

Two-stage methods for the analysis of pooled data

Therese A. Stukel^{*,†}, Eugene Demidenko, James Dykes and Margaret R. Karagas

*Department of Community and Family Medicine, Section of Biostatistics and Epidemiology,
Dartmouth Medical School, Hanover, NH 03755-3863, U.S.A.*

SUMMARY

Epidemiologic studies of disease often produce inconclusive or contradictory results due to small sample sizes or regional variations in the disease incidence or the exposures. To clarify these issues, researchers occasionally pool and reanalyse original data from several large studies. In this paper we explore the use of a two-stage random-effects model for analysing pooled case-control studies and undertake a thorough examination of bias in the pooled estimator under various conditions. The two-stage model analyses each study using the model appropriate to the design with study-specific confounders, and combines the individual study-specific adjusted log-odds ratios using a linear mixed-effects model; it is computationally simple and can incorporate study-level covariates and random effects. Simulations indicate that when the individual studies are large, two-stage methods produce nearly unbiased exposure estimates and standard errors of the exposure estimates from a generalized linear mixed model. By contrast, joint fixed-effects logistic regression produces attenuated exposure estimates and underestimates the standard error when heterogeneity is present. While bias in the pooled regression coefficient increases with interstudy heterogeneity for both models, it is much smaller using the two-stage model. In pooled analyses, where covariates may not be uniformly defined and coded across studies, and occasionally not measured in all studies, a joint model is often not feasible. The two-stage method is shown to be a simple, valid and practical method for the analysis of pooled binary data. The results are applied to a study of reproductive history and cutaneous melanoma risk in women using data from ten large case-control studies. Copyright © 2001 John Wiley & Sons, Ltd.

1. INTRODUCTION

1.1. Pooled analyses

To improve the precision of the risk estimates in studies of the effects of exposures on disease, investigators often pool original data from several independent studies to increase the overall sample size and the generalizability of the results [1–6]. Meta-analysis refers to the statistical combination of the analytic results of several independent studies that are comparable in outcome and exposure for the purpose of integrating the findings. Issues regarding potential

*Correspondence to: Therese A. Stukel, Department of Community and Family Medicine, Section of Biostatistics and Epidemiology, Dartmouth Medical School, Hanover, NH 03755-3863, U.S.A.

†E-mail: stukel@dartmouth.edu

Contract/grant sponsor: NCI; contract/grant numbers: CA52192, CA62345

biases, statistical methods, sources of heterogeneity, validity, reliability and generalizability are discussed by several authors [7–10]. An investigation of sources of heterogeneity across studies, analysis approaches and the consequences of ignoring it are given by Colditz [11] and Berlin [12]. Statistical methods for combining summary results that take account of interstudy heterogeneity are described by Cochran [13], DerSimonian and Laird [14], Berlin *et al.* [15] and Zhou *et al.* [16]. Methods that incorporate study-level covariates to explain some of the heterogeneity include linear mixed-effects models [17, 18] and Bayesian hierarchical models [19–21]. Owing to the difficulty of obtaining original patient data, most meta-analyses combine summary results from published studies. However, outcomes and exposures are often defined differently enough across studies as to make combination of the results problematic.

In this paper, we examine statistical methods for pooled analyses, or meta-analyses using original patient data. When original data are available, there is flexibility in defining outcomes and exposures uniformly across studies in more precise and meaningful ways. By increasing the sample size, the pooling of individual patient data from many studies also allows for the examination of the effects of rare exposures, the effects of interactions among risk factors, the effects on subtypes of disease, and has the potential to resolve previously conflicting research results, which may not be possible with small individual studies [2, 22].

Inter-study variability exists in most situations where data from different studies are combined, due to the clustering of subjects within studies. Responses within such clusters are typically more homogeneous than those across different clusters because subjects within clusters are exposed to common environmental and study factors. For example, the demographic, clinical and life-style characteristics, area characteristics, and disease and exposure prevalences of a study population as well as study quality are more alike within than between studies. Some of the interstudy variability can be explained by study-level covariates. The additional unexplained interstudy variability can be measured using a random-effect term. The increase in standard error of the pooled exposure estimate is related to the degree of heterogeneity among studies, so that mild heterogeneity causes only a slight increase in the standard errors. The National Research Council recommends the routine use of random-effects approaches for analyses that combine information across studies and for the exploration of sources of variation in study results [23].

We assume that the goal of a pooled analysis is to estimate the average exposure effect β across studies as accurately and precisely as possible. We examine a two-stage method of analysis for pooled case-control data that combines the individual study-specific adjusted log-odds ratios using a linear mixed-effects model, and compare it to joint fixed-effects logistic regression, a commonly-used marginal model, assuming that a mixed-effects model is correct. We undertake a comprehensive study of bias in the pooled estimate of the exposure log-odds ratio β and its standard error using different study designs, study parameters, statistical methods and data collection methods. The results are generally applicable to different types of data and studies, such as cohort data with binary outcomes.

1.2. Model describing joint case-control studies

The model describing the joint case-control studies can be expressed in two stages. The first stage consists of the usual case-control model for each study with a uniformly-defined exposure variable [24]. The second-stage model combines the study-specific estimates to obtain a pooled estimate and can include study-level covariates. To establish notation, we define Y as a binary

variable denoting presence ($Y = 1$) or absence of disease, $p = \Pr(Y = 1)$ the probability of disease, X a uniformly defined exposure variable, and Z_k a confounder that may differ in definition across studies ($k = 1, \dots, K$). Stratum j in study k consists of n_{jk} cases and controls that are either pair or frequency matched. The outcome for individual i in stratum j of study k is denoted Y_{ijk} ($i = 1, \dots, n_{jk}$, $j = 1, \dots, S_k$, $k = 1, \dots, K$). To simplify, we assume only one confounder per study but this is easily generalized.

For each study, the first-stage model is a logistic regression describing the effect of the exposure X on disease, controlling for confounders Z_k , and can differ for each study. For study k ($k = 1, \dots, K$), the model is written,

$$\text{logit}(p_{ijk}) = \alpha_{jk} + X_{ik}\beta_k + Z_{ik}\gamma_k, \quad i = 1, \dots, n_{jk}, \quad j = 1, \dots, S_k \quad (1)$$

The exposure X is uniformly defined across studies. However, the confounders Z_k may be specific to a particular study and may vary in definition across studies. The exposure log-odds ratio for study k is denoted β_k , the confounding log-odds ratio γ_k , and the α_{jk} are study- and stratum-specific intercepts considered nuisance parameters. Model (1) thus describes each study separately using the age stratification and confounders relevant to that study. The β_k are assumed to vary across studies according to the second-stage model

$$\beta_k = \beta + b_k, \quad k = 1, \dots, K \quad (2)$$

where β is the pooled exposure log-odds ratio, and b_k are random effects with zero means and variance θ^2 . The b_k are random effects and represent the variability of the study-specific exposure effects β_k about the population mean β . The magnitude of this variability is denoted θ^2 and represents the extent of heterogeneity among individual study effects.

We assume that the true model describing the joint data is the mixed-effects model (1, 2) which explicitly accounts for subject- and study-specific variability. With appropriate notation and with common uniformly-defined confounders Z_k , it can be written as a joint mixed-effects logistic regression model [25]. In generalized linear mixed models (GLMMs) of this sort, estimation of β is complicated by the presence in the likelihood function of integrals that usually have no closed form expression, necessitating computer-intensive numerical integration techniques [26]. Approximate estimation techniques were proposed by Breslow and Clayton [27] based on penalized quasi-likelihood (PQL), and by Zeger *et al.* [28] based on generalized estimating equations (GEE). An alternative class of models is the marginal or population-averaged model that specifies the marginal distribution of the responses together with a ‘working’ correlation structure to account for intrastudy correlation without explicitly modelling it. Consistent estimates of the marginal model parameter estimates and standard errors can be obtained using GEE [29].

Neuhaus [30, 31] compared the two classes of models in terms of parameter estimates, standard errors and interpretation. Specifically, if the mixed-effects model is correct, a marginal model fitted to the data will produce an attenuated estimate of β , with the degree of attenuation related to the intrastudy correlation. Mixed models provide more powerful tests of within-study covariates, but tests for study-level covariates are approximately equivalent using the two approaches. As far as interpretation, the estimate of β from the mixed model measures the change in risk of disease with respect to the exposure X within a particular study context, since it is conditional on b_i . The estimate of β from the marginal model measures the average change in risk of disease with respect to the exposure X over all study populations. We are

assuming interest in the former type of inference in this paper, since it is similar to what most meta-analyses provide.

In his work, Neuhaus [30, 31] estimated the extent of attenuation in the regression coefficient of a logistic model with simple random sampling, one explanatory variable, no confounders, and no matching. In this paper we undertake a comprehensive examination of bias in the pooled estimator β in the more complex setting of pooled case-control study analysis with different study designs, study parameters and data collected across studies. We investigate two general methods for analysing these types of data: a two-stage random-effects model and a joint marginal model.

2. STATISTICAL METHODS

When important confounders are in common (or at least, measured) by all studies and can be uniformly defined, (1, 2) can be written as a joint generalized linear mixed model and any of the approximate methods described above for GLMMs can be used to estimate β . However, this is rarely the case, since pooled studies are usually analysed *post hoc*, and not all confounders are in common or are defined in such a way as to permit uniform coding. However, it remains a requirement of any joint model that covariates be measured by all studies in order to be included in the model.

The two models we examine are also approximations to (1, 2) but the two-stage method lends itself to flexibility with respect to differences in design, confounders and data collection across studies.

2.1. Two-stage estimation

An approximate estimation procedure for linear and non-linear mixed models is two-stage estimation [32–36]. In addition to being computationally simple, the practical advantage of the two-stage method for pooled analyses is its flexibility, since each study can be analysed using the method appropriate to the design, and confounders can be defined differently for each study. Estimation of β is implemented in two stages. First, $\hat{\beta}_k$ and its variance $\hat{\sigma}_k^2$ are estimated from the first-stage model (1), separately for every study k ($k = 1, \dots, K$). These estimates are then substituted in model (2) to obtain the linear mixed-effects model

$$\hat{\beta}_k = \beta + b_k + e_k, \quad k = 1, \dots, K \quad (3)$$

where the e_k are independent errors with zero means and variance $\hat{\sigma}_k^2$ describing the within-study variation of the $\hat{\beta}_k$, and β and b_k are defined as in (2). In (3), the $\hat{\beta}_k$ and $\hat{\sigma}_k^2$ are considered fixed. The marginal variance of $\hat{\beta}_k$ is thus the sum of the study-specific variance $\hat{\sigma}_k^2$ and the variance of the exposure effect across studies θ^2 . In fact, (3) is the model used in most meta-analyses of published study results [14, 17, 18].

The two-stage estimator of the pooled exposure effect β is a weighted average of the $\hat{\beta}_k$, weighting by the inverse marginal variances of the $\hat{\beta}_k$, denoted $w_k = (\hat{\sigma}_k^2 + \theta^2)^{-1}$. Thus

$$\hat{\beta} = \left(\sum_k w_k \hat{\beta}_k \right) / \sum_k w_k \quad (4)$$

$$\text{var}(\hat{\beta}) = (\sum_k w_k)^{-1} \quad (5)$$

When there is homogeneity of exposure effects across studies ($\theta = 0$), w_k simplifies to the inverse study-specific variances, $w_k = (\hat{\sigma}_k^2)^{-1}$. These estimators are identical to those proposed for meta-analysis by DerSimonian and Laird [14]. The estimation process is intuitive and computationally simple.

When there are study-level covariates that explain some of the heterogeneity across studies, the second-stage model (3) can be generalized to

$$\hat{\beta}_k = \mathbf{m}'_k \boldsymbol{\alpha} + b_k + e_k$$

where \mathbf{m}_k are study-specific covariates, and b_k and e_k are as defined above [17]. Estimation details are given by Berkey *et al.* [17], Pocock *et al.* [32] and Stukel and Demidenko [35].

2.2. Estimation of the random effects variance θ^2

To compute the pooled estimator (4), an estimate of the random effects variance θ^2 is required. Two methods are frequently used: pseudo-maximum likelihood and moment estimation.

2.2.1. Pseudo-maximum likelihood estimation of θ^2 . Assuming normal errors, an estimator for the variance of the random effects θ^2 is based on maximum likelihood (ML) using model (3), conditional on $\hat{\beta}_k$ and $\hat{\sigma}_k^2$. It uses the following recurrence formula:

$$\hat{\theta}_{(r+1)}^2 = \hat{\theta}_{(r)}^2 \frac{\sum_k (\hat{\beta}_k - \hat{\beta}_{(r)})^2 (\hat{\sigma}_k^2 + \hat{\theta}_{(r)}^2)^{-2}}{\sum_k (\hat{\sigma}_k^2 + \hat{\theta}_{(r)}^2)^{-1}}, \quad r = 0, 1, \dots \quad (6)$$

where $\hat{\beta}_{(r)}$ is recomputed at the r th iteration, $\hat{\theta}_0^2 = K^{-1} \sum_k [(\hat{\beta}_k - \bar{\beta})^2 - \hat{\sigma}_k^2]$ is an initial estimate of θ^2 and $\bar{\beta} = (\sum \hat{\beta}_k / \hat{\sigma}_k^2) / (\sum 1 / \hat{\sigma}_k^2)$, the weighted average of the study-specific $\hat{\beta}_k$, is an initial estimate of β [32, 34]. To obtain $\hat{\beta}$, one first computes $\bar{\beta}$ and $\hat{\theta}_0^2$ and then iterates between computing $\hat{\theta}^2$ from (6) and $\hat{\beta}_{(r)} = \hat{\beta}$ from (4) until convergence. This estimate is pseudo-ML because it conditions on the first-stage study-specific estimates, $\hat{\beta}_k$ and $\hat{\sigma}_k^2$. Berkey and Laird [34] also pointed out that this procedure was equivalent to the EM algorithm applied to the linear mixed-effects model (3), assuming the $\hat{\sigma}_k^2$ are fixed.

2.2.2. Moment estimation of θ^2 . The moment estimator of θ^2 is based on the ordinary least squares residuals from the second stage model (3), denoted $r_k = \hat{\beta}_k - \bar{\beta}$, where $\bar{\beta}$ is the weighted average of the $\hat{\beta}_k$. The sum of the squared residuals r_k is equated to its mathematical expectation to obtain the unbiased estimator

$$\hat{\theta}^2 = \left\{ \sum_k [\hat{\sigma}_k^{-2} (\hat{\beta}_k - \bar{\beta})^2] - (K - 1) \right\} / \left\{ \sum_k \hat{\sigma}_k^{-2} - \left(\sum_k \hat{\sigma}_k^{-4} / \sum_k \hat{\sigma}_k^{-2} \right) \right\}, \quad (7)$$

with $\hat{\theta}$ set to zero if it becomes negative [14]. This estimator is unbiased and non-iterative.

2.3. Joint fixed-effects logistic regression

Joint logistic regression refers to an analysis of the case-control studies by combining the data in a single model that requires common, uniformly-defined confounders, with stratification by study and age group [2]. The joint fixed-effects model is a marginal model that ignores interstudy heterogeneity so that $\beta_k = \beta$ and $\theta = 0$. When the extent of confounding is similar across studies, then $\gamma_k = \gamma$, and the model is written

$$\text{logit}(p_{ijk}) = \alpha_{jk} + X_{ik}\beta + Z_{ik}\gamma, \quad i = 1, \dots, n_{jk}, \quad j = 1, \dots, S_k, \quad k = 1, \dots, K \quad (8)$$

When the extent of confounding differs across studies, a study by confounder interaction term is included and the model is written

$$\text{logit}(p_{ijk}) = \alpha_{jk} + X_{ik}\beta + Z_{ik}\gamma_k, \quad i = 1, \dots, n_{jk}, \quad j = 1, \dots, S_k, \quad k = 1, \dots, K \quad (9)$$

As before, the exposure X is uniformly defined across studies. The confounders Z_k must also be measured by all studies and uniformly defined. The pooled exposure log-odds ratio is denoted β ; the confounding log-odds ratio is either γ or γ_k , depending on whether confounding varies by study. Pooling data from studies with different exposure prevalences and case-control ratios can result in biased estimates of effect. Whittemore *et al.* [2] recommend stratifying the analyses by 'study' in addition to age group using the α_{jk} to account for this potential confounding. Estimation is via conditional or unconditional logistic regression, depending on the size of the age strata [37]. In general, pairs are broken to create uniform age strata but this is not always necessary.

The practical drawback of a joint model is that it can include only covariates Z that have been measured by every study and are uniformly defined and coded across studies. Thus, a confounder, such as sunlight exposure, that is defined differently in every study, or number of nevi, that has not been measured by every study, cannot be controlled for in a joint model. This model may consequently produce biased overall exposure estimates for two reasons: ignoring random effects and ignoring important confounders.

3. SIMULATIONS TO COMPARE METHODS

3.1. Properties of the estimators using different statistical methods

To investigate the relative bias in the estimator of β and its standard error for the different methods in different situations, we used simulation. The general investigation assessed the effects of increasing study size N and increasing heterogeneity θ in exposure effects using different study designs, study parameters and statistical methods. All simulations assumed model (1, 2) was correct and varied key parameters: β , the pooled exposure effect across studies expressed as a log-odds ratio; θ , the degree of heterogeneity of the exposure effects across studies; $\bar{\gamma}$, the average confounding effect across studies expressed as a log-odds ratio; $\text{SD}(\gamma_k)$, the extent of confounding across studies; and $\text{OR}(X, Z)$, the odds ratio between the exposure X and the confounder Z , as well as number of studies, prevalence of the exposure (0.05 to 0.95), prevalence of the confounder (0.05 to 0.35), and the confounding odds ratio, $\exp(\gamma)$, from 0.89 to 5.42. Each design was assessed under many combinations of these factors.

We simulated data corresponding to the designs of the ten melanoma case-control studies described in the example using model (1, 2) and analysed the studies first using a correctly specified model (all confounders were measured by all studies, uniformly defined and incorporated into the models). We also examined a misspecified model where confounder information that was important for some studies was not measured by every study. These confounders could easily be incorporated into the two-stage but not the joint models.

We compared five methods of estimation: (i) joint fixed-effects logistic regression with the same extent of confounding γ across studies (JLR1); (ii) joint logistic regression with different extent of confounding γ_k across studies (JLR2); (iii) two-stage fixed-effects models ($\theta=0$) (TS0), (iv) two-stage method with ML estimation of random effects (TSML); (v) two-stage method with moment estimation of random effects (TSMM). We used conditional logistic regression, with the study strata as defined, to analyse the data. For the sake of comparison, we also used unconditional logistic regression, creating uniform five-year age strata for all studies and breaking pairs when necessary. The algorithm used to generate data based on model (1, 2) is described in the appendix. The statistical optimization code for all methods was written in the programming language C. The algorithm for maximizing the conditional logistic likelihood function was based on Smith *et al.* [38] and Krailo [39].

For each design specification, we generated 5000 sets of ten studies, and for each set of studies and each method, we computed the pooled estimate $\hat{\beta}$ and the model-based standard error based on (5). From the sampling distributions, we computed the empirical means of the $\hat{\beta}$ and their model-based standard errors, and the empirical standard error of $\hat{\beta}$ over all 5000 simulations. Per cent relative bias of the empirical mean of the $\hat{\beta}$ is computed with respect to the true β , and per cent relative bias of the empirical mean of the model-based standard error of $\hat{\beta}$ is computed with respect to the empirical standard error of $\hat{\beta}$, considered the true standard error.

Figure 1, relating study bias to study sample size N , shows that the two-stage random-effects method produces consistent and efficient estimates of β and its standard error, whether or not heterogeneity is present. However, joint fixed-effects models are consistent only when there is homogeneity in the exposure effects ($\theta=0$); relative bias in $\hat{\beta}$ and the standard error do not decrease with N as expected when heterogeneity is present ($\theta=0.82$).

Figure 2 demonstrates the effect of increasing heterogeneity θ in exposure effects across studies when $N=700$. When the models are correct, bias in $\hat{\beta}$ is negative and increases with heterogeneity for both models but is much smaller for the two-stage random-effects models. The two-stage models show little bias in $\hat{\beta}$ and its standard error when the important confounders for each study are measured and can be incorporated into the study-specific models, even if they differ from study to study. However, when confounders that are important for some studies have not been collected by every study, joint models may not be feasible; they produce biased estimates of $\hat{\beta}$ and its standard error since they cannot properly control for confounding. Standard errors for the two-stage random-effects estimators are relatively unbiased compared to the large underestimation observed with methods that do not incorporate random effects. Parameter estimates using the moment estimator of the random effects variance (TSMM) were generally less biased than those using the pseudo-maximum likelihood estimates (TSML).

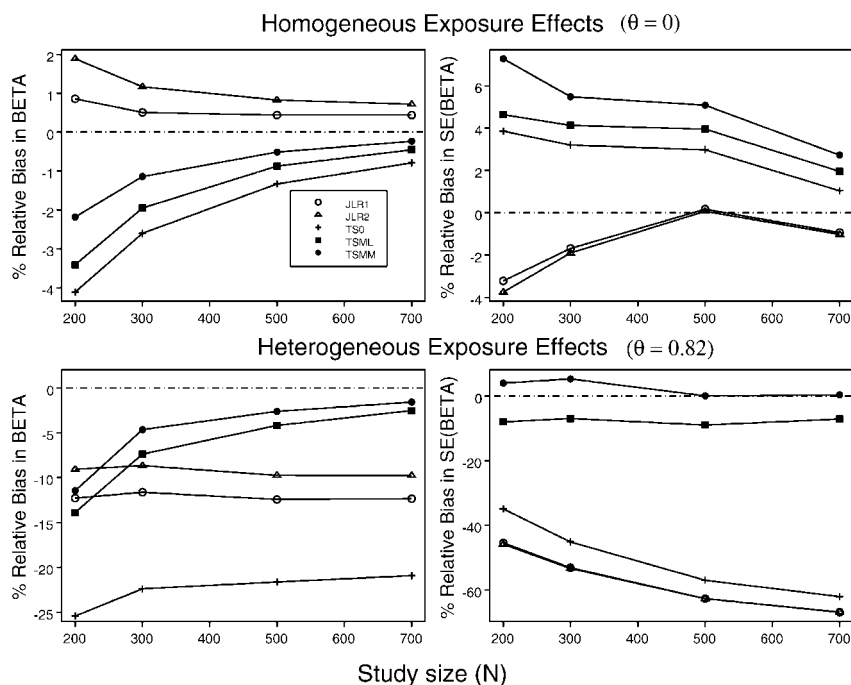


Figure 1. Plot of per cent relative bias in $\hat{\beta}$ and standard error of $\hat{\beta}$ against study size N when exposure effects are homogeneous ($\theta = 0$) or heterogeneous ($\theta = 0.82$) across studies, assuming the mixed model (1, 2) is correct. JLR1 refers to joint fixed-effects conditional logistic regression with similar extent of confounding across studies using model (8); JLR2 refers to joint fixed-effects conditional logistic regression with extent of confounding differing among studies using model (9); TS0 refers to the two-stage fixed-effects method; TSML and TSMML refer to the two-stage method of estimation with maximum likelihood and moment estimation of the random effects variance, respectively.

These results were similar for simulated data based on identically designed studies that were either all pair matched or all frequency matched, and were independent of the magnitude of β , $\bar{\gamma}$, $SD(\gamma_k)$ and $OR(X, Z)$ (results not shown).

4. EXAMPLE

Studies of reproductive events, use of oral contraceptives and non-contraceptive hormones and melanoma risk in women have yielded inconsistent results. To help clarify these issues, data from ten large melanoma case-control studies completed by July 1994 were pooled and analysed. The study designs ranged from pair-matched on age and area of residence to frequency-matched on five- or ten-year age groups and occasionally on area of residence. To facilitate comparability, we collected data only on studies that included at least 100 female melanoma cases and 100 controls, used personal interviews (in person or telephone), included melanoma cases treated on either an inpatient or outpatient basis, and collected information on pigimentary characteristics and sun exposure.

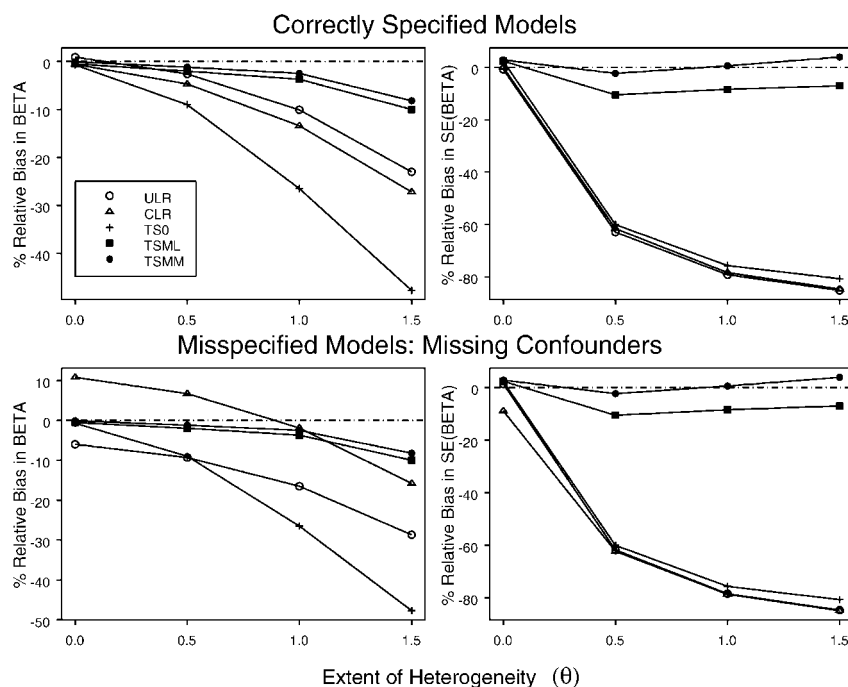


Figure 2. Plot of per cent relative bias in $\hat{\beta}$ and standard error of $\hat{\beta}$ against extent of heterogeneity θ when $\beta = 1$, for correctly specified and misspecified joint models, assuming the mixed model (1, 2) is correct. ULR refers to joint unconditional fixed-effects logistic regression with uniform five-year age strata; CLR refers to joint conditional fixed-effects logistic regression with stratification according to the study design; TS0 refers to the two-stage fixed-effects method; TSML and TSMM refer to the two-stage method of estimation with maximum likelihood and moment estimation of the random effects variance, respectively.

These ten studies represent six countries of varying latitudes with melanomas diagnosed between 1979 and 1987 and sample sizes ranging from 236 to 1387 subjects. Controls were selected using pair matching in five studies and frequency matching in five studies. In three frequency-matched studies, the matching criteria included age and area of residence and in two, age alone; some studies were more finely matched than others. Questions on oral contraceptives and reproductive history, the primary exposure variables, were collected in a similar fashion. While data on multiple potentially confounding or modifying factors were available in every study, some factors, such as number of nevi, were not measured in all studies. Even when covariates were collected by all studies, the categories used were not always the same. For example, all studies obtained hair colour and most studies classified it as red, blonde, brown and black; however, a few studies grouped light brown with blonde, or dark brown with black hair. A full description of these studies appears in Karagas *et al.* [40].

We analysed the risk of melanoma with respect to age at first birth, classified as <20 or ≥ 20 years of age. Since this was based on the subset of parous women, many strata were incomplete, necessitating a regrouping of the reference age stratum to three categories across all studies (<35 , $35-44$, $45+$ years of age). Two studies asked only about age at pregnancy

Table I. Effect of age at first birth on melanoma, adjusted for hair colour and number of Nevi.

Study	Unadjusted * odds ratio*	95 per cent CI		Confounders	Adjusted odds ratio†	95 per cent CI	
Elwood [41]	0.64	0.30	1.35	HAIR, NEVI	0.58	0.26	1.29
Gallagher [42]	2.22	1.31	3.78	HAIR	2.14	1.25	3.66
Holly [43]	0.91	0.62	1.35	HAIR, NEVI	0.77	0.50	1.18
Kirkpatrick [44]	0.61	0.27	1.38	HAIR, NEVI	0.53	0.23	1.25
Mack	1.12	0.70	1.80	HAIR, NEVI	1.30	0.77	2.10
Osterlind [45]	2.06	1.32	3.24	HAIR, NEVI	2.26	1.40	3.63
Swerdlow [46]	3.35	0.79	14.26	HAIR, NEVI	5.29	1.00	28.0
Zanett [47]	0.69	0.23	2.05	HAIR	0.67	0.22	2.07
<i>Pooled exposure estimates</i>							
Homogeneity	1.23	0.96	1.58		1.23	0.95	1.59
Heterogeneity	1.18	0.63	2.20		1.17	0.55	2.52

* Adjusted for reference age.

† Adjusted for reference age, hair colour and nevi, where possible.

and were excluded from these analyses. We used unconditional logistic regression in the first-stage models, controlling for study by reference age group strata, to analyse the effect of age at first birth on risk of melanoma. Controlling for age alone, the individual study-specific exposure estimates ranged from an odds ratio of 0.61 to 3.35; only two studies produced statistically significant positive odds ratios (Table I). There was significant heterogeneity among the exposure estimates ($\chi^2(7)=21.04$; $p=0.004$) so that the pooled exposure estimate was computed using a second-stage random-effects model (3), with a moment estimate of the random effects variance. The pooled exposure odds ratio (OR) was 1.18 (95 per cent confidence interval (CI), 0.63 to 2.20). Controlling simultaneously for age and hair colour did not substantially change the pooled estimator (OR 1.19; 95 per cent CI, 0.58 to 2.45).

However, controlling simultaneously for age, hair color and nevi was possible only in the six studies that collected information on nevi (Table I). Incorporating the remaining two studies in the pooled odds ratio was possible by adjusting these for age and hair colour alone. The adjusted pooled odds ratio using a two-stage random-effects estimator was similar (OR=1.17, 95 per cent CI, 0.55 to 2.52), since neither hair colour nor nevi was a strong confounder. However, if nevi were a confounding variable, then combining results across these studies would not produce a valid pooled estimate. This is a standard problem encountered by pooled analyses that cannot be overcome by modelling. Many meta-analyses report results both with and without these studies. Olkin [48] recommends sensitivity analyses to determine whether meta-analysis results hold up under modest perturbations of the data. To examine the influence of the missing confounder on the pooled estimator, a sensitivity analysis was performed. In the two studies that were missing the nevi variable, the log-odds ratios $\hat{\beta}$ were simultaneously shifted by ± 15 per cent to simulate moderate confounding. The resulting ORs and 95 per cent CIs were (1.15; 0.54 to 2.45) and (1.20; 0.56 to 2.58), indicating that in this case, moderate confounding would minimally change the overall pooled estimate. As Olkin [48] points out, caution is required in any meta-analysis to ensure that a few studies are not unduly influential.

5. DISCUSSION

If confounding can be controlled in a pooled analysis, the remaining important factors affecting bias are study sample size and degree of heterogeneity across studies. When the individual studies are large, two-stage estimation methods are shown to be valid for the analysis of pooled heterogeneous data and to produce nearly unbiased parameter estimates and standard errors, without the practical drawbacks of a joint regression model. This is due to their better ability to control for confounding in individual studies and the inclusion of random effects in the models. Large sample properties of logistic regression ensure that the individual exposure log-odds ratios $\hat{\beta}_k$ are nearly unbiased provided each case-control study is large and models properly control for confounding [49, 50]. A weighted average of study-specific estimates, weighted by the inverse marginal variances (5), is the best linear unbiased estimate (BLUE) of the pooled exposure estimate when the variances are known. As such, the two-stage pooled estimator is expected to be unbiased when the individual $\hat{\beta}_k$ are unbiased, and have minimum variance. Although it is not maximum likelihood, it is asymptotically equivalent when the individual studies are large [51]. Although it does not require the uniform recoding of confounders that is required for joint models, one cannot obtain a valid pooled estimate when some studies are missing important confounders. It is unclear how well the two-stage method would perform if individual studies were smaller, however, a few small studies would likely not influence these results since their weights in (4) and (5) would be relatively small.

Parameter estimates using the moment estimator of the random effects variance (TSMM) were slightly less biased than those using the ML estimator (TSML) possibly due to its unbiasedness and robustness to normality assumptions. However, when there are few studies, ML underestimates the random effects variance θ , and moment estimation is inefficient, leading to biased estimates of β . In this situation, Morris [52] derives a Stein-type estimator using a parametric empirical Bayesian approach when several independent unbiased estimates of $\hat{\beta}_k$ and their variances are available, as in the meta-analysis situation.

Joint fixed-effects logistic regression produces a maximum likelihood estimator that has optimal statistical properties when the marginal model is correct and the sample size is large [49]. In other words, it produces unbiased estimates only when studies are homogeneous with respect to exposure effects, and all important confounders are measured and can be uniformly defined by all studies, a situation that rarely occurs in practice. The bias in the joint model estimates is a combination of ignoring heterogeneity (negative bias) [30, 31] and residual confounding (positive bias), when all relevant confounders are not included in the model.

Although linear regression models that omit random effects still produce unbiased regression estimates, this is not true for generalized linear mixed models. The omission of random effects in such models leads to an attenuation of the regression parameter, producing a negative bias [30, 31], similar to ignoring measurement error in a covariate in a non-linear model [53]. Using GLMMs would be preferable, but currently, the estimation methods are all approximations to maximum likelihood, and are computationally intensive due to the presence in the likelihood function of an integral having no closed form solution, necessitating numerical integration [25]. Inclusion of confounders that are not measured or uniformly defined by all studies is still a practical drawback to any joint model.

Berkey *et al.* [17] investigated the sampling distribution of $\hat{\beta}$ via simulations and found that for models such as (3) with no study-level covariates, it followed a t_{K-4} distribution. We used the sampling distribution of the two-stage estimator of $\hat{\beta}$ to evaluate the tail probabilities

of the Wald statistic, $t = \hat{\beta}/\text{SE}(\hat{\beta})$ using a model-based estimator of the standard error (5). We obtained inconclusive results using both the normal and various t -distributions and were not able to confirm these results.

Two-stage methods are simple and provide a flexibility for the analysis of pooled data that joint models do not allow. Analysing studies one at a time allows for better control of confounding when studies have different confounders coded in a unique manner or when they have different designs, such as in combining pair- and frequency-matched case-control studies. In practice, the individual models can be refined until residual confounding is removed using data from that study. It is then computationally straightforward to run and store the individual model estimates and to combine the results in a second stage when all the first-stage models are satisfactory. This is particularly useful when there are many exposure variables of interest, necessitating the examination of large numbers of different models. The two-stage method is thus a practical and valid method for the analysis of pooled data within the constraints of the data collected; like any meta-analysis, it cannot control for confounding with respect to variables that were not measured. These results are applicable to other types of data and studies such as cohort data with binary or continuous outcomes.

ACKNOWLEDGEMENTS

We are grateful to the International Melanoma Analysis Group (IMAGE) for use of their data, Debra Whitney and Matt Siano for programming assistance, Barbara Moskalenko for assistance with manuscript preparation, and two reviewers for their helpful comments. This research was supported by NCI grants CA52192 and CA62345.

APPENDIX: DATA GENERATION FOR CASE-CONTROL STUDIES

A1. Model for data generation

The objective is to provide an algorithm to generate binary exposure and confounder variables for cases ($y=1$) and controls ($y=0$), given known values of β, γ , prevalence of the exposures and confounders in each stratum, and the odds ratio between X and Z . For each study and stratum, we wish to generate n_0 control pairs (x_i, z_i) and $n - n_0$ case pairs. For a pair-matched study, $n_0 = 1$, $n = 2$. The underlying model for a cohort study is

$$\Pr(y=1 | x, z) = \frac{\exp(\alpha + \beta x + \gamma z)}{1 + \exp(\alpha + \beta x + \gamma z)} \quad (\text{A1})$$

where x (exposure) and z (confounder) are binary variables. We generate pairs (x, z) based on the well-known formulae for case-control studies [49, 54]

$$\begin{aligned} \Pr(x, z | y=1) &= \frac{\Pr(y=1 | x, z)\Pr(x, z)}{1 - \Pr(y=0)} \\ \Pr(x, z | y=0) &= \frac{\Pr(y=0 | x, z)\Pr(x, z)}{\Pr(y=0)} \end{aligned} \quad (\text{A2})$$

where $\Pr(y=1|x,z)$ is specified in (A1), $\Pr(x,z)$ is the marginal distribution of the (x,z) in the stratum and $\Pr(y=1)=w=(n-n_0)/n$. Parameters α and $\Pr(x,z)$ are determined below based on the marginal probabilities for x and z and the odds ratio between them.

A2. Determination of α and $\Pr(x,z)$

To generate the data, we must know the intercept α which is based on the known incidence rate w , the log-odds ratio between X and Z , denoted ρ , and the marginal probabilities, $p_x = \Pr(x=1)$ and $p_z = \Pr(z=1)$. First, we define the association between z and x as the conditional probability

$$\Pr(z=0|x) = [1 + \exp(v + \rho x)]^{-1} \quad (\text{A3})$$

where v is based on ρ , p_x , p_z . Since

$$\begin{aligned} 1 - p_z &= \Pr(z=0) \\ &= \Pr(z=0|x=0)\Pr(x=0) + \Pr(z=0|x=1)\Pr(x=1) \\ &= (1 - p_x)/(1 + \exp(v)) + p_x/(1 + \exp(v + \rho)) \end{aligned} \quad (\text{A4})$$

v can be found from the above, which collapses to a quadratic equation for $\exp(v)$. We find α based on the known incidence rate w . Using standard formulae for conditional probability, we obtain

$$\begin{aligned} \Pr(y=0) &= \Pr(y=0|x=0,z=0)\Pr(x=0,z=0) \\ &\quad + \Pr(y=0|x=1,z=0)\Pr(x=1,z=0) \\ &\quad + \Pr(y=0|x=0,z=1)\Pr(x=0,z=1) \\ &\quad + \Pr(y=0|x=1,z=1)\Pr(x=1,z=1) = 1 - w \end{aligned}$$

The required probabilities are found from (A3)

$$\begin{aligned} \pi_{11} &= \Pr(x=1,z=1) = \Pr(z=1|x=1)\Pr(x=1) = \frac{p_x \exp(v + \rho)}{1 + \exp(v + \rho)} \\ \pi_{10} &= \Pr(x=1,z=0) = \Pr(z=0|x=1)\Pr(x=1) = \frac{p_x}{1 + \exp(v + \rho)} \\ \pi_{01} &= \Pr(x=0,z=1) = \Pr(z=1|x=0)\Pr(x=0) = \frac{(1 - p_x) \exp(v)}{1 + \exp(v)} \\ \pi_{00} &= \Pr(x=0,z=0) = \Pr(z=0|x=0)\Pr(x=0) = \frac{(1 - p_x)}{1 + \exp(v)} \end{aligned}$$

Therefore, denoting $\kappa = \exp(\alpha)$, $B = \exp(\beta)$, $C = \exp(\gamma)$ we come to the following equation for κ :

$$\frac{\pi_{00}}{1 + \kappa} + \frac{\pi_{10}}{1 + \kappa B} + \frac{\pi_{01}}{1 + \kappa C} + \frac{\pi_{11}}{1 + \kappa CB} = 1 - w \quad (\text{A5})$$

This equation is solved by using the Newton algorithm beginning with $\kappa_0 = 0$. Since the left side of the equation is a decreasing convex function, the sequence of approximations $\{\kappa_s, s = 0, 1, 2, \dots\}$ is increasing and converges to the solution.

A3. Algorithm for data generation

To generate $(x, z | y = 1)$ and $(x, z | y = 0)$, we use formulae (A2) which lead to

$$\Pr(x, z | y = 1) = \frac{n}{n - n_0} \frac{\exp(\alpha + \beta x + \gamma z)}{1 + \exp(\alpha + \beta x + \gamma z)} \Pr(x, z) \quad (\text{A6})$$

and

$$\Pr(x, z | y = 0) = \frac{n}{n_0} \frac{1}{1 + \exp(\alpha + \beta x + \gamma z)} \Pr(x, z) \quad (\text{A7})$$

where $\Pr(x, z)$ are defined via the π . Thus, the generation of binary data for case-control studies consists of four steps:

1. Given p_x, p_z and ρ , find v from (A4) using the quadratic equation for $\exp(v)$.
2. Find $\kappa = \exp(\alpha)$ solving (A5), where $w = (n - n_0)/n$.
3. Calculate the conditional probabilities (A6) and (A7).
4. Generate $n - n_0$ binary pairs (x, z) with distribution (A6) for cases, and n_0 pairs with distribution (A7) for controls.

REFERENCES

1. Negri E, La Vecchia C, Bruzzi P, Dardanoni G, Decarli A, Palli D, Parazzini F, Rosselli del Turco M. Risk factors for breast cancer: pooled results from three Italian case-control studies. *American Journal of Epidemiology* 1988; **128**:1207–1215.
2. Whittemore AS, Harris R, Itnyre J, Halpern J and the Collaborative Ovarian Cancer Group. Characteristics related to ovarian cancer risk: collaborative analysis of 12 case-control studies. I. Methods. *American Journal of Epidemiology* 1992; **136**:1175–1183.
3. Whittemore AS, Harris R, Itnyre J and the Collaborative Ovarian Cancer Group. Characteristics related to ovarian cancer risk: Collaborative analysis of 12 case-control studies. II. Intensive epithelial ovarian cancers in white women. *American Journal of Epidemiology* 1992; **136**:1184–1203.
4. Howe GR, Benito E, Castelletto R. Dietary intake of fiber and decreased risk of cancers of the colon and rectum: evidence from the combined analysis of 13 case-control studies. *Journal of the National Cancer Institute* 1992; **84**:1887–1896.
5. Skegg DC, Noonan EA, Paul C, Spears GF, Meirik O, Thomas DB. Depot medroxyprogesterone acetate and breast cancer. A pooled analysis of the World Health Organization and New Zealand studies. *Journal of the American Medical Association* 1995; **273**:799–804.
6. Smith-Warner SA, Spiegelman D, Yuan S-S, van den Brandt PA, Folsom AR, Goldbohm RA, Graham S, Holmberg L, Howe GR, Marshall JR, Miller AB, Potter JD, Speizer FE, Willett WC, Wolk A, Hunter DJ. Alcohol and breast cancer in women. A pooled analysis of cohort studies. *Journal of the American Medical Association* 1998; **279**:535–540.
7. Dickersin K, Berlin JA. Meta-analysis: state-of-the-science. *Epidemiologic Reviews* 1992; **14**:154–176.
8. Olkin I. Statistical and theoretical considerations in meta-analysis. *Journal of Clinical Epidemiology* 1995; **48**:133–146.
9. Olkin I. Meta-analysis: current issues in research synthesis. *Statistics in Medicine* 1996; **15**:1253–1257.

10. Egger M, Smith GD. Meta-analysis: potentials and promise. *British Medical Journal* 1997; **315**:1371–1374.
11. Colditz GA, Burdick E, Mosteller F. Heterogeneity in meta-analysis of data from epidemiologic studies: A commentary. *American Journal of Epidemiology* 1995; **142**:371–382.
12. Berlin JA. Invited commentary: Benefits of heterogeneity in meta-analysis of data from epidemiologic studies. *American Journal of Epidemiology* 1995; **142**:383–387.
13. Cochran, WG. The combination of estimates from different experiments. *Biometrics* 1954; **10**:101–129.
14. DerSimonian R, Laird, N. Meta-analysis in clinical trials. *Controlled Clinical Trials* 1986; **7**:177–188.
15. Berlin JA, Laird NM, Sacks HS, Chalmers TC. A comparison of statistical methods for combining event rates from clinical trials. *Statistics in Medicine* 1989; **8**:141–151.
16. Zhou XH, Brizendine EJ, Pritz MB. Methods for combining rates from several studies. *Statistics in Medicine* 1999; **18**:557–566.
17. Berkey CS, Hoaglin DC, Mosteller F, Colditz GA. A random-effects regression model for meta-analysis. *Statistics in Medicine* 1995; **14**:395–411.
18. Stram DO. Meta-analysis of published data using a linear mixed-effects model. *Biometrics* 1996; **52**:536–544.
19. Morris CN, Normand SL. Hierarchical models for combining information and for meta-analyses. In *Bayesian Statistics 4*. Oxford University Press: Oxford, UK, 1992; 321–344.
20. Carlin JB. Meta-analysis for 2×2 tables – a Bayesian approach. *Statistics in Medicine* 1992; **11**:141–158.
21. Schmid CH. Exploring heterogeneity in randomized trials via meta-analysis. *Drug Information Journal* 1999; **33**:211–224.
22. Stewart LA, Clarke MJ. Practical methodology of meta-analyses (overviews) using updated individual patient data. *Statistics in Medicine* 1995; **14**:2057–2079.
23. National Research Council. *Combining Information, Statistical Issues and Opportunity for Research*. National Academy Press: Washington, D.C. 1992; 41–46.
24. Breslow NE, Day NE. *Statistical Methods in Cancer Research, Volume I. The Analysis of Case-Control Studies*. IARC: Lyon, 1980.
25. Stiratelli R, Laird N, Ware JH. Random-effects models for serial observations with binary response. *Biometrics* 1984; **40**:961–971.
26. Diggle PJ, Liang KY, Zeger SL. *Analysis of Longitudinal Data*. Clarendon Press: Oxford, 1995.
27. Breslow NE, Clayton DG. Approximate inference in generalized linear mixed models. *Journal of the American Statistical Association* 1993; **88**:9–25.
28. Zeger SL, Liang KY, Albert PS. Models for longitudinal data: a generalized estimating equation approach. *Biometrics* 1988; **44**:1049–1060.
29. Zeger SL, Liang KY. Longitudinal data analysis for discrete and continuous outcomes. *Biometrics* 1986; **42**:121–130.
30. Neuhaus JM, Kalbfleisch JD, Hauck WW. A comparison of cluster-specific and population-averaged approaches for analyzing correlated binary data. *International Statistical Review* 1991; **59**:25–35.
31. Neuhaus JM. Estimation efficiency and tests of covariate effects with clustered binary data. *Biometrics* 1993; **49**:989–996.
32. Pocock SJ, Cook DG, Beresford SAA. Regression of area mortality rates on explanatory variables: What weighting is appropriate? *Applied Statistics* 1981; **30**:286–295.
33. Steimer JL, Mallet A, Golmard JL, Boisivieux JF. Alternative approaches to estimation of population pharmacokinetic parameters: Comparison with the nonlinear mixed-effect model. *Drug Metabolism Reviews* 1984; **15**:265–292.
34. Berkey CS, Laird NM. Nonlinear growth curve analysis: estimating the population parameters. *Annals of Human Biology* 1986; **13**:111–128.
35. Stukel TA, Demidenko E. Two-stage method of estimation for general linear growth curve models. *Biometrics* 1997; **53**:720–728.
36. Stukel TA, Demidenko E. Comparison of methods for general nonlinear mixed-effects models. In *Modelling Longitudinal and Spatially Correlated Data: Methods, Applications, and Future Directions*. Springer Lecture Notes in Statistics. Springer-Verlag: New York, 1997; 135–146.
37. Pike MC, Hill AP, Smith PG. Bias and efficiency in logistic analyses of stratified case-control studies. *International Journal of Epidemiology*, 1980; **9**:89–95.
38. Smith PG, Pike MC, Hill P, Breslow NE, Day NE. Multivariate conditional logistic analysis of stratum-matched case-control studies. Algorithm AS 162. *Applied Statistics*, 1981; **30**:190–197.
39. Kralio MD, Pike MC. Conditional multivariate logistic analysis of stratified case-control studies. Algorithm AS 196. *Applied Statistics*, 1984; **33**:95–103.
40. Karagas MR, Dykes J, Stukel TA, Armstrong B, Elwood M, Gallagher R, Green A, Holly E, Kirkpatrick C, Langholz B, Mack T, Osterlind AS 162. *Applied Statistics*, 1981; **30**:190–197.
41. Elwood JM, Whitehead SM, Davison J, Stewart M, Galt M. Malignant melanoma in England: risks associated with naevi, freckles, social class, hair colour, and sunburn. *International Journal of Epidemiology* 1990; **19**:801–810.

42. Gallagher RP, Elwood JM, Hill GB, Coldman AJ, Threlfall WJ, Spinelli JJ. Reproductive factors, oral contraceptives and risk of malignant melanoma: Western Canada Melanoma Study. *British Journal of Cancer*, 1985; **52**:901–907.
43. Holly EA, Cress RD, Ahn DK. Cutaneous melanoma in women. III. Reproductive factors and oral contraceptive use. *American Journal of Epidemiology*, 1995; **141**:943–950.
44. Kirkpatrick CS, White E, Lee JA. Case-control study of malignant melanoma in Washington State. II. Diet, alcohol, and obesity. *American Journal of Epidemiology*, 1994; **139**:869–880.
45. Osterlind A, Tucker MA, Stone BJ, Jensen OM. The Danish case-control study of cutaneous malignant melanoma. III. Hormonal and reproductive factors in women. *International Journal of Cancer*, 1988; **42**: 821–824.
46. Swerdlow AJ, English J, MacKie RM, O'Doherty CJ, Hunter JA, Clark J. Benign melanocytic naevi as a risk factor for malignant melanoma. *British Medical Journal (Clinical Research Education)*, 1986; **292**:1555–1559.
47. Zanetti R, Franceschi S, Rosso S, Bidoli E, Colonna S. Cutaneous malignant melanoma in females: the role of hormonal and reproductive factors. *International Journal of Epidemiology*, 1990; **19**:522–526.
48. Olkin I. Diagnostic statistical procedures in medical meta-analyses. *Statistics in Medicine* 1999; **18**:2331–2341.
49. Prentice RL, Pyke R. Logistic disease incidence models and case-control studies. *Biometrika*, 1979; **66**:403–411.
50. Breslow N. Regression analysis of the log-odds ratio: a method for retrospective studies. *Biometrics*, 1976; **32**:409–416.
51. Demidenko E. Asymptotic properties of nonlinear mixed-effects models. In *Modelling Longitudinal and Spatially Correlated Data: Methods, Applications, and Future Directions*. Springer Lecture Notes in Statistics, Springer-Verlag: New York, 1997; 49–62.
52. Morris CN. Parametric empirical Bayes inference: theory and applications. *Journal of the American Statistical Association*, 1983; **78**:47–55.
53. Carroll RJ, Ruppert D, Stefanski LA. *Measurement Error in Nonlinear Models*. Chapman and Hall: New York, 1995.
54. Farewell VT. Some results on the estimation of logistic models based on retrospective data. *Biometrika*, 1979; **66**:27–32.