratio, relative risk, and risk difference), the beta-binomial model (BBIN) was under the two methods that were found to behave satisfactorily under the null hypothesis of no treatment effect. In the simulations with a medium treatment effect, the three BBIN models were comparable or superior in terms of convergence, empirical power, and empirical coverage as compared with their counterparts. As such, and also to keep things simple, we

recommend to use the beta-binomial model for the three common treatment estimators as the preferred method in the meta-analysis with binary responses and double-zero or single-zero studies. If the arcsine difference should be used to assess the treatment effect, we recommend to use it with the conservative variance estimate. A SAS macro that computes the recommended procedures is given in the Supporting information.

It is probably not that surprising that true random effects models came up as the recommended methods from our simulation because these are the more general models as compared with their fixed effect counterparts. As such, our work is also an important extension of the two most complete simulation studies for meta-analysis with rare events up to now, those of Sweeting *et al.* [8] and Bradburn *et al.* [9]. Both authors mainly concentrated on fixed effect methods, although Sweeting *et al.* [8] initially considered 'classical and Bayesian random effects models' but did not give numerical results, whereas Bradburn *et al.* [9] included the standard inverse-variance random effects model with the variance estimator of DerSimonian and Laird and found it behaving poorly in their simulations. It is interesting that our recommended methods are more simple random effects methods that use a conjugacy relation to arrive at a closed-form likelihood. Opposed to this, the more complicated random effect models that assume normally distributed random effects, and thus require more challenging estimation methods (GHQ, MCMC, and PQL), lead to less valid and especially less robust methods. In terms of Bayesian/MCMC estimation, the simulation confirmed the expected problems with the prior distributions which, although determined as being noninformative, dominate the results in sparse data situations. Especially the inverse gamma prior for the random effects variance has been criticized for compromising empirical coverage for effect estimates [50]. However, it is without any doubt that in situations with prior evidence from external sources this information can (and probably should) be included, and this is of course only possible by using Bayesian methods. It might be interesting future work if extensions of the beta-binomial model would yield additional insights. For example, one could model separate beta-binomial distributions for control and treatment groups and link treatment and control groups from the same study by a random effect. Other possible extensions are the zero-inflated beta-binomial model of Hall and Berenhaut [51] and the random-clumped beta-binomial model of Morel and Neerchal [52].

It has been stated that double-zero studies contain no information and thus can be safely removed from the analysis when using the odds ratio [8, Chapter 3.2] or both relative effects measures [53, Chapter 5.2.2]. As stated before, supported by the simulation results as well as with reference to the points given in Section 1, we object to this view. We feel that the recommendations of previous authors concerning the odds ratio and the relative risk are mainly driven by the fact that the estimates' variances in double-zero studies are infinite, and thus their weights in a standard inverse-variance meta-analysis are zero. However, this recommendation overlooks two aspects. First, it points to a peculiarity of the inverse-variance method, because when using the risk difference to describe the treatment effect, double-zero studies would have variance zero, therefore, an indefinite weight, and thus dominate each meta-analysis [43]. This is one more argument against the standard inverse-variance method, because no one would like to have a method that gives zero weight to a study with the odds ratio and the relative risk but infinite weight to the very same study with the risk difference as the treatment effect estimator. Of course, we take a rather dogmatic point of view here, because in applied work researchers would not allow variances of risk differences to become zero but apply some sort of correction. However, this dogmatic view was our study's starting point where we wanted to check how far one could go with methods that work without continuity corrections and use information from all studies regardless of the number of events. Second, the recommendation of deleting double-zero studies overlooks that it relies on a peculiarity of the maximum likelihood principle (which, of course, works pretty well in lots of other situations). When using the median unbiased estimate (MUE) principle for deriving odds ratios [54], relative risks [55], or risk differences [56], balanced double-zero studies indeed give estimates of 1 (OR, RR) or 0 (RD), which is what we expect intuitively. Moreover, in both cases, we also obtain sensible CIs. As such, methods for meta-analysis that use CIs to arrive at summary estimates, for example [24], could also be used to include the information from double-zero studies.

It is important to point to some limitations of our simulation study. We explicitly included only methods that can be computed within a manageable number of code lines and computing time within one major statistical package (SAS). It is possible that there are some other easily accessible methods within other major statistical packages (R or STATA) that would outperform the methods here. Moreover, to enable a fair comparison between methods, we always used the default options in the various SAS procedures and Wald CIs for assessing coverage. It is possible that some methods would have performed better after some tuning of the estimation methods, with better starting values, or with applying other principles (*t*-

distributions with a lower number of degrees of freedom, or profile likelihood) for CI estimation. Finally, we restricted the simulation study to the typical RCT situation with roughly balanced sample sizes within single studies. As such, our recommendations might not be valid for meta-analyses including studies with different sample sizes in treatment arms.

As seen in the simulation study, convergence is an issue for the recommended methods, and it is legitimate to ask what the applied researcher should do in a situation where the recommended methods do not converge or give nonsense estimates for the preferred effect measure. As a work-around for such an extreme situation where most probably only a very low number of events are given, one could always resort to the arcsine difference with the conservative variance; this method guarantees convergence. If results should be communicated with one of the more familiar effect measures, it might be wise to compute this effect measure from the collapsed fourfold table and construct a CI by inverting the  $p$ -value from the test of no effect from the arcsine difference. For the ultimate situation of a meta-analysis *without any event in all studies*, one could collapse the study data and compute MUEs for the odds ratio [54], the relative risk [55], or the risk difference [56]. For example, Salpeter *et al.* [57, Analysis 1.2, page 24] gave nine studies that compared beta-blocker and placebo for the occurrence of single-dose respiratory symptoms. In all studies, no single event was observed from 177 observations in the beta-blocker and 124 in the placebo groups. Collapsing the studies gives a Parzen odds ratio of 0.70 (95% CI: 0.11 to 4.45 by a 'conservative' bootstrap method) and a Carter relative risk of 0.70 (95% CI: 0.11 to 4.40). The MUE estimator for the risk difference amounts to  $-0.0024$  (95% CI:  $-0.0213$  to  $0.0166$ ).

To finally conclude, in meta-analyses with binary outcomes and double-zero or single-zero studies, we recommend to use beta-binomial regression methods to arrive at summary estimates for the odds ratio, the relative risk, or the risk difference. Methods that ignore information from double-zero studies or use continuity corrections should no longer be used. A reviewer pointed to a sensible exception from that rule, which is when plotting effect measures from single studies in forest or funnel plots. In these plots, a continuity correction might be desirable so that each study is guaranteed to appear in the graph. But of course, these corrections should only be made for plotting and removed again for further analysis.

## Acknowledgements

This work was supported by grant KU 1443/3-1 from the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation), Bonn, Germany. We are grateful to Rebecca Turner (MRC Biostatistics Unit, Institute of Public Health, Cambridge, UK) who did some additional analyses from her review data set and shared the results with us.

## References

1. Sutton AJ, Cooper NJ, Lambert PC, Jones DR, Abrams KR, Sweeting MJ. Meta-analysis of rare and adverse event data. *Expert Review of Pharmacoeconomics & Outcomes Research* 2002; **2**(4):367–379.
2. Huang HY, Andrews E, Jones J, Skovron ML, Tilson H. Pitfalls in meta-analyses on adverse events reported from clinical trials. *Pharmacoepidemiology & Drug Safety* 2011; **20**(10):1014–1020.
3. Vandermeer B, Bialy L, Hooton N, Hartling L, Klassen TP, Johnston BC, Wiebe N. Meta-analyses of safety data: a comparison of exact versus asymptotic methods. *Statistical Methods in Medical Research* 2009; **18**(4):421–432.
4. Kuss O, Wandrey M, Kunze M. How frequent are meta-analyses with “double-zero” studies in systematic reviews? 2009. Available from: <http://www.egms.de/static/de/meetings/gmds2009/09gmds155.shtml> [Accessed 7 November 2014].
5. Salpeter SR, Greyber E, Pasternak GA, Salpeter EE. Risk of fatal and nonfatal lactic acidosis with metformin use in type 2 diabetes mellitus. *Cochrane Database Syst Reviews* 2010; **4**:CD002967. DOI: 10.1002/14651858.
6. Brookes G, Ahmed AG. Pharmacological treatments for psychosis-related polydipsia. *Cochrane Database of Systematic Reviews* 2006; **4**:CD003544. DOI: 10.1002/14651858.CD003544.pub2.
7. Nietert PJ, Wahlquist AE, Herbert TL. Characteristics of recent biostatistical methods adopted by researchers publishing in general/internal medicine journals. *Statistics in Medicine* 2013; **32**(1):1–10.
8. Sweeting MJ, Sutton AJ, Lambert PC. What to add to nothing? Use and avoidance of continuity corrections in meta-analysis of sparse data. *Statistics in Medicine* 2004; **23**(9):1351–1375.
9. Bradburn MJ, Deeks JJ, Berlin JA, Russell Localio A. Much ado about nothing: a comparison of the performance of meta-analytical methods with rare events. *Statistics in Medicine* 2007; **26**(1):53–77.
10. Friedrich JO, Adhikari NK, Beyene J. Inclusion of zero total event trials in meta-analyses maintains analytic consistency and incorporates all available data. *BMC Medical Research Methodology* 2007; **7**:5. DOI: 10.1186/1471-2288-7-5.

11. Liu D. Combining information for heterogeneous studies and rare events studies: a confidence distribution approach, 2012. Available from: <https://rucore.libraries.rutgers.edu/rutgers-lib/37435/> [Accessed 7 November 2014].
12. Keus F, Wetterslev J, Gluud C, Gooszen HG, van Laarhoven CJ. Robustness assessments are needed to reduce bias in meta-analyses that include zero-event randomized trials. *The American Journal of Gastroenterology* 2009; **104**(3): 546–551.
13. Agresti A, Hartzel J. Strategies for comparing treatments on a binary response with multi-centre data. *Statistics in Medicine* 2000; **19**(8):1115–1139.
14. Hirji KF. *Exact Analysis of Discrete Data*. Chapman & Hall/ CRC: Boca Raton, FL, USA, 2006.
15. Kuss O, Gummert JF, Borgermann J. Meta-analyses with rare events should use adequate methods. *The Journal of Thoracic and Cardiovascular Surgery* 2008; **136**(1):241.
16. Eaton S, Hall NJ, Pierro A. Zero-total event trials and incomplete pyloromyotomy. *The Journal of Pediatric Surgery* 2009; **44**(12):2434–2435.
17. Greenland S. Simpson's paradox from adding constants in contingency tables as an example of Bayesian noncollapsibility. *The American Statistician* 2010; **64**(4):340–344.
18. Breslow N. Odds ratio estimators when the data are sparse. *Biometrika* 1981; **68**(1):73–84.
19. O'Gorman TW, Woolson RF, Jones MP, Lemke JH. A Monte Carlo study of three odds ratio estimators and four tests of association in several  $2 \times 2$  tables when the data are sparse. *Communication in Statistics - Simulation and Computation* 1988; **17**(3):813–835.
20. Sankey SS, Weissfeld LA, Fine MJ, Kapoor W. An assessment of the use of the continuity correction for sparse data in meta-analysis. *Communication in Statistics - Simulation and Computation* 1996; **25**(4):1031–1056.
21. Lui KJ. A Monte Carlo evaluation of five interval estimators for the relative risk in sparse data. *Biometrical Journal* 2006; **48**(1):131–143.
22. Kuss O, Gromann C. An exact test for meta-analysis with binary endpoints. *Methods of Information in Medicine* 2007; **46**(6):662–668.
23. Higgins JPT, Green S (eds). *Cochrane Handbook for Systematic Reviews of Interventions Version 5.1.0 [updated March 2011]*. The Cochrane Collaboration, 2011. Available from: [www.cochrane-handbook.org](http://www.cochrane-handbook.org) [Accessed on 26 November 2014].
24. Tian L, Cai T, Pfeffer MA, Piankov N, Cremieux PY, Wei LJ. Exact and efficient inference procedure for meta-analysis and its application to the analysis of independent  $2 \times 2$  tables with all available data but without artificial continuity correction. *Biostatistics* 2009; **10**(2):275–281.
25. Stijnen T, Hamza TH, Ozdemir P. Random effects meta-analysis of event outcome in the framework of the generalized linear mixed model with applications in sparse data. *Statistics in Medicine* 2010; **29**(29):3046–3067.
26. Möller CH, Penninga L, Wetterslev J, Steinbrüchel DA, Gluud C. Off-pump versus on-pump coronary artery bypass grafting for ischaemic heart disease. *Cochrane Database Systematic Reviews* 2012; **3**:CD007224. DOI: 10.1002/14651858.CD007224.pub2.
27. Altman DG, Deeks JJ. Meta-analysis, Simpson's paradox, and the number needed to treat. *BMC Medical Research Methodology* 2002; **2**:3. DOI: 10.1186/1471-2288-2-3.
28. Lièvre M, Cucherat M, Leizorovicz A. Pooling, meta-analysis, and the evaluation of drug safety. *Current Controlled Trials in Cardiovascular Medicine* 2002; **3**(1):6.
29. Platt RW, Leroux BG, Breslow N. Generalized linear mixed models for meta-analysis. *Statistics in Medicine* 1999; **18**(6):643–654.
30. McCulloch CE, Searle SR, Neuhaus JM. *Generalized, Linear, and Mixed Models*. John Wiley & Sons, Inc.: Hoboken, New Jersey, USA, 2008.
31. Gao S. Combining binomial data using the logistic normal model. *Journal of Statistical Computation and Simulation* 2006; **74**(4):293–306.
32. Agresti A. *Categorical Data Analysis*. John Wiley & Sons, Inc.: Hoboken, New Jersey, USA, 2002.
33. Molenberghs G, Verbeke G. *Models for Discrete Longitudinal Data*, Springer Series in Statistics. Springer: New York, USA, 2005.
34. Fleiss JL, Levin B, Paik MC. *Statistical Methods for Rates and Proportions*. John Wiley & Sons, Inc.: Hoboken, New Jersey, USA, 2003.
35. Lin DY, Wei LJ. The robust inference for the proportional hazards model. *Journal of the American Statistical Association* 1989; **84**(408):1074–1078.
36. Firth D. Bias reduction of maximum likelihood estimates. *Biometrika* 1993; **80**(1):27–38.
37. Heinze G. A comparative investigation of methods for logistic regression with separated or nearly separated data. *Statistics in Medicine* 2006; **25**(24):4216–4226.
38. Cai T, Parast L, Ryan L. Meta-analysis for rare events. *Statistics in Medicine* 2010; **29**(20):2078–2089.
39. Greenland S, Robins JM. Estimation of a common effect parameter from sparse follow-up data. *Biometrics* 1985; **41**(1): 55–68.
40. Sato T. On the variance estimator for the Mantel-Haenszel risk difference (letter). *Biometrics* 1989; **45**(4):1323–1324.
41. Kuhnert R, Böhning D. The failure of meta-analytic asymptotics for the seemingly efficient estimator of the common risk difference. *Statistical Papers* 2005; **46**(4):541–554.
42. Newman SC. *Biostatistical Methods in Epidemiology*. John Wiley & Sons, Inc.: Hoboken, New Jersey, USA, 2001.
43. Rücker G, Schwarzer G, Carpenter J, Olkin I. Why add anything to nothing? The arcsine difference as a measure of treatment effect in meta-analysis with zero cells. *Statistics in Medicine* 2009; **28**(5):721–738.
44. Kuss O, Hoyer A, Solms A. Meta-analysis for diagnostic accuracy studies: a new statistical model using beta-binomial distributions and bivariate copulas. *Statistics in Medicine* 2014; **33**(1):17–30.



45. Turner RM, Davey J, Clarke MJ, Thompson SG, Higgins JP. Predicting the extent of heterogeneity in meta-analysis, using empirical data from the Cochrane Database of Systematic Reviews. *International Journal of Epidemiology* 2012; **41**(3):818–827.
46. Moher D, Tetzlaff J, Tricco AC, Sampson M, Altman DG. Epidemiology and reporting characteristics of systematic reviews. *PLoS Medicine* 2007; **4**(3):e78. DOI: 10.1371/journal.pmed.0040078.
47. Fleishman AI. A method for simulating non-normal distributions. *Psychometrika* 1978; **43**(4):521–532.
48. Fan X, Felsövályi A, Sivo SA, Kennan SC. *SAS for Monte Carlo Studies: A Guide for Quantitative Researchers*. SAS Institute: Cary, NC, USA, 2002.
49. Sato T. Bias in the Peto one-step estimator for the common odds ratio. *Bulletin of Informatics and Cybernetics* 2005; **37**:13–18.
50. Lambert PC, Sutton AJ, Burton PR, Abrams KR, Jones DR. How vague is vague? A simulation study of the impact of the use of vague prior distributions in MCMC using WinBUGS. *Statistics in Medicine* 2005; **24**(15):2401–2428.
51. Hall DB, Berenhaut KS. Score tests for heterogeneity and overdispersion in zero-inflated Poisson and binomial regression models. *The Canadian Journal of Statistics* 2002; **30**(3):415–430.
52. Morel JG, Neerchal NK. Clustered binary logistic regression in teratology data using a finite mixture distribution. *Statistics in Medicine* 1997; **16**(24):2843–2853.
53. Fu R, Gartlehner G, Grant M, Shamliyan T, Sedrakyan A, Wilt TJ, Griffith L, Oremus M, Raina P, Ismaila A, Santaguida P, Lau J, Trikalinos TA. Conducting quantitative synthesis when comparing medical interventions: AHRQ and the Effective Health Care Program. *Journal of Clinical Epidemiology* 2011; **64**(11):1187–1197.
54. Parzen M, Lipsitz S, Ibrahim J, Klar N. An estimate of the odds ratio that always exists. *Journal of Computational and Graphical Statistics* 2002; **11**(2):420–436.
55. Carter RE, Lin Y, Lipsitz SR, Newcombe RG, Hermayer KL. Relative risk estimated from the ratio of two median unbiased estimates. *Applied Statistics* 2010; **59**(4):657–671.
56. Lin Y, Newcombe RG, Lipsitz S, Carter RE. Fully specified bootstrap confidence intervals for the difference of two independent binomial proportions based on the median unbiased estimator. *Statistics in Medicine* 2009; **28**(23):2876–2890.
57. Salpeter S, Ormiston T, Salpeter E. Cardiorespective beta-blockers for chronic obstructive pulmonary disease. *Cochrane Database Syst Reviews* 2005; **4**:CD003566. DOI: 10.1002/14651858.CD003566.pub2.

## Supporting information

Additional supporting information may be found in the online version of this article at the publisher's web-site.