

A GENERAL PARAMETRIC APPROACH TO THE META-ANALYSIS OF RANDOMIZED CLINICAL TRIALS

ANNE WHITEHEAD AND JOHN WHITEHEAD

Department of Applied Statistics, University of Reading, Whiteknights, P.O. Box 217, Reading RG6 2AN, U.K.

SUMMARY

Meta-analysis provides a systematic and quantitative approach to the summary of results from randomized studies. Whilst many authors have published actual meta-analyses concerning specific therapeutic questions, less has been published about comprehensive methodology. This article presents a general parametric approach, which utilizes efficient score statistics and Fisher's information, and relates this to different methods suggested by previous authors. Normally distributed, binary, ordinal and survival data are considered. Both the fixed effects and random effects model for treatments are described.

INTRODUCTION

The results from a collection of independent randomized studies can be summarized in a systematic and quantitative way using a meta-analysis. The main objective of such an analysis is to obtain information about treatment effects that cannot be ascertained from any of the studies alone. Any individual study may either be too small to detect moderate treatment effects, say on mortality, or too limited to allow generalization to other patient populations. We should like to know, overall, whether a treatment has a beneficial or harmful effect. Reviews of treatments or therapeutic areas are frequently carried out when new compounds are developed: they could perhaps benefit from such a quantitative approach in addition to qualitative and subjective summaries.

A meta-analysis can be viewed as an extreme form of multi-centre study. There is a continuum from the true multi-centre study, in which all centres follow an identical protocol, to a collection of studies addressing the same general therapeutic question but with different protocols, different treatments, and different primary response variables. The methods discussed in this paper may be applied with differing levels of validity, across this continuum.

The hardest part of combining results from different studies is deciding which to include. Selection of studies has been discussed by various authors including Yusuf¹ and Light.² It is sensible to consider carefully the quality of each study before deciding on inclusion. A good example of this procedure has been given by Goodwin and Boyd³ in an appraisal of epidemiological studies investigating the association between breast cancer and mammographic parenchymal pattern. Each study was given a point for meeting each of seven methodological standards including avoidance of bias, adjustment for potential confounders and the handling of follow-up. Goodwin and Boyd noted that there was an overall consistency in the relative risk estimate from studies meeting all the standards, whereas heterogeneity was introduced when the remaining studies were included. Goodwin and Boyd's meta-analysis is an example of an

overview of non-randomized studies, demonstrating the utility of the method beyond the context of clinical trials. We shall not consider further the choice of studies for inclusion in a meta-analysis: instead we shall concentrate on the methodology of the analysis.

The objective of this paper is to present a general parametric approach to estimation and hypothesis testing in meta-analyses, and to the identification and treatment of heterogeneity. First a model in which treatment effects are fixed is considered, and later a random effects model allowing treatment effects to vary between studies. The general approach utilizes efficient score statistics and Fisher's information. Normally distributed, binary and ordinal responses will be considered, as well as survival data. The relationship with different methods suggested by other authors^{4, 5} will be clarified. We illustrate the methods in the binary case using results from a series of studies concerning the occurrence of stroke in hypertensive patients.⁶

A GENERAL FIXED EFFECTS PARAMETRIC APPROACH

Suppose that we have k studies each comparing an experimental treatment (E) with a control (C). The response variable of interest in each study is the same. Let θ_i denote the value of the chosen measure of treatment effect, that is the effect of the experimental treatment relative to the control, in the i th study. This may, for example, be the log-odds-ratio for binary data or the difference between treatment means for normally distributed data. Denote by $\hat{\theta}_i$, an estimate of θ_i from the i th study and let w_i be the inverse of the asymptotic variance of $\hat{\theta}_i$. If we can assume that asymptotically

$$\hat{\theta}_i \sim N(\theta_i, w_i^{-1}), \quad \text{for } i = 1, \dots, k, \quad (1)$$

then $\hat{\theta}_i w_i \sim N(\theta_i w_i, w_i)$ and under the null hypothesis $H_{0i}: \theta_i = 0$, $\hat{\theta}_i w_i \sim N(0, w_i)$. Furthermore under the combined null hypothesis $H_0: \theta_1 = \dots = \theta_k = 0$, $\sum \hat{\theta}_i w_i \sim N(0, \sum w_i)$ and so the statistic $U = (\sum \hat{\theta}_i w_i)^2 / \sum w_i$ follows a χ^2 distribution with 1 degree of freedom. In the meta-analysis, U can be used to test H_0 . More generally, assuming homogeneity of treatment effects over all studies, that is $\theta_1 = \dots = \theta_k = \theta$, $\sum \hat{\theta}_i w_i \sim N(\theta \sum w_i, \sum w_i)$ and we can estimate θ by $\hat{\theta}$, where $\hat{\theta} = \sum \hat{\theta}_i w_i / \sum w_i$. An approximate 95 per cent confidence interval for θ is given by $\hat{\theta} \pm 1.96 \sqrt{1 / \sum w_i}$.

To test the homogeneity of the treatment effect across all studies we use a large sample test based on the statistic Q where

$$Q = \sum w_i (\hat{\theta}_i - \hat{\theta})^2, \quad (2)$$

which is a weighted sum of squares of deviations. When treatment effects are homogeneous, Q follows a χ^2 distribution with $(k - 1)$ degrees of freedom.⁷

A general method for choosing the estimate $\hat{\theta}_i$ of θ_i can be given. If Z_i is the efficient score for θ_i and V_i is Fisher's information, both evaluated for $\theta = 0$, then asymptotically $Z_i \sim N(\theta_i V_i, V_i)$.⁸ The choice $\hat{\theta}_i = Z_i / V_i$ is an approximate maximum likelihood estimate. As (approximately) $Z_i / V_i \sim N(\theta_i, V_i^{-1})$, w_i will be taken to equal V_i .

To be precise, let $l_i(\theta, \phi)$ denote the log-likelihood of θ and ϕ based on the data from the i th study, where ϕ is a (possibly vector-valued) nuisance parameter. If $\hat{\phi}(\theta)$ denotes the maximum likelihood estimate of ϕ conditional on θ , then $l_i(\theta, \hat{\phi}(\theta))$ is a profile likelihood. The statistic Z_i is the first derivative of $l_i(\theta, \hat{\phi}(\theta))$ with respect to θ , evaluated at $\theta = 0$; V_i is minus the second derivative of $l_i(\theta, \hat{\phi}(\theta))$ with respect to θ , evaluated at $\theta = 0$, which in turn is equal to minus the inverse of the leading term of the inverse of the matrix of second derivatives of $l_i(\theta, \phi)$, evaluated at $\theta = 0$ and $\phi = \hat{\phi}(0)$. For details see Chapter 3 of Reference 8.

Table I. Binary data from the i th study

ith study	Treatment		Overall
	E	C	
Model: $p(\text{success})$	p_{Ei}	p_{Ci}	p_i
Data: success	s_{Ei}	s_{Ci}	s_i
failure	f_{Ei}	f_{Ci}	f_i
total	n_{Ei}	n_{Ci}	n_i

Whether the estimate $\hat{\theta}_i = Z_i/V_i$ is chosen or not, the important point is that the measure of treatment effect and its estimate should fulfil equation (1) with reasonable accuracy. Optimal parameterizations for normality are discussed by Anscombe⁹ and Spratt.¹⁰

In the remainder of this section the estimate based on the efficient score and Fisher's information is derived for binary, ordinal, normally distributed and survival data. Relationships with different published methods, most of which depend on the asymptotic normality of the estimate are discussed.

Binary data

Two commonly used measures of treatment effect are the probability difference and the log-odds-ratio. The log-odds-ratio has the advantage that all finite values can be interpreted whereas the probability difference is constrained to lie between -1 and $+1$: the normal approximation (1) is more accurate for the former. Unless there is no treatment effect at all, homogeneity of treatment effect across all studies in one scale implies heterogeneity in the other. Heterogeneity in the probability difference scale is likely to arise if the control rates take a wide range of values or if all the rates are close to zero or one.

The two possible outcomes of a binary response will be referred to as 'success' and 'failure'. For each study we have the information shown in Table I.

Consider first the log-odds-ratio:

$$\theta_i = \log \left\{ \frac{p_{Ei}(1 - p_{Ci})}{p_{Ci}(1 - p_{Ei})} \right\}.$$

Treating s_i and f_i as fixed, the efficient score for θ_i is

$$Z_i = s_{Ei} - n_{Ei}s_i/n_i$$

and Fisher's information is

$$V_i = n_{Ei}n_{Ci}s_i f_i / n_i^2 (n_i - 1).$$

Notice that Z_i can be expressed as $O_i - E_i$ where O_i and E_i are the observed and expected number of successes on the experimental treatment. Notice also that if $V'_i = (n_i - 1)V_i/n_i$, then Z_i^2/V_i is the familiar Pearson χ^2 -statistic for the analysis of a 2×2 contingency table. The parameter θ_i can be estimated by $\hat{\theta}_i = Z_i/V_i$. Now the statistics U and Q have already been defined as

$$U = (\sum \hat{\theta}_i w_i)^2 / \sum w_i \quad \text{and} \quad Q = \sum w_i (\hat{\theta}_i - \hat{\theta})^2$$

Table II. Meta-analysis of antihypertensive studies using the log-odds-ratio approach

Study	Treatment		Control		$\hat{\theta}_i$	w_i	$\hat{\theta}_i w_i$	$\hat{\theta}_i^2 w_i$
	Number of strokes	Total number	Number of strokes	Total number				
VA-NHLBI	0	508	0	504	—	—	—	—
HDFP (Stratum I)	59	3903	88	3922	-0.40	36.1	-14.3	5.69
Oslo	0	406	5	379	-2.08	1.2	-2.6	5.38
ANBPS	13	1721	22	1706	-0.53	8.7	-4.6	2.42
MRC	60	8700	109	8654	-0.59	41.8	-24.7	14.61
VAII	5	186	20	194	-1.24	5.9	-7.2	8.95
USPHS	1	193	6	196	-1.44	1.7	-2.5	3.55
HDFP (Stratum II)	25	1048	36	1004	-0.42	14.8	-6.2	2.56
HSCSG	43	233	52	219	-0.32	18.8	-6.0	1.90
VAI	1	68	3	63	-1.10	1.0	-1.1	1.19
WOLFF	2	45	1	42	0.61	0.7	0.4	0.27
Barracough	0	58	0	58	—	—	—	—
Carter	10	49	21	48	-1.06	5.3	-5.7	6.01
HDFP (Stratum III)	18	534	34	529	-0.66	12.4	-8.1	5.33
EWPHS	32	416	48	424	-0.42	18.1	-7.6	3.20
Coope	20	419	39	465	-0.58	13.7	-8.0	4.62
Total						180.2	-98.0	65.68

Test for heterogeneity

$$Q = 12.4 \quad \chi^2_{13} (5 \text{ per cent}) = 22.36$$

Test for treatment effect

$$U = 53.3 \quad \chi^2_1 (5 \text{ per cent}) = 3.84$$

$$\hat{\theta} = -0.544 \quad 95 \text{ per cent CI} = (-0.690, -0.398)$$

where $\hat{\theta} = \sum \hat{\theta}_i w_i / \sum w_i$. For the choice $\hat{\theta}_i = Z_i / V_i$ it follows that $w_i = V_i$. Thus $\hat{\theta} = \sum Z_i / \sum V_i$, $U = (\sum Z_i)^2 / \sum V_i$ and

$$Q = \sum V_i \left(\frac{Z_i}{V_i} - \frac{\sum Z_i}{\sum V_i} \right)^2$$

$$= \sum \frac{Z_i^2}{V_i} - \frac{(\sum Z_i)^2}{\sum V_i}.$$

These statistics can be seen to be the same as those quoted in the Statistical Appendix of Yusuf *et al.*⁴ DerSimonian and Laird⁵ explore the same problem, and for the log-odds-ratio choose the maximum likelihood estimate $\hat{\theta}_i = \log \{s_{Ei} f_{Ci} / (s_{Ci} f_{Ei})\}$, which has variance w_i^{-1} where $w_i^{-1} = (s_{Ei}^{-1} + s_{Ci}^{-1} + f_{Ei}^{-1} + f_{Ci}^{-1})$.

In the probability difference approach, $\theta_i = p_{Ei} - p_{Ci}$. Choosing the maximum likelihood estimate $\hat{\theta}_i = (s_{Ei}/n_{Ei}) - (s_{Ci}/n_{Ci})$ gives $w_i^{-1} = (s_{Ei} f_{Ei} n_{Ei}^{-3} + s_{Ci} f_{Ci} n_{Ci}^{-3})$ as discussed by DerSimonian and Laird⁵ and Berlin *et al.*¹¹ This does not follow from the efficient score and Fisher's information.

As an illustration of the general parametric method for binary data consider the results from 16 randomized trials of antihypertensive drugs presented by Collins *et al.*⁶ Table II shows a meta-analysis for the occurrence of strokes based on the log-odds-ratio approach using the efficient score and Fisher's information. For comparison, a meta-analysis based on the probability difference approach mentioned above is shown in Table III. Perversely, occurrence of a stroke is

Table III. Meta-analysis of antihypertensive studies using the probability difference approach

Study	Treatment		Control		$\hat{\theta}_i$	w_i	$\hat{\theta}_i w_i$	$\hat{\theta}_i^2 w_i$
	Number of strokes	Total number	Number of strokes	Total number				
VA-NHLBI	0	508	0	504	—	—	—	—
HDFP (Stratum I)	59	3903	88	3922	-0.0073	106303	-778	5.70
Oslo	0	406	5	379	-0.0132	29112	-384	5.07
ANBPS	13	1721	22	1706	-0.0053	84620	-452	2.41
MRC	60	8700	109	8654	-0.0057	449571	-2562	14.60
VAII	5	186	20	194	-0.0762	1620	-123	9.41
USPHS	1	193	6	196	-0.0254	5614	-143	3.63
HDFP (Stratum II)	25	1048	36	1004	-0.0120	17651	-212	2.54
HSCSG	43	233	52	219	-0.0529	679	-36	1.90
VAI	1	68	3	63	-0.0329	1072	-35	1.16
WOLFF	2	45	1	42	0.0206	668	14	0.28
Barracough	0	58	0	58	—	—	—	—
Carter	10	49	21	48	-0.2334	118	-28	6.45
HDFP (Stratum III)	18	534	34	529	-0.0306	5725	-175	5.35
EWPHE	32	416	48	424	-0.0363	2454	-89	3.23
Coope	20	419	39	465	-0.0361	3653	-132	4.77
Total						708861	-5136	66.51

Test for heterogeneity

$$Q = 29.3 \quad \chi^2_{13} (5 \text{ per cent}) = 22.36$$

Test for treatment effect

$$U = 37.2 \quad \chi^2_1 (5 \text{ per cent}) = 3.84$$

$$\hat{\theta} = -0.0072 \quad 95 \text{ per cent CI} = (-0.0095, -0.0049)$$

taken to be a 'success', thus a negative θ_i implies that fewer strokes occur on treatment than on control. In two of the constituent studies no strokes were observed in either treatment or control and consequently they provide no information on the treatment effect: the meta-analysis is performed on the remaining 14. The percentage of patients in the control group, who had a stroke, varies from 1.3 to 43.8 across the 14 studies and this goes some way to explain why the test for heterogeneity is significant for the probability difference estimates. The test for heterogeneity is not significant for the log-odds-ratio estimates, suggesting that the meta-analysis based on the log-odds-ratio is preferable.

Ordinal data

Patient responses fall into one of m categories C_1, \dots, C_m which are ordered in terms of desirability: C_1 is the worst and C_m the best. For each study we have the information shown in Table IV.

The measure of treatment effect is the log-odds-ratio

$$\theta_i = \log \left\{ \frac{Q_{jCi}(1 - Q_{jEi})}{Q_{jEi}(1 - Q_{jCi})} \right\}$$

where

$$Q_{jEi} = p_{1Ei} + \dots + p_{jEi}, \quad Q_{jCi} = p_{1Ci} + \dots + p_{jCi}, \quad j = 1, \dots, m-1.$$

Table IV. Ordinal data from the i th study

ith study	Treatment		Overall
	E	C	
Model:			
Probability of falling into category			
C_1	p_{1Ei}	p_{1Ci}	p_{1i}
\vdots	\vdots	\vdots	\vdots
C_m	p_{mEi}	p_{mCi}	p_{mi}
Data:			
Number of patients in category			
C_1	n_{1Ei}	n_{1Ci}	n_{1i}
\vdots	\vdots	\vdots	\vdots
C_m	n_{mEi}	n_{mCi}	n_{mi}
Total	n_{Ei}	n_{Ci}	n_i

The value of θ_i is assumed not to depend on j . This is the proportional odds model, described by McCullagh.¹² Now define the following statistics:

$$L_{jEi} = n_{1Ei} + \dots + n_{(j-1)Ei}, \quad j = 2, \dots, m,$$

$$U_{jEi} = n_{(j+1)Ei} + \dots + n_{mEi}, \quad j = 1, \dots, m-1,$$

$L_{1Ei} = U_{mEi} = 0$, and make similar definitions for the control group. Then the efficient score is given by

$$Z_i = \frac{1}{n_i + 1} \sum_{j=1}^m n_{jEi} (L_{jCi} - U_{jCi})$$

and Fisher's information by

$$V_i = \frac{Z_i^2}{n_i + 2} + \frac{W_i}{(n_i + 1)(n_i + 2)}$$

where

$$W_i = \sum_{j=1}^m \{ n_{jEi}(n_{Ci} - n_{jCi}) + n_{jEi}n_{jCi}(n_i - n_{ji}) + 2n_{jEi}L_{jCi}U_{jCi} + 2n_{jCi}L_{jEi}U_{jEi} \}.$$

These statistics are given in Section 1.2.3 of Whitehead and Brunier¹³ and in Section 3.6 of Whitehead.⁸ Furthermore, Z_i is proportional to the Mann-Whitney statistic.¹⁴ Notice that the implied model is of stratified proportional odds. Proportionality of odds between centres is not assumed.

Normally distributed data

For normally distributed data, two measures of treatment effect, which might be considered are the absolute and the standardized difference between means. The advantages of using stand-

Table V. Normally distributed data from the i th study

ith study	Treatment		Overall
	E	C	
Model: mean	μ_{Ei}	μ_{Ci}	μ_i
s.d.	σ_i	σ_i	σ_i
Data: number	n_{Ei}	n_{Ci}	n_i
sum	s_{Ei}	s_{Ci}	s_i
sum of squares	q_{Ei}	q_{Ci}	q_i
mean	\bar{x}_{Ei}	\bar{x}_{Ci}	\bar{x}_i
s.d.	D_{Ei}	D_{Ci}	D_i

s.d. = standard deviation

ardized differences are that studies using different scales can be combined without conversion of units. One example would be the combination of studies with exercise tests carried out using a treadmill in some centres and an exercise cycle in others. Another example would be where a laboratory in one centre may produce internally consistent measurements which cannot, however, be directly compared with those from other laboratories. Standardization of differences allows their pooling in a meta-analysis.

For each study we have the information shown in Table V. The statistic D_i^2 is the usual pooled variance estimate,

$$D_i^2 = \frac{(q_{Ei} - s_{Ei}^2/n_{Ei}) + (q_{Ci} - s_{Ci}^2/n_{Ci})}{(n_i - 2)}.$$

Parameterizing by the standardized difference between means, $\theta_i = (\mu_{Ei} - \mu_{Ci})/\sigma_i$, the efficient score and Fisher's information are

$$Z_i = \frac{n_{Ei}n_{Ci}}{n_i D_i^*} (\bar{x}_{Ei} - \bar{x}_{Ci}) \quad \text{and} \quad V_i = \frac{n_{Ei}n_{Ci}}{n_i} - \frac{Z_i^2}{2n_i},$$

respectively, where $(D_i^*)^2 = (q_i - s_i^2/n_i)/n_i$. Now $D_i^* \doteq D_i$ if both θ_i and Z_i are small. The estimate $\hat{\theta}_i = Z_i/V_i \doteq (\bar{x}_{Ei} - \bar{x}_{Ci})/D_i$, and it follows that $w_i \doteq n_{Ei}n_{Ci}/n_i \doteq V_i$. The normal approximation (1) is improved if we use the estimate $\hat{\theta}_i = J(n_i - 2)(\bar{x}_{Ei} - \bar{x}_{Ci})/D_i$ where the function J is defined by $J(m) = 1 - 3/(4m - 1)$ when m is large, and is obtainable from Table 2 of Section 5.A.2, Hedges and Olkin,¹⁵ for small m . The corresponding weight is w_i where $w_i^{-1} = (n_i n_{Ei}^{-1} n_{Ci}^{-1} + 0.5 \hat{\theta}_i^2 n_i^{-1})$. Alternatively, treatment difference can be expressed as the absolute difference between means, $\theta_i = \mu_{Ei} - \mu_{Ci}$, and estimated by $\hat{\theta}_i = \bar{x}_{Ei} - \bar{x}_{Ci}$. For this choice, $w_i = n_{Ei}n_{Ci}/(n_i D_i^2)$. The normal approximation (1) is not nearly as accurate for the absolute difference parameterization and furthermore problems arise if different centres have different scales of measurement.

Survival data

Suppose that each study concerns the time between randomization to treatment and the occurrence of some event. In the i th study, events occur at the $m(i)$ distinct times $t_{1i}, \dots, t_{m(i)i}$. More than one event may be recorded at each of these times, due to rounding of the data or due to

Table VI. Survival data from the i th study

ith study	Treatment		Overall
	E	C	
Model: hazard function	$h_{Ei}(t)$	$h_{Ci}(t)$	
survivor function	$S_{Ei}(t)$	$S_{Ci}(t)$	
Data: number of events	e_{Ei}	e_{Ci}	e_i
number of survival times equal to			
t_{1i}	o_{1Ei}	o_{1Ci}	o_{1i}
t_{2i}	o_{2Ei}	o_{2Ci}	o_{2i}
\vdots	\vdots	\vdots	\vdots
$t_{m(i)i}$	$o_{m(i)Ei}$	$o_{m(i)Ci}$	$o_{m(i)i}$
number of survival times greater than or equal to			
t_{1i}	r_{1Ei}	r_{1Ci}	r_{1i}
t_{2i}	r_{2Ei}	r_{2Ci}	r_{2i}
\vdots	\vdots	\vdots	\vdots
$t_{m(i)i}$	$r_{m(i)Ei}$	$r_{m(i)Ci}$	$r_{m(i)i}$

the schedule of patient examinations. Thus o_{ji} events occur at time t_{ji} , $j = 1, \dots, m(i)$, the total number of events being $e_i = o_{1i} + \dots + o_{m(i)i}$. Some survival times may be right-censored, and so it is necessary to record separately for each j the number of patients r_{ji} with survival times of t_{ji} or more. These are the patients 'at risk' at time t_{ji} .

The model and data for the i th survival study are shown in Table VI. The hazard functions for the experimental treatment and control groups are given by h_{Ei} and h_{Ci} , respectively. Under the assumption of proportional hazards, the treatment effect could be measured by the log-hazard-ratio

$$\theta_i = -\log \left\{ \frac{h_{Ei}(t)}{h_{Ci}(t)} \right\}, \quad \text{for all } t > 0.$$

Equivalently, this is the difference between complementary log transformations of the survivor functions:

$$\theta_i = -\log \{ -\log S_{Ei}(t) \} + \log \{ -\log S_{Ci}(t) \}, \quad \text{for all } t > 0.$$

The efficient score is given by

$$Z_i = e_{Ci} - \sum_{j=1}^{m(i)} \frac{o_{ji} r_{jCi}}{r_{ji}}$$

and Fisher's information by

$$V_i = \sum_{j=1}^{m(i)} \frac{o_{ji}(r_{ji} - o_{ji}) r_{jEi} r_{jCi}}{(r_{ji} - 1) r_{ji}^2}.$$

These statistics are given in References 8 and 13, Sections 3.4 and 1.2.5 respectively. Furthermore Z_i is the logrank statistic and V_i its null variance (denoted by $O_A - E_A$ and V_A , respectively in

Section 14.6 of Armitage and Berry¹⁶). The model is stratified proportional hazards, with proportionality between centres not being assumed.

Published reports of survival studies are unlikely to report sufficient detail for the data of Table VI to be determined. The value of the χ^2 statistic for the logrank test may be quoted: approximately $Z_i^2/V_i = \chi^2$. Furthermore, e_i , the total number of events in the study may be given: approximately $V_i = \frac{1}{4}e_i$. From this information, Z_i and V_i may be roughly determined. Alternatively, Cox's proportional hazards model might have been fitted, and the coefficient corresponding to treatment might be quoted, together with its standard error. This coefficient is an estimate $\hat{\theta}_i$ of θ_i , and approximately

$$Z_i = \hat{\theta}_i \text{SE}(\hat{\theta}_i)^{-2}, \quad V_i = \text{SE}(\hat{\theta}_i)^{-2}.$$

Caution should be exercised in using these derived statistics. The authors of the published study may not be using the same conventions of terminology used here and so identification of the appropriate statistics may be difficult. Further, reliance is placed on the accuracy of their calculations.

DEALING WITH HETEROGENEITY

The variation of treatment effect from study to study will be of interest in a meta-analysis. A graphical presentation of the results is a useful first step. The radial plot described by Galbraith¹⁷ is a bivariate scatter plot (x, y) of the 'standardized estimate' of treatment effect ($\hat{\theta}_i\sqrt{w_i}$) against 'precision' ($\sqrt{w_i}$) for each study and is consistent with the methodology described in this paper. Radial plots of the log-odds-ratio estimates from Table II and of the probability difference estimates from Table III are presented in Figures 1 and 2, respectively. In a radial plot, uninformative small trials correspond to points lying close to the origin whereas informative large trials provide influential points on the right-hand edge of the plot. The least squares estimate of the slope of the regression line through the origin is equal to the estimate $\hat{\theta} = \sum \hat{\theta}_i w_i / \sum w_i$. The residual from the fitted regression line is equal to $(\hat{\theta}_i - \hat{\theta})\sqrt{w_i}$, which has a variance of $(1 - w_i / \sum w_j)$. As this variance is approximately equal to 1, a plot of the parallel lines $y = \hat{\theta}x \pm 2$ provides an approximate 95 per cent confidence band for individual study results. Trials which are not consistent with the overall picture are easily identified because they correspond to points falling outside this confidence band. The radial plots in Figures 1 and 2 show consistency among the log-odds-ratio estimates but not among the probability difference estimates.

The formal test of homogeneity is based on the statistic Q given in equation (2). If this is non-significant one may choose to ignore heterogeneity and calculate the overall estimate of θ as described in the previous section. If heterogeneity is found the reason for its presence could then be explored.

One reason for heterogeneity might be an inappropriate choice for the measure of treatment effect θ : this was illustrated in the antihypertensives example. A second reason could be a genuine study by treatment interaction caused by different design features or clinical procedures. Further investigation along these lines will require the use of linear models packages such as Genstat, GLIM or SAS: the elementary computations described in this paper are no longer sufficient. By including covariates pertaining to whole centres (such as country, average age or proportion of one racial group) in the model, we might ask whether a covariate by treatment interaction explained away most of the study by treatment interaction. If this were the case we might consider presenting individual estimates for different groups or strata. Imbalances in prognostic features between patients could cause heterogeneity. If baseline recordings pertaining to individual

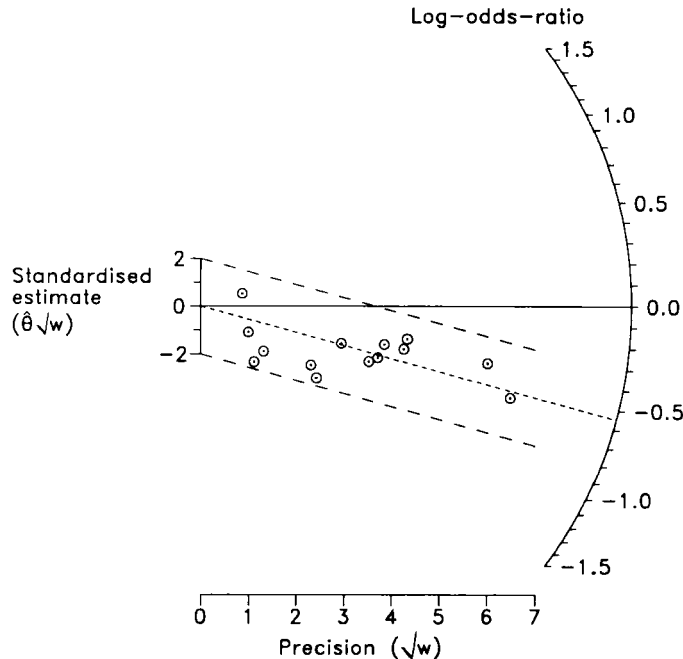


Figure 1. A radial plot of the log-odds-ratios in Table II, with approximate 95 per cent confidence band for individual study results. For any point, the log-odds-ratio is obtained by extrapolating a line from the origin (centre of circle) through that point to the circular scale

patients were available then these could be used as covariates. We might ask whether the study by treatment interaction was still significant once the baseline by treatment interaction was in the model. Analyses of subgroups of patients and of related variables can also throw light on possible heterogeneity. Canner¹⁸ presents an investigation of possible heterogeneity in six clinical trials of aspirin in coronary heart disease by performing a number of these statistical procedures.

If no reason is found to explain the heterogeneity we may still wish to answer the question 'Will the treatment achieve benefit on average?'. At this point we need to consider a model which allows for treatment effects to vary from study to study, namely a random effects model.

A GENERAL RANDOM EFFECTS PARAMETRIC APPROACH

Assume that the treatment effects from the k studies ($\theta_1, \dots, \theta_k$) are a sample of independent observations from $N(\theta, \tau^2)$. Suppose, as before, that the estimate $\hat{\theta}_i$ satisfies the distributional relationship $\hat{\theta}_i \sim N(\theta_i, w_i^{-1})$ where now $\theta_i \sim N(\theta, \tau^2)$. It follows that the marginal distribution of $\hat{\theta}_i$ is $\hat{\theta}_i \sim N(\theta, w_i^{-1} + \tau^2)$. The fixed effects estimate of θ , $\hat{\theta} \sim \sum \hat{\theta}_i w_i / \sum w_i$ still has mean θ , but its variance is now given by

$$\text{var } \hat{\theta} = \frac{\sum w_i^2 \text{var } \hat{\theta}_i}{(\sum w_i)^2} = \frac{\sum w_i^2 (w_i^{-1} + \tau^2)}{(\sum w_i)^2} = \frac{1}{\sum w_i} + \frac{\tau^2 \sum w_i^2}{(\sum w_i)^2}. \quad (3)$$

The following considerations provide an estimate for τ . The homogeneity test statistic Q is

$$Q = \sum w_i (\hat{\theta}_i - \hat{\theta})^2 = \sum w_i (\hat{\theta}_i - \theta)^2 - (\sum w_i) (\hat{\theta} - \theta)^2$$

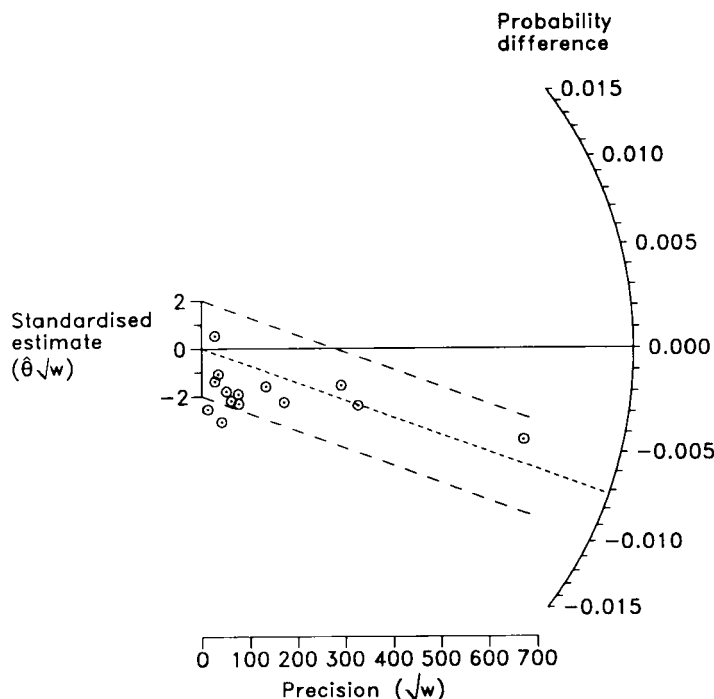


Figure 2. A radial plot of the probability differences in Table III, with approximate 95 per cent confidence band for individual study results

so that

$$\begin{aligned} E(Q) &= \sum w_i \text{var } \hat{\theta}_i - (\sum w_i) \text{var } \hat{\theta} \\ &= \sum w_i (w_i^{-1} + \tau^2) - (\sum w_i) \left\{ \frac{1}{\sum w_i} + \frac{\tau^2 \sum w_i^2}{(\sum w_i)^2} \right\}. \end{aligned}$$

That is,

$$E(Q) = (k - 1) + \tau^2 \left(\sum w_i - \frac{\sum w_i^2}{\sum w_i} \right).$$

This motivates use of the method of moments estimate $\hat{\tau}^2$ for τ^2 , where

$$\hat{\tau}^2 = \frac{Q - (k - 1)}{\sum w_i - (\sum w_i^2) / \sum w_i},$$

(Berlin *et al.*¹¹). The test of homogeneity, using Q , is a test of $H_0: \tau^2 = 0$. Should $\hat{\tau}^2 \leq 0$, a fixed effects analysis is more appropriate because this happens when $Q < E(Q; \tau^2 = 0) = k - 1$. If $\hat{\tau}^2 > 0$ we may use the approximate result $\hat{\theta}_i \sim N(\theta, w_i^{-1} + \hat{\tau}^2) \equiv N(\theta, (w_i^*)^{-1})$ where $w_i^* = (w_i^{-1} + \hat{\tau}^2)^{-1}$. This provides the test statistic $U^* = (\sum \hat{\theta}_i w_i^*)^2 / (\sum w_i^*)$ which follows a χ^2 distribution with 1 degree of freedom under the null hypothesis, $H_0: \theta = 0$. An improved estimate of θ is given by $\hat{\theta}^* = \sum \hat{\theta}_i w_i^* / (\sum w_i^*)$. Now $\hat{\theta}^*$ is asymptotically unbiased for θ , with variance approximately equal to $(\sum w_i^*)^{-1}$. Comparison with equation (3) and application of the Cauchy-Schwarz inequality shows that $\text{var } \hat{\theta}^* < \text{var } \hat{\theta}$ when $\tau^2 > 0$, and so $\hat{\theta}^*$ is to be preferred. The corresponding approximate 95 per cent confidence interval is $\hat{\theta}^* \pm 1.96 \sqrt{(1/\sum w_i^*)}$.

Table VII Example of a random effects model

Study	$\hat{\theta}_i$	w_i	Confidence interval				
			0	1	2	3	4
1	0.6	22	←→				
2	3.0	15				←→	
3	0.5	30	←→				

As an illustration of the random effects model consider the (rather extreme) example given in Table VII. This fictitious example has been exaggerated to clarify the consequences of a random effects analysis.

Here $\sum w_i = 67$, $\sum w_i^2 = 1609$, $\sum \hat{\theta}_i w_i = 73.2$ and $\sum \hat{\theta}_i^2 w_i = 150.4$. The fixed effects test statistic $U = 73.2^2/67 = 80.0$, which is highly significant. The fixed effects estimate of θ , $\hat{\theta} = 73.2/67 = 1.09$, with a 95 per cent confidence interval (0.85, 1.33). The test for homogeneity shows that $Q = 150.4 - (1.09)^2 67 = 70.8$ is highly significant. Proceeding to the random effects model, $\hat{\tau}^2 = (70.8 - 2)/(67 - 1609/67) = 1.60$, $\sum w_i^* = 1.82$, and $\sum \hat{\theta}_i w_i^* = 2.471$. The test statistic $U^* = 2.471^2/1.82 = 3.35$ which is not significant at the 5 per cent level and $\hat{\theta}^* = 2.471/1.82 = 1.36$ with a 95 per cent confidence interval (−0.09, 2.81). The estimate $\hat{\theta}^*$ does not even lie within the previous 95 per cent confidence interval. The apparently significant treatment effect has been overwhelmed by the variability between centres.

We can use the random effects model to consider the question 'What is the probability that the experimental treatment is superior to the control at a fourth centre?'. Now the treatment effect at this fourth centre, $\theta_4 \sim N(\theta, \tau^2) \equiv N(1.36, 1.60)$ so that $P(\theta_4 > 0) = \Phi\{1.36/\sqrt{1.60}\} = \Phi(1.075) = 0.86$. This would be an overestimate if the study centres had been chosen to be similar, rather than at random.

DISCUSSION

Whilst many authors have published actual meta-analyses of specific therapeutic questions, less has been published about comprehensive methodology. In this paper we have produced a general framework for the meta-analysis of clinical trials.

The methods of analysis discussed in this paper are appropriate for approximate calculations pooling the results from a large number of patients. For smaller total samples problems of biased estimation mentioned by Greenland and Salvan¹⁹ should be considered.

The suggestions given here allow application to response types other than those illustrated in the paper. The authors intend to do further work on trials which explore a common therapeutic question but use different patient response types. Another area deserving of further work is the robustness of the random effects approach, and in particular of the assumption of a normal random effects distribution. We feel that when random effects are seen to be present, it is likely that the approximate allowance proposed here will be far preferable to no allowance at all. In keeping with the likelihood approach to the fixed effects model, a likelihood-based random effects model could be devised. The methodology described by Ezzet and Whitehead²⁰ could be applied. Unfortunately this would be computationally intensive and is infeasible without special software.

We believe that the question of heterogeneity should be carefully examined in any meta-analysis. If heterogeneity is present but no explanation is found then a random effects model

should be used. In our account of the fixed-effects approach the assumption that all studies share the same treatment effect (that is $\theta_1 = \dots = \theta_k = \theta$) was made to justify the methods of estimation. If the θ_i are not equal, then the quantity $\hat{\theta} = \sum \hat{\theta}_i w_i / \sum w_i$ may remain an informative average measure of treatment comparison. However, out of the setting of a specific mathematical model, it is no longer an estimate of any parameter, nor can its standard error or associated confidence interval be found. In the context of the random effects model, $\hat{\theta}$ is an unbiased estimate of θ – now redefined as the population mean of the θ_i . For random effects, the unbiased estimate $\hat{\theta}^*$ of θ has smaller variance and is therefore preferable.

REFERENCES

1. Yusuf, S. 'Obtaining medically meaningful answers from an overview of randomised clinical trials', *Statistics in Medicine*, **6**, 281–286 (1987).
2. Light, R. J. 'Accumulating evidence from independent studies: what we can win and what we can lose', *Statistics in Medicine*, **6**, 221–228 (1987).
3. Goodwin, P. J. and Boyd, N. F. 'Mammographic parenchymal pattern and breast cancer risk: a critical appraisal of the evidence', *American Journal of Epidemiology*, **127**, 1097–1108 (1988).
4. Yusuf, S., Peto, R., Lewis, J., Collins, R. and Sleight, T. 'Beta blockade during and after myocardial infarction: an overview of the randomised trials', *Progress in Cardiovascular Disease*, **27**, 335–371 (1985).
5. DerSimonian, R. and Laird, N. 'Meta-analysis in clinical trials', *Controlled Clinical Trials*, **7**, 177–188 (1986).
6. Collins, R., Peto, R., MacMahon, S., Herbert, P., Fiebach, N. H., Eberlein, K. A., Godwin, J., Qizilbash, N., Taylor, J. O. and Hennekens, C. H. 'Blood pressure, stroke, and coronary heart disease Part 2, short-term reductions in blood pressure: overview of randomized drug trials in their epidemiological context', *Lancet*, **335**, 827–838 (1990).
7. Cochran, W. G. 'The combination of estimates from different experiments', *Biometrics*, **10**, 101–129 (1954).
8. Whitehead, J. *The Design and Analysis of Sequential Clinical Trials*, (2nd edition), Ellis Horwood, Chichester, 1991.
9. Anscombe, F. J. 'Normal likelihood functions', *Annals of the Institute of Statistical Mathematics*, **16**, 1–17 (1964).
10. Sprott, D. A. 'Normal likelihoods and their relation to large sample theory of estimation', *Biometrika*, **60**, 457–465 (1973).
11. Berlin, J. A., Laird, N. M., Sacks, H. S. and Chalmers, T. C. 'A comparison of statistical methods for combining event rates from clinical trials', *Statistics in Medicine*, **8**, 141–151 (1989).
12. McCullagh, P. 'Regression models for ordinal data', *Journal of the Royal Statistical Society, Series B*, **42**, 109–142 (1980).
13. Whitehead, J. and Brunier, H. *PEST 2.0 Planning and Evaluation of Sequential Trials, Operating Manual*, University of Reading, 1989.
14. Mann, H. B. and Whitney, D. R. 'On a test of whether one of two random variables is stochastically larger than the other', *Annals of Mathematical Statistics*, **18**, 50–60 (1947).
15. Hedges, L. V. and Olkin, J. *Statistical Methods for Meta-Analysis*, Academic Press, Orlando, 1985.
16. Armitage, P. and Berry, G. *Statistical Methods in Medical Research*, Blackwell Scientific Publications, Oxford, 1987.
17. Galbraith, R. F. 'A note on graphical presentation of estimated odds ratios from several clinical trials', *Statistics in Medicine*, **7**, 889–894 (1988).
18. Canner, P. 'An overview of six clinical trials of aspirin in coronary heart disease', *Statistics in Medicine*, **6**, 255–263 (1987).
19. Greenland, S. and Salvan, A. 'Bias in the one-step method for pooling study results', *Statistics in Medicine*, **9**, 247–252 (1990).
20. Ezzet, F. and Whitehead, J. 'Models for nested binary and ordinal data', in Decarli, A., Francis, B. J., Gilchrist, R. and Seeber, G. U. H. (eds.) *Lecture Notes in Statistics*, **57**, Springer, New York, 1989.