

Bayesian random effects meta-analysis of trials with binary outcomes: methods for the absolute risk difference and relative risk scales

D. E. Warn^{*,†}, S. G. Thompson and D. J. Spiegelhalter

*MRC Biostatistics Unit, Institute of Public Health, University Forvie Site, Robinson Way,
Cambridge, CB2 2SR, U.K.*

SUMMARY

When conducting a meta-analysis of clinical trials with binary outcomes, a normal approximation for the summary treatment effect measure in each trial is inappropriate in the common situation where some of the trials in the meta-analysis are small, or the observed risks are close to 0 or 1. This problem can be avoided by making direct use of the binomial distribution within trials. A fully Bayesian method has already been developed for random effects meta-analysis on the log-odds scale using the BUGS implementation of Gibbs sampling. In this paper we demonstrate how this method can be extended to perform analyses on both the absolute and relative risk scales. Within each approach we exemplify how trial-level covariates, including underlying risk, can be considered. Data from 46 trials of the effect of single-dose ibuprofen on post-operative pain are analysed and the results contrasted with those derived from classical and Bayesian summary statistic methods. The clinical interpretation of the odds ratio scale is not straightforward. The advantages and flexibility of a fully Bayesian approach to meta-analysis of binary outcome data, considered on an absolute risk or relative risk scale, are now available. Copyright © 2002 John Wiley & Sons, Ltd.

KEY WORDS: meta-analysis; binary data; absolute risk difference; relative risk; Bayesian methods

1. INTRODUCTION

Recent advances in computational methods have encouraged the application of Bayesian methods to random effects meta-analysis. The key difference between the Bayesian approach and the classical, or frequentist, viewpoint lies in the incorporation of subjective, or data-based, prior beliefs into the analysis of data. There are strong reasons for favouring a fully Bayesian approach over non-Bayesian methods, not least the ability to account for all parameter uncertainty [1–3]. Bayesian methods offer a unified modelling framework which overcomes issues such as the appropriate treatment of small trials, and a flexibility which allows the approach to be extended to consider distributions other than Normal for the random effects, or to adjust

*Correspondence to: David Warn, MRC Biostatistics Unit, Institute of Public Health, University Forvie Site, Robinson Way, Cambridge, CB2 2SR, U.K.

†E-mail: david.warn@mrc-bsu.cam.ac.uk

for covariates through regression models. Other advantages include the opportunity to ‘borrow strength’ from other trials [4], and to make predictions about the outcomes of future trials. Furthermore, because the posterior distribution is produced by simulation, via Markov chain Monte Carlo techniques, inference regarding non-standard functions of the parameters is possible.

It is common that many, if not most, of the trials included in a meta-analysis are small. Small trials, and trials for which the observed risks are close to 0 or 1, present problems for methods based on summary statistics (such as log-odds ratios). The usual approach to meta-analysis of trials with binary outcome data, which facilitates estimation, is to assume that these summary statistics have an approximate normal likelihood. However, when a trial is small, such an assumption may not be tenable, and in practice it is rarely validated. Clearly, any method which models binomial outcome data directly is preferable to a summary statistic approach. Smith *et al.* [1] do exactly this, though their approach is presented for the log-odds ratio scale only. Efforts to extend the method to embrace other treatment effect scales, such as the absolute risk difference and the log relative risk, have hitherto proved difficult. Thompson *et al.* [5] stated that such analyses on scales other than the log-odds ratio, though possible in principle, were ‘currently difficult in practice’.

Carlin [6] has indicated a possible solution on the absolute risk scale. This method places a bivariate normal distribution on the treatment effect, defined as the log-odds ratio, and the underlying risk, the average of the treatment and control group log-odds. Transforming the average treatment effect and underlying risk back to the absolute scale yields an estimate of the risk difference. An estimate of the relative risk could be obtained in a similar manner. A weakness of this approach is that the estimates are based on the average underlying risk on the log-odds scale, assuming that the average log-odds ratio is constant across trials. More critically, the heterogeneity of the treatment effect remains on the log-odds scale throughout. Consequently, no estimate of the between-trial variation on the risk difference or relative risk scale is available.

In this paper, we demonstrate how the method of Smith *et al.* [1] can be modified to make such analyses possible. The three common treatment effect measures for binomial outcome data, namely the absolute risk difference (RD), the (log) relative risk (LRR) and (log) odds ratio (LOR), vary in their mathematical properties, their ease of interpretation and their consistency, that is the stability of the treatment effect estimates over the various populations from which the trials have been drawn [7]. The choice of which effect scale is most appropriate in a given situation is rarely straightforward, and so it is important to have available methodology for all possible scales.

In contrast to the odds ratio and relative risk which measure relative effects, the risk difference is a measure of absolute effect. It can be interpreted as the difference in the proportion of events observed with treatment or, at an individual level, as an estimate of the change in a patient’s probability of experiencing the event. Its reciprocal, the number needed to treat (NNT) (to prevent one outcome) is widely used to describe treatment effects in practice since it has a clinically attractive interpretation [8–11]. An advantage of the fully Bayesian approach is that it is possible to obtain a posterior distribution for the NNT. The risk ratio describes the factor by which the risk of the event is multiplied by the use of the treatment. The odds ratio is more difficult to interpret but benefits from mathematical advantages. For example, values of both the risk difference and relative risk are constrained, the risk difference between -1 and 1 , and the risk ratio bounded above in a manner dependent on control group risk.

The odds ratio can take any value between 0 and infinity, and on the logarithmic scale it is unbounded in both directions.

In this paper, we show how the risk difference and relative risk can be constrained within their permissible bounds thereby enabling the extension of Smith *et al.*'s method to these treatment effect scales. In Section 2, we describe a real example for which the most appropriate effect measure might be any of the three scales. Section 3 describes the general Bayesian random effects model and the necessary modifications of the existing method. The results of the analyses are presented for the example in Section 4. Section 5 describes how each of the models may be extended to adjust for trial-level covariates, including the special case of underlying risk, and incorporates a brief discussion on informative priors for the between-trial variance within the random effects model. Finally, Section 6 reviews all the methods presented here and discusses their application within the context of general meta-analysis methodology.

2. AN EXAMPLE: SINGLE-DOSE IBUPROFEN FOR POST-OPERATIVE PAIN

Data are available from a Cochrane Review investigating the effectiveness of single-dose ibuprofen in reducing post-operative pain [12]. Ibuprofen is one of a class of non-steroidal anti-inflammatory (NSAID) analgesics and it is important to know which drug and dose should be recommended for post-operative pain relief. The review comprises 46 small trials of single-dose ibuprofen against placebo with binomial outcome data. The dose used in different trials ranges from 50 mg to 800 mg. A measure of 'at least 50 per cent pain relief' in the 4–6 hours after administration of the dose is used as the common descriptor of analgesic efficacy. Since the length of follow-up is the same across trials, it is appropriate to consider the patient's 'risk' of experiencing pain relief. The data from the trials are given in Table I.

For the majority of this paper, we consider just the 31 trials with dose 400 mg. Not surprisingly, there is considerable evidence that ibuprofen improves pain relief, but heterogeneity of its effect is evident on all three scales (Figure 1). In Section 5, we will consider all the trials as we describe how to investigate the relationship between treatment effect and dose.

3. METHODS

The methods described in this section are presented in the context of meta-analysis of clinical trials. They are also appropriate for other applications of hierarchical modelling to binary outcome data, for example cluster randomized trials or multi-centre trials. The notation used in the following methods is given in Table II.

The L'Abbé plot [13] is widely used as a graphical tool in meta-analysis. The observed treatment group risk is plotted against the observed control group risk for each trial. The distribution of the plotted points indicates the amount of heterogeneity present [14] and can help decide which effect measure provides the best overall summary for a meta-analysis [7]. A L'Abbé plot of the 31 trials of 400 mg ibuprofen is shown in Figure 2.

Our approach involves modelling the joint distribution of the points of a L'Abbé plot of the 'true' treatment and control group risks. To achieve this in a Bayesian framework, we require a joint prior distribution which is non-informative about the 'true' risks but allows

Table I. Data from 46 trials of single-dose ibuprofen on post-operative pain. var is the estimated variance of the quantity in the previous column; x = dose in mg; other notation given in Table II.

Trial	r^T	n^T	r^C	n^C	RD	var	LRR	var	LOR	var	x
1 Forbes 1991a	16	57	0	51	0.275	0.004	3.387	2.024	3.712	2.104	50
2 Forbes 1991a	13	49	0	51	0.260	0.004	3.335	2.035	3.640	2.121	100
3 Jain 1986	3	39	0	47	0.077	0.002	2.128	2.240	2.209	2.334	100
4 Cooper 1977	15	38	6	40	0.245	0.009	0.968	0.182	1.307	0.306	200
5 Forbes 1991a	18	48	0	51	0.368	0.005	3.670	2.014	4.135	2.106	200
6 Hersch 1993a	13	51	0	51	0.250	0.004	3.296	2.036	3.587	2.119	200
7 Jain 1986	7	47	0	47	0.146	0.003	2.708	2.092	2.867	2.179	200
8 Kiersch 1993	37	81	4	42	0.362	0.005	1.568	0.241	2.078	0.326	200
9 McQuay 1996a	2	31	0	11	0.036	0.006	0.629	2.285	0.667	2.521	200
10 Nelson 1994a	43	75	8	40	0.373	0.007	1.053	0.110	1.682	0.211	200
11 Seymour 1996	18	35	2	19	0.409	0.012	1.586	0.474	2.197	0.673	200
12 Ahlstrom 1993	19	32	2	30	0.527	0.010	2.187	0.488	3.019	0.665	400
13 Arnold 1990	2	15	0	14	0.123	0.010	1.545	2.271	1.681	2.543	400
14 Bakshi 1994	57	80	31	82	0.334	0.005	0.634	0.025	1.405	0.113	400
15 Cooper 1977	20	40	6	40	0.350	0.009	1.204	0.167	1.735	0.296	400
16 Cooper 1982	22	38	5	46	0.470	0.009	1.673	0.197	2.423	0.332	400
17 Cooper 1988a	19	37	6	43	0.374	0.010	1.303	0.169	1.873	0.302	400
18 Cooper 1989	37	61	9	64	0.466	0.006	1.462	0.106	2.243	0.198	400
19 Forbes 1984	21	28	3	28	0.643	0.010	1.946	0.310	3.219	0.564	400
20 Forbes 1990	15	32	0	34	0.455	0.008	3.493	2.006	4.113	2.151	400
21 Forbes 1991b	18	37	1	39	0.461	0.007	2.943	1.003	3.584	1.135	400
22 Forbes 1992	20	38	0	38	0.513	0.007	3.714	1.997	4.446	2.129	400
23 Frame 1989	26	42	0	38	0.603	0.006	3.873	1.989	4.818	2.124	400
24 Fricke 1993	40	81	2	39	0.443	0.004	2.265	0.487	2.893	0.576	400
25 Gay 1996	9	41	7	39	0.040	0.008	0.201	0.204	0.251	0.316	400
26 Heidrich 1985	15	40	5	40	0.250	0.009	1.099	0.217	1.435	0.335	400
27 Hersch 1993a	19	49	0	51	0.380	0.005	3.703	2.012	4.187	2.103	400
28 Hersch 1993b	9	12	6	16	0.375	0.030	0.693	0.132	1.609	0.711	400
29 Jain 1986	9	49	0	47	0.180	0.003	2.904	2.064	3.104	2.151	400
30 Jain 1988	33	49	17	48	0.319	0.009	0.643	0.048	1.325	0.184	400
31 Laska 1986	39	39	14	37	0.606	0.007	0.951	0.043	4.852	2.137	400
32 Lavenziana 1996	29	42	24	41	0.105	0.011	0.165	0.028	0.458	0.212	400
33 McQuay 1996a	6	30	0	11	0.168	0.009	1.616	2.038	1.809	2.282	400
34 Mehilisch 1990	124	306	5	85	0.346	0.001	1.930	0.193	2.389	0.226	400
35 Mehilisch 1995	67	98	1	40	0.659	0.003	3.309	0.980	4.434	1.073	400
36 Pagnoni	13	30	5	32	0.277	0.012	1.020	0.212	1.418	0.373	400
37 Schachtel 1989	27	36	13	38	0.408	0.011	0.785	0.060	1.753	0.265	400
38 Seymour 1991(I)	42	63	10	32	0.354	0.010	0.758	0.077	1.482	0.217	400
39 Seymour 1991(II)	42	62	7	30	0.444	0.009	1.066	0.117	1.932	0.260	400
40 Seymour 1996	22	31	2	19	0.604	0.012	1.908	0.461	3.034	0.715	400
41 Sunshine 1983	21	30	3	30	0.600	0.010	1.946	0.314	3.045	0.529	400
42 Sunshine 1987	16	38	11	40	0.146	0.011	0.426	0.102	0.651	0.233	400
43 Parker 1986	35	44	22	33	0.129	0.010	0.177	0.021	0.665	0.276	600
44 Laska 1986	36	36	14	37	0.605	0.007	0.950	0.043	4.773	2.139	600
45 Seymour 1986	19	34	2	19	0.454	0.012	1.669	0.471	2.376	0.678	600
46 Laska 1986	39	39	14	37	0.606	0.007	0.951	0.043	4.852	2.137	800

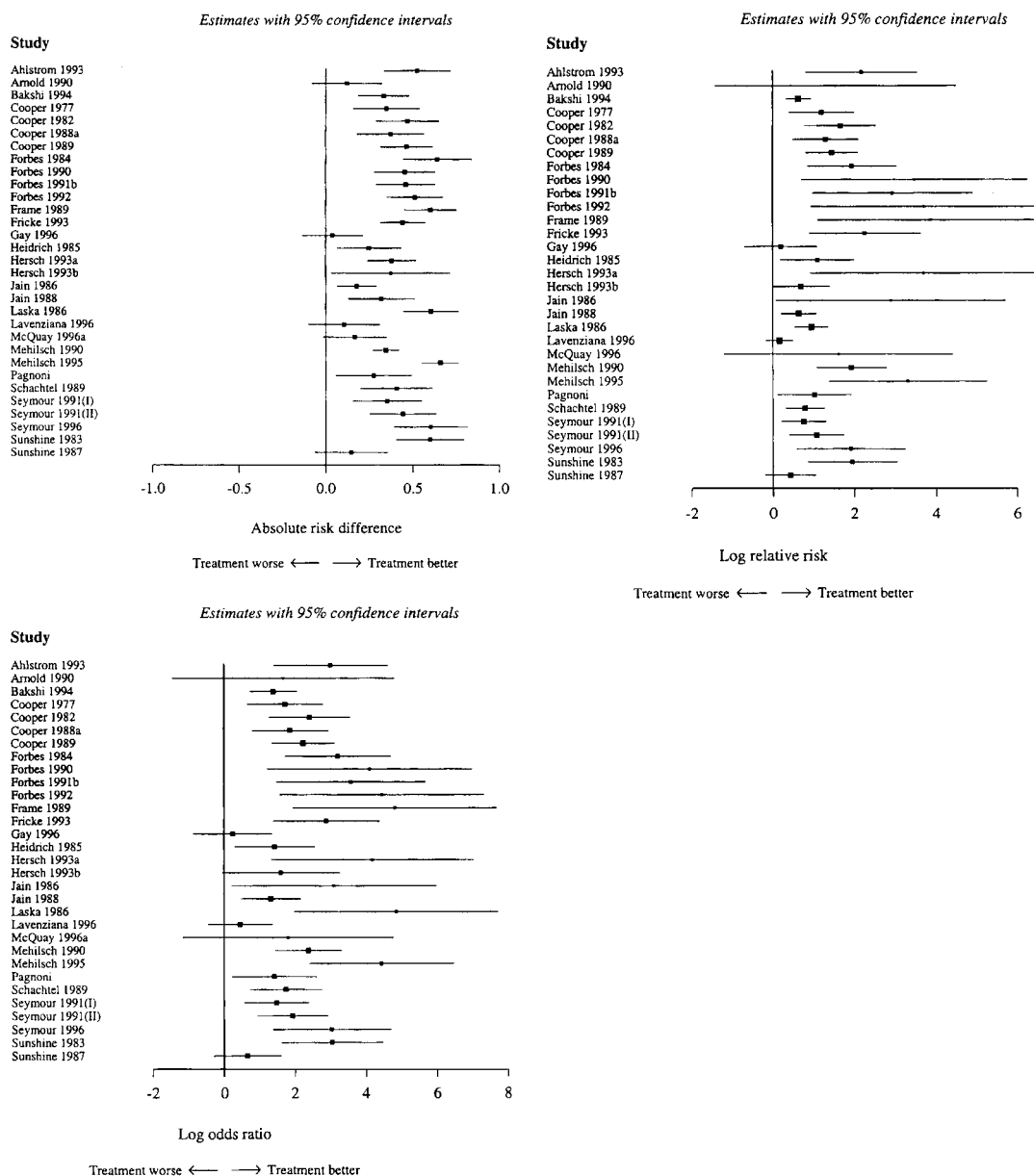


Figure 1. Confidence interval plot of results for 400 mg ibuprofen trials on risk difference, log relative risk and log-odds ratio scales.

for some unknown correlation between them. On the absolute risk difference scale, ideally, this equates to a structure which provides a correlated prior on the unit square with uniform marginal distributions. We describe an approach that closely approximates this situation.

Table II. Table of notation.

i	indexes trials
k	number of trials
r_i^C	number of events in the control group of trial i
r_i^T	number of events in the treatment group of trial i
n_i^C	number in control group of trial i
n_i^T	number in treatment group of trial i
p_i^C	observed control group risk in trial i
p_i^T	observed treatment group risk in trial i
d_i	observed treatment effect in trial i
π_i^C	true control group risk in trial i
π_i^T	true treatment group risk in trial i
δ_i	true treatment effect in trial i
δ_i^L	lower bound for δ_i
δ_i^U	upper bound for δ_i
δ	average treatment effect
τ^2	variance of δ_i
μ_i	true underlying risk in trial i
α, β	parameters of beta distribution
γ	regression slope

3.1. The absolute risk difference scale

The Bayesian random effects model of Smith *et al.* [1, 15] may be written as follows:

$$\begin{aligned}
 r_i^C &\sim \text{Bin}(n_i^C, \pi_i^C) \\
 r_i^T &\sim \text{Bin}(n_i^T, \pi_i^T) \\
 \mu_i &= \text{logit}(\pi_i^C) \\
 \text{logit}(\pi_i^T) &= \mu_i + \delta_i \\
 \delta_i &\sim \text{N}(\delta, \tau^2)
 \end{aligned} \tag{1}$$

in which $\delta_i = \text{logit}(\pi_i^T) - \text{logit}(\pi_i^C)$ is the log-odds ratio. Adopting the same approach and notation, an equivalent model on the absolute risk scale replaces the logits of the treatment and control group risks by the risks themselves:

$$\begin{aligned}
 \mu_i &= \pi_i^C \\
 \pi_i^T &= \mu_i + \delta_i \\
 \delta_i &\sim \text{N}(\delta, \tau^2)
 \end{aligned} \tag{2}$$

δ_i is now the absolute risk difference. To ensure the model has appropriate parameter values, δ_i needs to be constrained so that $\pi_i^T \in [0, 1]$, achieved by confining δ_i to the interval

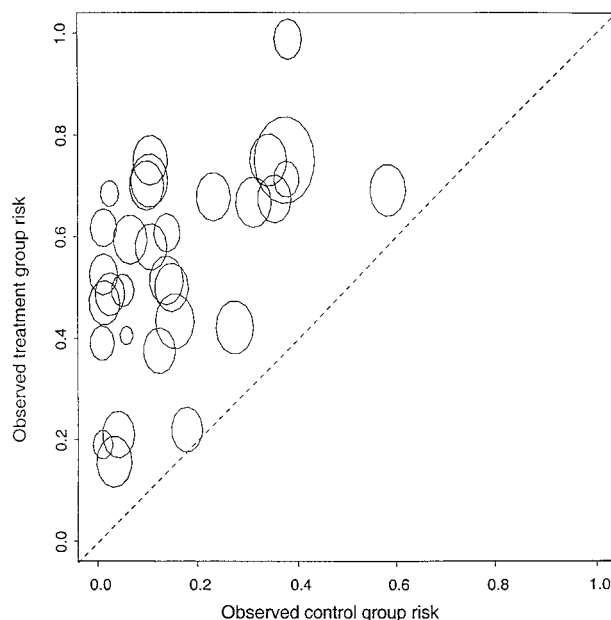


Figure 2. L'Abbé plot of observed risks for 31 trials of 400mg ibuprofen on the absolute risk scale. Each trial is represented by a symbol of area proportional to its precision (inverse variance of absolute risk difference).

$[-\pi_i^C, 1 - \pi_i^C]$. Since each δ_i is drawn from a normal distribution with mean δ and variance τ^2 , it can take any value in the range $(-\infty, \infty)$. We define two new parameters, δ_i^U and δ_i^L , corresponding, respectively, to upper and lower bounds for δ_i . If we let δ_i^L be the maximum of δ_i and $-\pi_i^C$, then δ_i^L can take any value in the range $[-\pi_i^C, \infty)$. Similarly, if we let δ_i^U be the minimum of δ_i^L and $1 - \pi_i^C$, then δ_i^U is confined to the required range $[-\pi_i^C, 1 - \pi_i^C]$. The full model is then

$$\begin{aligned}
 \mu_i &= \pi_i^C \\
 \pi_i^T &= \mu_i + \delta_i^U \\
 \delta_i^U &= \min(\delta_i^L, 1 - \pi_i^C) \\
 \delta_i^L &= \max(\delta_i, -\pi_i^C) \\
 \delta_i &\sim N(\delta, \tau^2)
 \end{aligned} \tag{3}$$

Dispensing with the new parameters, we may rewrite (3) as

$$\begin{aligned}
 \mu_i &= \pi_i^C \\
 \pi_i^T &= \mu_i + \min(\max(\delta_i, -\pi_i^C), 1 - \pi_i^C) \\
 \delta_i &\sim N(\delta, \tau^2)
 \end{aligned} \tag{4}$$

A fully Bayesian analysis requires priors to be placed on the unknown parameters δ and τ^2 . δ is given a $\text{Unif}(-1, 1)$ prior distribution which gives equal weight to all possible values of the parameter. A $\text{Unif}(0, 2)$ prior was placed on τ since this range contains all plausible values of τ , the between-trial standard deviation of the true risk difference. We must also place prior distributions on the μ_i , in this case equal to π_i^C . For consistency across the different treatment effect scales considered in this paper, we will always consider the prior distribution of the π_i^C . Two approaches are adopted here: method (a) places independent $\text{Unif}(0, 1)$ prior distributions on the π_i^C , equivalent to treating them as ‘fixed effects’ so that each π_i^C is unrelated to the others; method (b) assumes the π_i^C are drawn from a common ‘random effects’ distribution [16–18]. For (b), we use a beta prior distribution with hyperparameters α and β , with a $\text{Unif}(1, 100)$ hyperprior on each. This ensures that the distribution of the π_i^C is unimodal, but not unduly constrained otherwise. Method (b) is analogous to the bivariate approach of van Houwelingen *et al.* [19] which also models the binary outcome data directly and imposes a common random-effects distribution on the (log-odds of the) true group risks.

The purpose of the $\min(\max(\delta_i, -\pi_i^C), (1 - \pi_i^C))$ construction in (4) is to ensure that the probability π_i^T is bounded between 0 and 1. Whenever δ_i is sampled outside the range $[-\pi_i^C, 1 - \pi_i^C]$, π_i^T is set to 0 or 1, and so we observe spikes of probability at these values. The size of each spike is determined by the location of δ and the size of τ . For method (a), the implied prior distribution for π_i^T is the sum of a uniform and a normal distribution, and is approximately uniform on the interval $(0, 1)$. The spikes at 0 and 1 are undesirable insofar that support for these values in the likelihood leads to spikes of probability in the posterior distribution for the π_i^T . This occurs when the observed treatment group risk of a trial r_i^T/n_i^T is exactly 0 or 1. However, since it is implausible that the true treatment group risk is precisely 0 or 1, it is reasonable to adjust the likelihood so that the support for these values is withdrawn. In practice this might be achieved by substituting, in the data, 0.0001 for the number of events when $r_i^T = 0$ and $n_i^T - 0.0001$ when $r_i^T = n_i^T$. In this way, we can ensure that the prior is approximately uniform over the range plausibly supported by the likelihood.

We note that an observed treatment group risk of 1 occurs in the ibuprofen trials. In trial 31 (Table I), all 39 patients in the treatment group experienced the event. By monitoring the posterior distribution of π_{31}^T , we may examine sensitivity to the choice of substitute value $n_{31}^T - \varepsilon$ for n_{31}^T in the data. For $\varepsilon = 0$, π_{31}^T takes the value 1 around 73 per cent of the time and so has posterior median 1. For $\varepsilon = 0.01$, 0.0001 and 0.000001, the posterior median of π_{31}^T is in each case 0.971. This indicates that the posterior mean is robust to different small $\varepsilon > 0$. Moreover, the effect of the prior probability spike at 1 is clear from the larger posterior median obtained when no substitution is made in the data.

3.2. The log relative risk scale

For the log relative risk scale, we proceed as above, but replace the logits in (1) by logs so that δ_i is the log relative risk:

$$\begin{aligned}\mu_i &= \log(\pi_i^C) \\ \log(\pi_i^T) &= \mu_i + \delta_i \\ \delta_i &\sim N(\delta, \tau^2)\end{aligned}\tag{5}$$

As before δ_i needs to be constrained so that $\pi_i^T \in [0, 1]$. This is equivalent to constraining $\log(\pi_i^T)$ to the interval $(-\infty, 0]$, achieved by confining δ_i to be less than $-\log(\pi_i^C)$. As in Section 3.1, if we let δ_i^U be the minimum of δ_i and $-\log(\pi_i^C)$, then δ_i^U can take any value in the range $(-\infty, -\log(\pi_i^C))$. The full model is then

$$\begin{aligned}\mu_i &= \log(\pi_i^C) \\ \log(\pi_i^T) &= \mu_i + \delta_i^U \\ \delta_i^U &= \min(\delta_i, -\log(\pi_i^C)) \\ \delta_i &\sim N(\delta, \tau^2)\end{aligned}\tag{6}$$

which we may rewrite as

$$\begin{aligned}\mu_i &= \log(\pi_i^C) \\ \log(\pi_i^T) &= \mu_i + \min(\delta_i, -\log(\pi_i^C)) \\ \delta_i &\sim N(\delta, \tau^2)\end{aligned}\tag{7}$$

Prior distributions are needed for δ and τ^2 . δ is given a vague $N(0, 10)$ distribution and a $\text{Unif}(0, 2)$ prior placed on τ . This distribution covers all plausible values of τ in this case though other examples may require a wider range. The priors for the μ_i are dealt with as before by considering the π_i^C .

3.3. The log-odds ratio scale

For comparison, we include a method for the log-odds scale. This is a version of Smith *et al.*'s [1] approach (1) modified to enable direct comparison with the two methods described above. Whereas Smith *et al.* specify a prior distribution for the μ_i directly, we employ the same procedure as above by placing priors on the π_i^C . We also replace Smith's inverse-gamma $\text{IG}(3, 1)$ prior for τ^2 with a $\text{Unif}(0, 2)$ for τ , which again covers the range of plausible values for τ in this example.

3.4. A Bayesian method for summary statistic data

We compare the results from the Bayesian methods for binomial outcome data with both a Bayesian method, which we will call method (c), and classical methods for summary statistic data. A simple Bayesian random effects meta-analysis is possible within the framework of a two-level normal-normal hierarchical model by summarizing the results from each trial with an approximate normal likelihood for the treatment effect [20–23]. We use standard approaches for producing a point estimate and standard error for the treatment effect in each trial, which can be regarded as approximating a normal mean and standard deviation.

When either the treatment or control group of a trial contains either no events or no non-events, the problem of a zero cell arises. Zero cells prevent computation of both the relative treatment effect measures and the standard errors of both these and that of the risk difference. The usual way [24] to circumvent this problem, and that adopted in this paper, is to add 0.5 to the count in each cell of all 2×2 tables containing zero cells prior to analysis. The formulae [7] used to derive the effect estimates and standard errors are given in Table III.

Table III. Calculating summary effect estimates and their standard errors.

δ_i	d_i	$SE(d_i)$
RD	$\frac{r_i^T}{n_i^T} - \frac{r_i^C}{n_i^C}$	$\sqrt{\left\{ \frac{r_i^T(n_i^T - r_i^T)}{(n_i^T)^3} + \frac{r_i^C(n_i^C - r_i^C)}{(n_i^C)^3} \right\}}$
LRR	$\log\left(\frac{r_i^T}{n_i^T}\right) - \log\left(\frac{r_i^C}{n_i^C}\right)$	$\sqrt{\left(\frac{1}{r_i^T} + \frac{1}{r_i^C} - \frac{1}{n_i^T} - \frac{1}{n_i^C}\right)}$
LOR	$\log\left(\frac{r_i^T(n_i^C - r_i^C)}{r_i^C(n_i^T - r_i^T)}\right)$	$\sqrt{\left(\frac{1}{r_i^T} + \frac{1}{n_i^T - r_i^T} + \frac{1}{r_i^C} + \frac{1}{n_i^C - r_i^C}\right)}$

Let d_i be the point estimate of the true treatment effect δ_i (RD, LRR, LOR) in trial i . We assume that the standard error of d_i is known and that

$$d_i \sim N(\delta_i, SE(d_i)^2) \quad (8)$$

Following the general random effects approach [25], the δ_i are assumed to be drawn from a common normal distribution

$$\delta_i \sim N(\delta, \tau^2) \quad (9)$$

Such a model is easily implemented. We use the same prior distributions for δ and τ^2 as in the binomial outcome data methods.

3.5. Implementation

In this paper, the Bayesian analyses were performed using the software package BUGS [26, 27], developed for carrying out Bayesian inference on statistical problems using Gibbs sampling. The BUGS code for implementing each of the models described above is given in the Appendix.

BUGS is based on a graphical modelling strategy in which model quantities and their conditional independence structure are specified using a directed acyclic graph (DAG). In a DAG, all the data and parameters in the model are represented by linked nodes. Each node is represented by an object in the BUGS language, and a model specified by declaring all the nodes and describing the links between them. From these statements, BUGS automatically generates the code to carry out the necessary Gibbs sampling. A DAG for the absolute risk scale analysis is shown in Figure 3.

For each BUGS analyses initial values were specified. Following a burn-in of 5000 iterations, 15 000 iterations were monitored and used to estimate the posterior quantities. The length of burn-in was determined by assessing convergence using the methods of Gelman and Rubin [28], Geweke [29] and Raftery and Lewis [30].

The classical random effects methods were implemented in Stata using the routine `metareg` [31, 32], which allows adjustment for covariates within a random effects meta-analysis. The regression model relates treatment effect to trial-level covariates, assuming a normal distribution for the residual errors with both a within-trial and an additive between-trial component of variance, τ^2 . The within-trial variances or standard errors must be supplied by the user, and are considered known. τ^2 is estimated either by an iterative procedure, using an estimate which

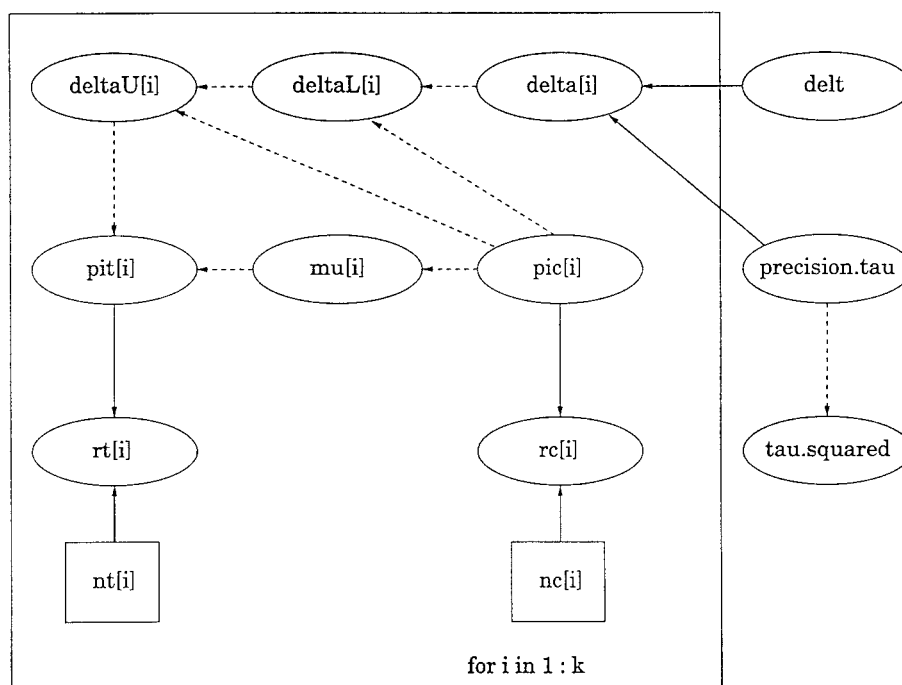


Figure 3. Directed acyclic graph (DAG) for the absolute risk scale showing the conditional independence assumptions of the model. The dashed arrows represent functional relationships and the full arrows stochastic relationships. The corresponding notation in the text is: $\text{delt} = \delta$, $\text{delta}[i] = \delta_i$, $\text{precision.tau} = 1/\tau^2$, $\text{tau.squared} = \tau^2$, $\mu[i] = \mu_i$, $\text{pic}[i] = \pi_i^C$, $\text{pit}[i] = \pi_i^T$.

is based on one of restricted maximum likelihood (REML), maximum likelihood or empirical Bayes methods, or by a non-iterative procedure using a method of moments estimator. If no covariate is specified, *metareg* carries out a random effects meta-analysis.

In this paper, we employ two classical methods of estimating τ^2 : method (d) uses restricted maximum likelihood (REML), an approach advocated by Thompson and Sharp [33]; method (e) uses the moment estimator, the most common approach in practice.

4. RESULTS

For each of the three treatment effect scales, five analyses were carried out (top of Table IV). Of the three Bayesian analyses, two used the binomial data directly, the other using summary estimates derived using the formulae in Table III. Two classical analyses using the summary estimates were performed for comparison. In Table IV, for the Bayesian methods (a), (b) and (c), posterior medians and standard deviations are given for δ and τ . These quantities provide the equivalent of the estimates and standard errors for the classical analysis methods (d) and (e). The 95 per cent intervals are the Bayesian credible intervals for methods (a), (b) and (c) which, unlike the 95 per cent confidence intervals for methods (d) and (e) calculated in the usual way as $\hat{\delta} \pm 1.96\text{SE}(\hat{\delta})$, do not rely on asymptotic standard errors.

Table IV. Summary of methods implemented and results from 31 trials of 400 mg ibuprofen on post-operative pain.

Method	Approach		Data		Note	
(a)	Bayesian		Binary		π_i^C treated as 'fixed effects'	
(b)	Bayesian		Binary		π_i^C treated as 'random effects'	
(c)	Bayesian		Summary		—	
(d)	Classical		Summary		τ^2 estimated using REML	
(e)	Classical		Summary		τ^2 estimated using method of moments	

	δ	(SD)	τ	(SD)	Treatment effect	(95 per cent CI)
RD					RD	
(a)	0.375	0.032	0.147	0.028	0.375	(0.312, 0.440)
(b)	0.393	0.031	0.145	0.028	0.393	(0.333, 0.456)
(c)	0.391	0.032	0.148	0.027	0.391	(0.328, 0.453)
(d)	0.391	0.031	0.140	—	0.391	(0.331, 0.451)
(e)	0.391	0.031	0.140	—	0.391	(0.331, 0.451)
log RR					RR	
(a)	1.336	0.154	0.674	0.146	3.864	(2.870, 5.263)
(b)	1.349	0.134	0.530	0.134	3.853	(3.045, 5.143)
(c)	1.235	0.164	0.584	0.156	3.438	(2.602, 4.933)
(d)	1.237	0.136	0.541	—	3.445	(2.635, 4.500)
(e)	1.219	0.131	0.509	—	3.384	(2.617, 4.380)
log OR					OR	
(a)	2.133	0.209	0.923	0.203	8.670	(5.570, 12.94)
(b)	2.210	0.197	0.822	0.201	9.340	(6.300, 13.82)
(c)	2.086	0.213	0.784	0.220	8.051	(5.481, 12.57)
(d)	2.095	0.192	0.764	—	8.125	(5.573, 11.86)
(e)	2.082	0.187	0.727	—	8.020	(5.557, 11.58)

The results from applying these analyses to the 31 trials of 400 mg ibuprofen are shown in Table IV. On the risk difference scale there is very little difference between the results from all the methods for both δ and τ . The posterior median for δ is smaller for method (a) than (b), a pattern repeated for the other treatment effect scales. There is no clear theoretical reason why this should occur and it may just be a feature of this data set. On the log relative risk and log-odds ratio scale both δ and τ behave similarly. Though the point estimates are fairly similar for all the methods, the treatment effect is larger for the methods (a) and (b) which model the binomial outcome data directly than for the summary statistic methods. This is probably due to the substantial proportion of trials with zero cell counts and the consequences that this entails for the calculation of the summary statistics. It is worth noting the difference between the values of the overall relative risk and odds ratio. The odds ratio is often interpreted as if it was the relative risk and the disparity between the two quantities for this data set not only demonstrates the dangers of this misinterpretation, but also the value of having a method for the relative risk scale itself.

The estimate of τ from method (a) is considerably larger, for the absolute risk difference and log-odds ratio scales, than from the other four methods. The difference between method

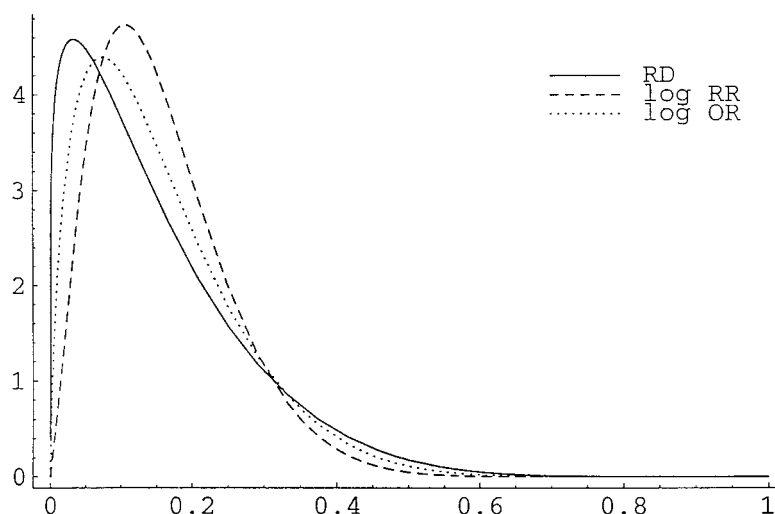


Figure 4. The estimated distribution of the π_i^C from method (b).

(a) and (b) is to be expected, since the imposition of a common distribution on the π_i^C shrinks the δ_i towards δ in such a way as to decrease the amount of variation around the common treatment effect. Methods (c), (d) and (e) assume that the within-trial variances are known, which also tends to decrease the estimated value of τ . The heterogeneity parameter τ on an absolute risk scale is much easier to interpret than on other effect scales. For example, the value of $\tau=0.15$ around an average risk difference of $\delta=0.40$ (Table IV) corresponds to a wide distribution of true absolute risk differences across trials, a 95 per cent range of around 0.1 to 0.7.

Figure 4 shows the estimated distribution of the π_i^C for each of the treatment scales using method (b). The curves are reassuringly similar.

5. EXTENSIONS

5.1. The use of informative prior distributions for τ

So far in this paper we have used a locally uniform prior for τ . The conjugate prior for the variance of a normal distribution is an inverse-gamma distribution, $IG(\alpha, \lambda)$. The standard Jeffreys's non-informative prior takes $p(\tau^2) \propto \frac{1}{\tau^2}$, corresponding to $IG(0, 0)$. This prior is not suitable for use at the second level of a hierarchical model, since it can lead to an improper posterior distribution [34]. Though progress has been made in identifying non-informative priors for τ^2 in normal-normal hierarchical models [35, 36], the issue of what represents a suitable reference prior in the binomial-normal case remains problematic.

However, given that a Bayesian framework seeks to encourage the application of prior knowledge, it makes sense to consider using more informative priors, especially if the analyst has knowledge of the particular medical context. Informative prior distributions may be formulated by 'introspection', which amounts to choosing a distribution that provides support to all reasonable parameter values. On the log-odds scale, Smith [15] proposed the prior $IG(3, 1)$

Table V. Results for methods (a) and (b) using the alternative prior $IG(0.001, 0.001)$ for τ^2 .

	δ	(SD)	τ	(SD)	Treatment effect	(95 per cent CI)
RD					RD	
(a)	0.376	0.032	0.141	0.027	0.376	(0.315, 0.438)
(b)	0.394	0.030	0.140	0.027	0.394	(0.333, 0.453)
log RR					RR	
(a)	1.315	0.149	0.642	0.141	3.724	(2.848, 5.115)
(b)	1.320	0.133	0.474	0.143	3.742	(2.969, 5.016)
log OR					OR	
(a)	2.124	0.202	0.854	0.194	8.362	(5.708, 12.61)
(b)	2.192	0.189	0.749	0.182	8.957	(6.255, 13.30)

for τ^2 after considering intervals likely to contain first the treatment difference and second the ratios of the two plausible extreme treatment differences. An informative prior distribution for τ^2 could also be obtained by empirical study. By collating the estimates of τ^2 from a large number of meta-analyses and examining the shape of the empirical cumulative distribution function, an informative data-based prior could be derived for each treatment effect scale [4, 37].

We can examine the sensitivity of the estimates of the parameters of interest to prior specification for τ . An alternative prior for τ that has been used in the past is an inverse-gamma $IG(0.001, 0.001)$ distribution on τ^2 . This distribution is a proper approximation to a flat prior for $\log(\tau)$. Consequently, it places more weight on values of τ close to zero than a prior which is uniform on τ itself. Results for methods (a) and (b) using this prior are given in Table V. The estimates of δ are very similar though, as we would expect, the values for τ are smaller for the inverse-gamma prior on all three treatment effect scales.

5.2. Relating treatment effect to trial-level covariates

The case for investigating potential sources of heterogeneity is well established [2, 5, 33, 38–41]. Heterogeneity of trial results is an indication of systematic differences between the trials included in a meta-analysis. These variations may be attributable to patient-level covariates or induced by trial-level characteristics. In either case, the analyst is interested in assessing the relationship between treatment efficacy and those factors giving rise to between-study variation. Indeed, the clinician must be informed if, and to what extent, net treatment benefit varies according to certain measurable characteristics [14]. These relationships can be explored either through the creation of relatively homogeneous subgroups of trials, or by regression analysis techniques such as meta-regression. In meta-regression, the aim is to detect any graded associations between the outcome of each clinical trial and some characteristic of the trial or its patients.

In the case of single-dose ibuprofen, we might be interested in whether the size of the dose x_i , a trial-level covariate, influences the size of the treatment effect. On the risk difference scale, we can investigate this by fitting a model in which

$$\delta_i = \delta_i^* + \gamma \left(\log(x_i) - \frac{1}{k} \sum_{i=1}^k \log(x_i) \right) \quad (10)$$

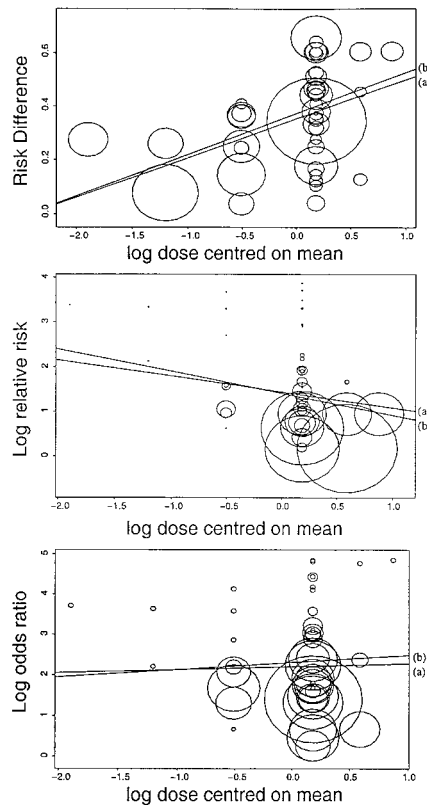


Figure 5. Plots of treatment effect against log dose on the three different effect scales. The size of each circle representing a trial is proportional to the precision of each trial's results (inverse variance of treatment effect estimate). The regression lines estimated by methods (a) and (b) are overlaid.

where $\delta_i^* \sim N(\delta, \tau^2)$. Here τ^2 measures the amount of residual heterogeneity not explained by the covariate. The use of $\log(x_i)$ rather than x_i is arbitrary, but the transformation of scale is motivated by the fact that the values of x_i (apart from 600 mg) increase multiplicatively. Given that we only have six different doses, an alternative way to analyse the data would be consider dose as a categorical variable. We also centre the covariate around its mean since this improves computational efficiency. Figure 5 shows plots of the observed treatment effect against $\log(x_i)$ for each of the three treatment effect scales.

A key advantage of the fully Bayesian approach is the ease with which one can adjust for covariates. Since the regression slope γ is the parameter of primary interest it is appropriate to place a vague prior on it, such as $\gamma \sim N(0, 10\,000)$. The minor changes in the BUGS code necessary for implementing (10) are given in the Appendix. As before, a burn-in of 5000 iterations was used, and the the posterior distribution of the parameters derived from the following 15 000 iterations.

Table VI. Relating treatment effect to trial-level covariates.

		Log dose			Underlying risk		
		Slope	(SD)	95 per cent CI	Slope	(SD)	95 per cent CI
RD	(a)	0.146	(0.050)	(0.033, 0.237)	-0.217	(0.212)	(-0.635, 0.186)
	(b)	0.153	(0.048)	(0.059, 0.246)	-0.135	(0.251)	(-0.645, 0.374)
log RR	(a)	-0.494	(0.294)	(-1.031, 0.082)	-0.817	(0.059)	(-0.931, -0.707)
	(b)	-0.357	(0.316)	(-1.028, 0.230)	-0.778	(0.060)	(-0.883, -0.646)
log OR	(a)	0.075	(0.367)	(-0.645, 0.813)	-0.629	(0.118)	(-0.850, -0.389)
	(b)	0.180	(0.380)	(-0.553, 0.955)	-0.596	(0.132)	(-0.826, -0.306)

It can be seen from Table VI that method (a) yields a more negative slope than method (b), though the difference is not sufficiently substantial to change the interpretation. There is evidence of a significant positive link between risk difference and dose, so that the absolute risk benefit of ibuprofen increases with dose (Figure 5). A negative (although not significant) relationship yields the opposite conclusion for relative risks, while there is almost no relationship with dose on an odds ratio scale (Figure 5). One reason for this is that the trials that contribute the most weight to the analysis on one scale are not always those that contribute most on another scale, as shown by the sizes of the circles in Figure 5. This is because the same binary outcome data can, for example, give rise to a precise estimate of the absolute risk difference and an imprecise estimate of the log relative risk, or indeed vice versa.

5.3. Relating treatment effect to underlying risk

The *underlying risk* of the event occurring, for all the patients in each trial, has been widely discussed as a potential source of between-study heterogeneity. Also known as the population [42] or baseline [43] risk, underlying risk represents a convenient aggregate measure of the health of a patient population in a trial. It provides a trial-level summary of a number of patient level attributes, or risk factors, which though measurable may not be available as individual patient data for some or all of the trials to be included in the meta-analysis. The existence of a relationship between treatment benefit and underlying risk has profound implications for the interpretation of the results of a meta-analysis both in the direct evaluation of the net treatment benefit, and with regard to broader health economic considerations associated with the appropriate use of treatment [5, 44, 45]. In particular, the nature of any relationship will help identify which patients will have most, or indeed least, to gain from a medical intervention [14].

With the increasing popularity of meta-analyses attempting to investigate this relationship in a range of medical areas, a strong impetus exists behind the development of statistically valid approaches for obtaining unbiased estimates of the true relationship between treatment effect and underlying risk. Naive analyses are biased by regression to the mean [14], which can be circumvented by Bayesian analyses for both binomial outcome on the log-odds scale [5] and summary statistic data [42] on all scales. We can extend the methods presented in this paper for the risk difference and log relative risk scales for binomial outcome data to allow analysis

on these scales. This is achieved by replacing the previous assumption $\delta_i \sim N(\delta, \tau^2)$ by the following:

$$\begin{aligned}\delta_i &= \delta_i^* + \gamma(\mu_i - \bar{\mu}) \\ \delta_i^* &\sim N(\delta, \tau^2)\end{aligned}\tag{11}$$

Here $\bar{\mu}$ denotes the average of the μ_i across trials and δ_i^* is the treatment effect in trial i adjusted for underlying risk. A prior is required for γ , and is once again taken as non-informative, $N(0, 10000)$.

In the underlying risk analyses, we consider just the 31 trials with dose 400mg (Table VI). As with log dose, the slope for method (a) is more negative slope than for method (b). For underlying risk, this is a typical regression dilution effect. The slope for underlying risk is negative on all scales and substantially and significantly so for the log relative risk and the log-odds ratio. This implies that as the underlying chance of a 50 per cent pain reduction increases, the size of the treatment effect decreases. It should be noted that the choice of treatment effect scale clearly affects the strength of the relationship, an area open for empirical investigation [37].

6. DISCUSSION

This paper extends an existing method for performing a Bayesian random effects meta-analysis of trials with binary outcomes to the risk difference and relative risk scales. The methods presented allow the investigation of sources of heterogeneity by considering trial-level covariates including underlying risk. McIntosh [42] proposed a general approach using summary statistics for modelling risk differences, (log) relative risks and (log) odds ratios in the context of underlying risk meta-regression. This method was used by Schmid *et al.* [46] in a subsequent empirical study into the prevalence of a relationship between underlying risk and treatment effect on different treatment effect scales. It is now possible to repeat this empirical study with a method that uses the binomial outcome data directly.

Classical approaches to meta-analysis ignore the fact that τ^2 has been estimated from the data. One of the advantages of a fully Bayesian approach is that it allows for uncertainty in all the parameters of the model. Furthermore, we can obtain credible intervals for these parameters from the posterior distribution. Though asymptotic classical confidence intervals have been derived for τ^2 [47], the comparative ease with which a credible interval may be obtained from a fully Bayesian analysis is attractive. Alternatively, we might be interested in predicting the treatment effect in a new trial. Obtaining the predictive distribution for this is a straightforward procedure [1]. Figure 6 shows a contour plot of the joint predictive distribution for the true group risks in a new trial produced by simulation under model (b). Binomial error could be added to this to produce predictions of observed responses. Moore *et al.* [48] produced a similar plot for the same data based on the results of 10 000 simulated studies with the purpose of demonstrating where a study might occur at random on a L'Abbé plot. However, they did not allow for between-trial variation in the true effect.

Models including trial-level covariates are ecological models and can describe only between-trial, not between-patient, variation in characteristics. Consequently, they are most useful when studying characteristics that differ across studies. This presents problems for meta-

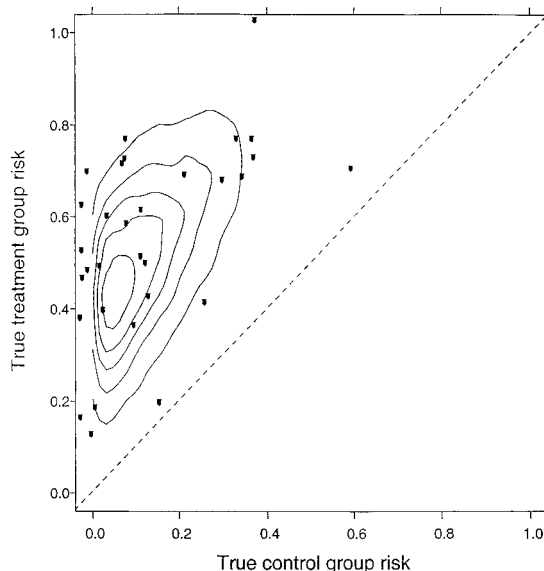


Figure 6. Contour plot of predictive distribution for a new trial superimposed on a L'Abbé plot of the observed data from the 31 trials.

regression since averages of patient characteristics which exhibit little between-study variation will provide only minimal information across the potential range of the factor [46]. A further complication lies in what Morgenstern [49] and Greenland [50] call aggregation or ecologic bias. This bias arises from the failure to account for within-trial variation through the use of aggregate values, and exists if the relation between group risks and means does not resemble the relation between individual outcomes and covariates. For example, since the underlying risk of all the patients in a trial is only an ecologic substitute for the actual risk of each individual, a regression may not truly reflect a genuine treatment effect interaction with underlying risk at the level of the individual.

The specification of prior distributions for parameters is one of the most difficult, and controversial, aspects of Bayesian inference. In this paper, we have attempted to use non-informative prior distributions. In the case of δ , this has been achieved by the use of a vague prior providing approximately uniform support over the range of the likelihood. Of course, it would be perfectly reasonable to use an alternative prior distribution for the effect parameter, for example, the 'sceptical' and 'enthusiastic' priors discussed by Spiegelhalter *et al.* [51]. For τ^2 we have placed a uniform prior on τ . Other possibilities include the use of priors formulated by considering the range in which treatment effects might reasonably lie, data-based prior distributions drawn from empirical studies, or subjective priors based on clinical judgement.

There is substantial value in modelling the binomial structure of the data directly. As this paper has shown, a fully Bayesian analysis avoids the need for simplifying approximations and provides the analyst with a powerful and flexible modelling framework.

APPENDIX: BUGS CODE

The BUGS software is freely available from <http://www.mrc-bsu.cam.ac.uk>. In BUGS, the normal distribution is parameterized in terms of the mean and the precision, that is, $1/\text{variance}$.

A1. Absolute risk difference scale

The code for method (a) is as follows:

```
model
{
  for (i in 1:k)
  {
    rc[i] ~ dbin(pic[i],nc[i])
    rt[i] ~ dbin(pit[i],nt[i])
    mu[i] <- pic[i]
    pit[i] <- mu[i] + min(max(delta[i],-pic[i]),(1-pic[i]))
    delta[i] ~ dnorm(delt,precision.tau)
    pic[i] ~ dunif(0,1)
  }
  delt ~ dunif(-1,1)
  precision.tau <- 1/tau.squared
  tau.squared <- tau*tau
  tau ~ dunif(0,2)
}
```

For method (b), replace

```
pic[i] ~ dunif(0,1)

with

pic[i] ~ dbeta(alpha,beta)
alpha ~ dunif(1,100)
beta ~ dunif(1,100)
```

To allow non-integer values in the binomial likelihood replace

```
rt[i] ~ dbin(pit[i],nt[i])
rc[i] ~ dbin(pic[i],nc[i])

with

onest[i] <- 1
onesc[i] <- 1
onest[i] ~ dbern(wt[i])
onesc[i] ~ dbern(wc[i])
wt[i] <- pow(pit[i],rt[i]) * pow((1-pit[i]), (nt[i]-rt[i]))
wc[i] <- pow(pic[i],rc[i]) * pow((1-pic[i]), (nc[i]-rc[i]))
```

A2. Relative risk scale

```

model
{
for (i in 1:k)
{
  rc[i] ~ dbin(pic[i],nc[i])
  rt[i] ~ dbin(pit[i],nt[i])
  mu[i] <- log(pic[i])
  log(pit[i]) <- mu[i] + min(delta[i],-log(pic[i]))
  delta[i] ~ dnorm(delt,precision.tau)
  pic[i] ~ dunif(0,1)
}
delt ~ dnorm(0,0.1)
precision.tau <- 1/tau.squared
tau.squared <- tau*tau
tau ~ dunif(0,2)
}

```

In order to make the code work in BUGS, it was necessary to alter

```
deltaU[i] <- min(delta[i],-log(pc[i]))
```

to ensure that $\log(\pi_i^T)$ is strictly less than 0. We achieve this by changing the upper limit of δ_i^U to $\frac{-\log(\pi_i^C)}{C}$, where $C = 1.0000001$. The line becomes

```
deltaU[i] <- min(delta[i],-log(pc[i])/C)
C <- 1.0000001
```

A3. Odds ratio scale

```

model
{
for (i in 1:k)
{
  rc[i] ~ dbin(pic[i],nc[i])
  rt[i] ~ dbin(pit[i],nt[i])
  mu[i] <- logit(pic[i])
  logit(pit[i]) <- mu[i] + delta[i]
  delta[i] ~ dnorm(delt,precision.tau)
  pic[i] ~ dunif(0,1)
}
delt ~ dnorm(0,0.1)
precision.tau <- 1/tau.squared
tau.squared <- tau*tau
tau ~ dunif(0,2)
}

```

A4. Bayesian summary statistic method

The code for method (c) is as follows:

```
model
{
  for(i in 1:k)
  {
    d[i] ~ dnorm(delta[i],precision.d[i])
    precision.d[i] <- 1/variance.d[i]
    delta[i] ~ dnorm(delt,precision.tau)
  }
  precision.tau <- 1/tau.squared
  tau.squared <- tau*tau
  tau ~ dunif(0,2)
  delt ~ dnorm(0,0.1)
}
```

For the risk difference scale use

```
delt ~ dunif(-1,1)
```

*A5. Introducing covariates and underlying risk**A5.1. Log dose*

Add the line

```
logdose[i]<-log(x[i])
```

Then change

```
delta[i] ~ dnorm(delt,precision.tau)
```

to

```
delta.star[i] ~ dnorm(delt,precision.tau)
delta[i] <- delta.star[i]+gamma*(logdose[i]-mean(logdose[]))
```

and add the prior

```
gamma ~ dnorm(0,0.0001)
```

A5.2. Underlying risk

Replace the line

```
delta[i] <- delta.star[i]+gamma*(logdose[i]-mean(logdose[]))
```

with

```
delta[i] <- delta.star[i]+gamma*(mu[i]-mean(mu[]))
```

ACKNOWLEDGEMENTS

The authors would like to thank Julian Higgins for his helpful comments at all stages of this work.

REFERENCES

1. Smith TC, Spiegelhalter DJ, Thomas A. Bayesian approaches to random-effects meta-analysis: a comparative study. *Statistics in Medicine* 1995; **14**:2685–2699.
2. Schmid CH. Exploring heterogeneity in randomized trials via meta-analysis. *Drug Information Journal* 1999; **33**:211–224.
3. Sutton AJ, Abrams KR, Jones DR, Sheldon TA, Song F. Systematic reviews of trials and other studies. *Health Technology Assessment* 1998; **2** report 19.
4. Higgins JPT, Whitehead A. Borrowing strength from external trials in a meta-analysis. *Statistics in Medicine* 1996; **15**:2733–2749.
5. Thompson SG, Smith TC, Sharp SJ. Investigating underlying risk as a source of heterogeneity in meta-analysis. *Statistics in Medicine* 1997; **16**:2741–2758.
6. Carlin JB. Letter to the Editor. Meta-analysis: formulating, evaluating, combining, and reporting. *Statistics in Medicine* 2000; **19**:753–761.
7. Deeks JJ, Altman DG. Effect measures for meta-analysis of trials with binary outcomes. In *Systematic Reviews in Health Care: Meta-analysis in Context*, Egger M, Davey Smith G, Altman DG (eds). BMJ Books: London, 2001; 313–336.
8. Maynard A. Evidence-based medicine: an incomplete method for informing treatment choices. *The Lancet* 1997; **349**:126–128.
9. Cook RJ, Sackett DL. The number needed to treat: a clinically useful measure of treatment effect. *British Medical Journal* 1995; **310**:452–454.
10. Altman DG. Confidence intervals for the number needed to treat. *British Medical Journal* 1998; **317**: 1309–1312.
11. Hutton JL. Number needed to treat: properties and problems. *Journal of the Royal Statistical Society, Series A* 2000; **163**:893–906.
12. Collins SL, Moore RA, McQuay HJ, Wiffen PJ, Edwards JE. Single dose oral ibuprofen and diclofenac for post-operative pain (Cochrane Review). The Cochrane Library, Oxford Update Software, 2, 2000.
13. L'Abbé KA, Detsky AS, O'Rourke K. Meta-analysis in clinical research. *Annals of Internal Medicine* 1987; **107**:224–233.
14. Sharp SJ, Thompson SG, Altman DG. The relation between treatment benefit and underlying risk in meta-analysis. *British Medical Journal* 1996; **313**:735–738.
15. Smith TC. Interpreting evidence from multiple randomised and non-randomised studies. *PhD thesis*, University of Cambridge, 1995.
16. van Houwelingen HC, Senn SJ. Letter to the Editor. Investigating underlying risk as a source of heterogeneity in meta-analysis. *Statistics in Medicine* 1999; **18**:110–113.
17. Thompson SG, Prevost TC, Sharp SJ. Author's Reply to Letter to the Editor. Investigating underlying risk as a source of heterogeneity in meta-analysis. *Statistics in Medicine* 1999; **18**:113–115.
18. Marshall EC, Spiegelhalter DJ. Comparing institutional performance using Markov chain Monte Carlo methods. In *Statistical Analysis of Medical Data: New Developments*, Everitt BS, Dunn G (eds). Arnold: London, 1999; 229–249.
19. van Houwelingen HC, Zwinderman KH, Stijnen T. A bivariate approach to meta-analysis. *Statistics in Medicine* 1993; **12**:2273–2284.
20. Carlin JB. Meta-analysis for 2×2 tables: a Bayesian approach. *Statistics in Medicine* 1992; **11**:141–158.
21. Gelmen A, Carlin JB, Stern HS, Rubin DB. *Bayesian Data Analysis*. Chapman and Hall: London, 1995.
22. Fleiss JL. The statistical basis of meta-analysis. *Statistical Methods in Medical Research* 1993; **2**:121–145.
23. Whitehead A, Whitehead J. A general parametric approach to meta-analysis of randomized clinical trials. *Statistics in Medicine* 1990; **10**:1665–1677.
24. Cox DR, Snell EJ. *Analysis of Binary Data*. Chapman and Hall: London, 1989.
25. DerSimonian R, Laird N. Meta-analysis in clinical trials. *Controlled Clinical Trials* 1986; **7**:177–188.
26. Gilks WR, Thomas A, Spiegelhalter DJ. A language and program for complex Bayesian modelling. *Statistician* 1994; **43**:169–177.
27. Spiegelhalter DJ, Thomas A, Best N, Gilks WR. *BUGS 0.5 Bayesian inference using Gibbs Sampling Manual (version ii)*, 1996.
28. Gelman A, Rubin DB. A single series from the Gibbs sampler provides a false sense of security. In *Bayesian Statistics 4*, Bernardo JM, Berger JO, Dawid AP, Smith AFM (eds). Oxford University Press: Oxford, 1992; 625–631.

29. Geweke J. Evaluating the accuracy of sampling-based approaches to calculating posterior moments. In *Bayesian Statistics 4*, Bernardo JM, Berger JO, Dawid AP, Smith AFM (eds). Oxford University Press: Oxford, 1992.
30. Raftery AE, Lewis SM. How many iterations in the Gibbs sampler? In *Bayesian Statistics 4*, Bernardo JM, Berger JO, Dawid AP, Smith AFM (eds). Oxford University Press: Oxford 1992; 763–775.
31. Sharp SJ. Meta-analysis regression. *Stata Technical Bulletin* 1998; **42**:16–22.
32. Sterne JAC, Bradburn MJ, Egger M. Meta-analysis in STATA. In *Systematic Reviews in Health Care: Meta-Analysis in Context*, Egger M, Davey Smith G, Altman DG (eds). BMJ Books: London, 2001; 347–372.
33. Thompson SG, Sharp SJ. Explaining heterogeneity in meta-analysis: a comparison of methods. *Statistics in Medicine* 1999; **18**:2693–2708.
34. DuMouchel W, Waternaux C. Comment on ‘Hierarchical models for combining information and for meta-analyses’. In *Bayesian Statistics 4*, Bernardo JM, Berger JO, Dawid AP, Smith AFM (eds). Oxford University Press: Oxford, 1992; 338–339.
35. Natarajan R, Kass RE. Reference Bayesian methods for generalized linear mixed models. *Journal of the American Statistical Association* 2000; **95**:227–237.
36. Daniels MJ. A prior for the variance in hierarchical models. *Canadian Journal of Statistics* 1999; **27**:569–580.
37. Engels EA, Schmid CH, Terrin N, Olkin I, Lau J. Heterogeneity and statistical significance in meta-analysis: an empirical study of 125 meta-analyses. *Statistics in Medicine* 2000; **19**:1707–1728.
38. Thompson SG. Why sources of heterogeneity in meta-analyses should be investigated. *British Medical Journal* 1994; **309**:1351–1355.
39. Lau J, Ioannidis JPA, Schmid CH. Summing up evidence: one answer is not always enough. *Lancet* 1998; **351**:123–127.
40. Berkey CS, Hoaglin DC, Mosteller F, Colditz GA. A random effects regression model for meta-analysis. *Statistics in Medicine* 1995; **14**:395–411.
41. Davey Smith G, Song F, Sheldon TA. Cholesterol lowering and mortality: the importance of considering the initial level of risk. *British Medical Journal* 1993; **306**:1367–1372.
42. McIntosh MW. The population risk as an explanatory variable in research synthesis of clinical trials. *Statistics in Medicine* 1996; **15**:1713–1728.
43. Walter SD. Variation in baseline risk as an explanation of heterogeneity in meta-analysis. *Statistics in Medicine* 1997; **16**:2883–2900.
44. Brand R, Kragt H. Importance of trends in the interpretation of an overall odds ratio in the meta-analysis of clinical trials. *Statistics in Medicine* 1992; **11**:2077–2082.
45. Davey Smith G, Egger M. Who benefits most from medical interventions? *British Medical Journal* 1994; **308**:72–74.
46. Schmid CH, Lau J, McIntosh MW, Cappelleri JC. An empirical study of the effect of the control rate as a predictor of treatment efficacy in meta-analysis of clinical trials. *Statistics in Medicine* 1998; **17**:1923–1942.
47. Biggerstaff BJ, Tweedie RL. Incorporating variability in estimates of heterogeneity in the random effects model in meta-analysis. *Statistics in Medicine* 1997; **16**:753–768.
48. Moore RA, Gavaghan D, Tramèr MR, Collins SL, McQuay HJ. Size is everything—large amounts of information are needed to overcome random effects in estimating direction and magnitude of treatment effects. *Pain* 1998; **78**:209–216.
49. Morgenstern H. Uses of ecological analysis in epidemiological research. *American Journal of Public Health* 1982; **72**(12):1336–1344.
50. Greenland S. Quantitative methods in the review of epidemiological literature. *Epidemiological Review* 1987; **9**:1–30.
51. Spiegelhalter DJ, Freedman LS, Parmar MKB. Bayesian approaches to randomized trials. *Journal of the Royal Statistical Society, Series A* 1994; **157**:357–387.