

# A comparison of meta-analysis using literature and using individual patient data

**Thomas Mathew**

Department of Mathematics and Statistics  
University of Maryland  
Baltimore, Maryland 21250

mathew@umbc.edu

Joint work with Kenneth Nordström, University of Oulu, Finland

**Example 1** (Taken from Whitehead, 2003)

A study comparing two anaesthetics A and B with respect to the recovery times of patients undergoing short surgical procedures.

Data available from nine centers (log-transformed recovery times)

$\mu_A, \mu_B$ : population mean log-recovery times for anaesthetics A and B, respectively.

To estimate  $\mu_A - \mu_B$ .

Center (Trial)	Anaesthetic A			Anaesthetic B		
	# patients	Mean	SD	# patients	Mean	SD
1	4	1.141	0.967	5	0.277	0.620
2	10	2.165	0.269	10	1.519	0.913
3	17	1.790	0.795	17	1.518	0.849
4	8	2.105	0.387	9	1.189	1.061
5	7	1.324	0.470	10	0.456	0.619
6	11	2.369	0.401	10	1.550	0.558
7	10	1.074	0.670	12	0.265	0.502
8	5	2.583	0.409	4	1.370	0.934
9	14	1.844	0.848	19	2.118	0.749

$\bar{y}_{A_j}, \bar{y}_{B_j}$ : sample mean log-recovery times for Anaesthetic A and Anaesthetic B for the  $j$ th center.

$n_{A_j}, n_{B_j}$ : corresponding sample sizes.

$\sigma_j^2$ : Variability in the  $j$ th center. Can be estimated by pooling the pair of sample variances for each center.

$\hat{\sigma}_j^2$ : the estimator so obtained.

$\bar{y}_{A_j} - \bar{y}_{B_j}$ : estimator of  $\mu_A - \mu_B$  from the  $j$ th center.

$\hat{\sigma}_j^2 \left( \frac{1}{n_{A_j}} + \frac{1}{n_{B_j}} \right)$ : estimated variance of  $\bar{y}_{A_j} - \bar{y}_{B_j}$ .

Meta-analysis estimator of  $\mu_A - \mu_B$  based on summary data: a weighted combination of  $\bar{y}_{A_j} - \bar{y}_{B_j}$ .

The estimator has value 0.627, with standard error 0.0990.

Meta-analysis estimator of  $\mu_A - \mu_B$  based on the individual patient data (IPD):

a weighted combination of  $\bar{y}_{A_j}$  – a weighted combination of  $\bar{y}_{B_j}$

The estimator has value 0.679, with standard error 0.0982.

95% confidence intervals for  $\mu_A - \mu_B$ :

The interval (0.433, 0.821) based on summary data.

The interval (0.486, 0.874) based on the IPD.

**Example 2** (Taken from Bower et al., 2003)

Data on the costs of counseling in primary care.

Data from different studies (trials) are available.

Short term costs for patients treated by counselors, and for those who remained under the care of a general practitioner.

$\mu_1, \mu_2$ : the population average cost for patients treated by counselors, and for patients who remained under the care of a general practitioner, respectively.

To estimate  $\mu_1 - \mu_2$ .

Study (Trial)	Treated by counselor			Care by general practitioner		
	# patients	Mean	SD	# patients	Mean	SD
1	58	304	170	57	226	480
2	87	221	157	45	140	97
3	82	283	142	79	171	291
4	53	322	285	49	166	329

Meta-analysis estimate of  $\mu_1 - \mu_2$  based on summary data: 94.09, with standard error 17.468.

95% confidence intervals for  $\mu_1 - \mu_2$ : (59.86, 128.33)

Meta-analysis estimate of  $\mu_1 - \mu_2$  based on IPD: 117.050, with standard error 15.918

95% confidence intervals for  $\mu_1 - \mu_2$ : (85.85, 148.25)

The estimates and confidence intervals are quite different.

For testing a hypothesis concerning  $\mu_1 - \mu_2$ , the conclusion based on summary data, and that based on the IPD, can be different.



Angelillo and Villari (2003): “One issue that merits closer scrutiny is whether meta-analysis of published data is sufficient or whether individual patient data are necessary”.

Olkin, I. and Sampson, A. (1998). Comparison of meta-analysis versus analysis of variance of individual patient data. *Biometrics* **54**, 317–322.

Mathew, T. and Nordström, K. (1999). On the equivalence of meta-analysis using literature and using individual patient data. *Biometrics* **55**, 1221–1223.

The above articles consider a two-way model without interaction, with fixed treatment effects and fixed study effects.

**Conclusion in the above articles:** For estimating the treatment contrasts, meta-analysis estimators based on summary data, and that based on the IPD, are the same.

What if there are covariates?

What if the study effects are random?

What if there are no study effects?

## A general formulation

Consider  $k$  independent studies (or trials).

$\mathbf{y}_j$ : vector of  $n_j$  responses from the  $j$ th trial.

$\mathbf{y}_1, \dots, \mathbf{y}_k$ : individual patient data (IPD).

Assume a linear model:

$$E(\mathbf{y}_j) = W_j\boldsymbol{\beta} + Z_j\boldsymbol{\delta}_j = (W_j, Z_j) \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\delta}_j \end{pmatrix} = X_j \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\delta}_j \end{pmatrix}$$

$$\text{Cov}(\mathbf{y}_j) = V_j, \quad j = 1, \dots, k$$

$\boldsymbol{\beta}, \boldsymbol{\delta}_j$ : vectors of unknown parameters of dimension  $p$  and  $q_j$

$X_j = (W_j, Z_j)$ : design matrix

Assume  $V_j$  is completely known.

$$E(\mathbf{y}_j) = X_j \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\delta}_j \end{pmatrix}, \text{ Cov}(\mathbf{y}_j) = V_j, \text{ where } X_j = (W_j, Z_j).$$

$\boldsymbol{\theta} = L\boldsymbol{\beta}$ : parameter of interest.

Weighted least squares estimator of  $(\boldsymbol{\beta}', \boldsymbol{\delta}_j')$ , based on the data from the  $j$ th trial:

$$\begin{pmatrix} \hat{\boldsymbol{\beta}}^{(j)} \\ \hat{\boldsymbol{\delta}}_j \end{pmatrix} = (X_j' V_j^{-1} X_j)^{-1} X_j' V_j^{-1} \mathbf{y}_j.$$

Known to be the best linear unbiased estimator (BLUE).

Write

$$\mathcal{I}^{(j)} = X_j' V_j^{-1} X_j = \begin{pmatrix} W_j' V_j^{-1} W_j & W_j' V_j^{-1} Z_j \\ Z_j' V_j^{-1} W_j & Z_j' V_j^{-1} Z_j \end{pmatrix} = \begin{pmatrix} \mathcal{I}_{11}^{(j)} & \mathcal{I}_{12}^{(j)} \\ \mathcal{I}_{21}^{(j)} & \mathcal{I}_{22}^{(j)} \end{pmatrix}$$

Then

$$\text{Cov} \begin{pmatrix} \hat{\beta}^{(j)} \\ \hat{\delta}_j \end{pmatrix} = (X_j' V_j^{-1} X_j)^{-1} = \mathcal{I}^{(j)^{-1}}.$$

$$\text{Cov}(\hat{\beta}^{(j)}) = \mathcal{I}_{11.2}^{(j)^{-1}}, \quad \text{where } \mathcal{I}_{11.2}^{(j)} = \mathcal{I}_{11}^{(j)} - \mathcal{I}_{12}^{(j)} \mathcal{I}_{22}^{(j)^{-1}} \mathcal{I}_{21}^{(j)}$$

The BLUE of  $\theta = L\beta$  from the  $j$ th trial is  $\hat{\theta}^{(j)} = L\hat{\beta}^{(j)}$ .

$$\text{Cov}(\hat{\theta}^{(j)}) = \text{Cov}(L\hat{\beta}^{(j)}) = L\mathcal{I}_{11.2}^{(j)^{-1}} L'.$$

Summary data from the  $j$ th trial:  $\hat{\theta}^{(j)}$  and its covariance matrix  $L \mathcal{I}_{11.2}^{(j)-1} L'$ .

Meta-analysis estimator of  $\theta = L\beta$  based on summary data:

$$\tilde{\theta} = \left( \sum_{j=1}^k (L \mathcal{I}_{11.2}^{(j)-1} L')^{-1} \right)^{-1} \sum_{j=1}^k (L \mathcal{I}_{11.2}^{(j)-1} L')^{-1} \hat{\theta}^{(j)}$$

$$\text{Cov}(\tilde{\theta}) = \left( \sum_{j=1}^k (L \mathcal{I}_{11.2}^{(j)-1} L')^{-1} \right)^{-1}.$$

Model for the IPD:

$$\mathbf{y} = (\mathbf{y}'_1, \dots, \mathbf{y}'_k)',$$

$$W = (W'_1, \dots, W'_k)',$$

$$Z = \text{diag}(Z_1, \dots, Z_k),$$

$$\boldsymbol{\delta} = (\boldsymbol{\delta}'_1, \dots, \boldsymbol{\delta}'_k)',$$

$$V = \text{diag}(V_1, \dots, V_k).$$

$$E(\mathbf{y}) = W\boldsymbol{\beta} + Z\boldsymbol{\delta} = (W, Z) \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\delta} \end{pmatrix} = X \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\delta} \end{pmatrix}$$

$$\text{Cov}(\mathbf{y}) = V, \text{ where } X = (W, Z).$$

Meta-analysis estimator of  $\boldsymbol{\theta} = L\boldsymbol{\beta}$  based on IPD: BLUE of  $\boldsymbol{\theta}$  obtained from the above model.

$$E(\mathbf{y}) = X \begin{pmatrix} \beta \\ \delta \end{pmatrix}, \quad \text{Cov}(\mathbf{y}) = V, \quad \text{where } X = (W, Z).$$

$$\begin{pmatrix} \hat{\beta} \\ \hat{\delta} \end{pmatrix} = (X'V^{-1}X)^{-1}X'V^{-1}\mathbf{y}, \quad \text{Cov} \begin{pmatrix} \hat{\beta} \\ \hat{\delta} \end{pmatrix} = (X'V^{-1}X)^{-1}.$$

$$\hat{\theta} = L\hat{\beta} = (L, 0) \begin{pmatrix} \hat{\beta} \\ \hat{\delta} \end{pmatrix} = (L, 0)(X'V^{-1}X)^{-1}X'V^{-1}\mathbf{y}$$

$$\text{Cov}(\hat{\theta}) = (L, 0)(X'V^{-1}X)^{-1}(L, 0)'$$



Recall:  $X_j = (W_j, Z_j)$  and

$$\mathcal{I}^{(j)} = X_j' V_j^{-1} X_j = \begin{pmatrix} W_j' V_j^{-1} W_j & W_j' V_j^{-1} Z_j \\ Z_j' V_j^{-1} W_j & Z_j' V_j^{-1} Z_j \end{pmatrix} = \begin{pmatrix} \mathcal{I}_{11}^{(j)} & \mathcal{I}_{12}^{(j)} \\ \mathcal{I}_{21}^{(j)} & \mathcal{I}_{22}^{(j)} \end{pmatrix}$$

$$\mathcal{I}_{11.2}^{(j)} = \mathcal{I}_{11}^{(j)} - \mathcal{I}_{12}^{(j)} \mathcal{I}_{22}^{(j)-1} \mathcal{I}_{21}^{(j)}$$

$$X' V^{-1} X = \sum_{j=1}^k (X_j' V_j^{-1} X_j) = \sum_{j=1}^k \mathcal{I}^{(j)}$$

$$\text{Cov}(\hat{\theta}) = (L, 0) (X' V^{-1} X)^{-1} (L, 0)'$$

$$= L \left( \sum_{j=1}^k \mathcal{I}_{11.2}^{(j)} \right)^{-1} L'$$

Based on summary data:

$$\tilde{\boldsymbol{\theta}} = \left( \sum_{j=1}^k (L \mathcal{I}_{11.2}^{(j)} L')^{-1} \right)^{-1} \sum_{j=1}^k (L \mathcal{I}_{11.2}^{(j)} L')^{-1} \hat{\boldsymbol{\theta}}^{(j)}$$

$$\text{Cov}(\tilde{\boldsymbol{\theta}}) = \left( \sum_{j=1}^k (L \mathcal{I}_{11.2}^{(j)} L')^{-1} \right)^{-1} .$$

Based on IPD:

$$\hat{\boldsymbol{\theta}} = (L, 0)(X'V^{-1}X)^{-1}X'V^{-1}\mathbf{y}$$

$$\text{Cov}(\hat{\boldsymbol{\theta}}) = L \left( \sum_{j=1}^k \mathcal{I}_{11.2}^{(j)} \right)^{-1} L'$$

For  $\theta = L\beta$ ,  $\hat{\theta}$  is the BLUE based on the linear model for the IPD.

$\tilde{\theta}$  is another unbiased estimator.

Hence  $\text{Cov}(\tilde{\theta}) - \text{Cov}(\hat{\theta})$  is nonnegative definite.

That is

$$\left( \sum_{j=1}^k (L \mathcal{I}_{11 \cdot 2}^{(j)} L')^{-1} \right)^{-1} - L \left( \sum_{j=1}^k \mathcal{I}_{11 \cdot 2}^{(j)} \right)^{-1} L'$$

is nonnegative definite.

Also follows from the concavity of the matrix function  $g(A) = (LA^{-1}L')^{-1}$ , for  $A$  p.d. (Marshall and Olkin, 1979, p. 469)

If the two matrices are equal, then the two estimators coincide, and there is no loss of information in using the summary data.

**Result:** The meta-analysis estimator based on the summary data, and that based on the IPD coincide if and only if the matrices  $(L \mathcal{I}_{11.2}^{(j)-1} L')^{-1} L \mathcal{I}_{11.2}^{(j)-1}$  are equal for  $j = 1, \dots, k$ .

The equality of the two estimators require a certain level of homogeneity across the trials.

The two estimators are equal if  $L = I$ . That is, if there are no common nuisance parameters across studies.

## Application to some special models

**To estimate a single treatment-control difference  $\mu_1 - \mu_2$ :**

Assume there are no other effects or covariates.

$\bar{y}_{1j}, \bar{y}_{2j}$ : sample means for the treatment and the control from the  $j$ th trial.

$n_{1j}, n_{2j}$ : corresponding sample sizes.

$\sigma_j^2$ : variability for the  $j$ th trial (usually an estimate).

Parameters in the mean:  $(\mu_1, \mu_2)$  or  $(\mu_1 - \mu_2, \mu_2)$

$\mu_2$  is a common nuisance parameter across the trials.

The two estimators of  $\mu_1 - \mu_2$  coincide if and only if  $\frac{n_{1j}}{n_{1j} + n_{2j}}$  are all the same for  $j = 1, 2, \dots, k$ .

The fraction of observations for the treatment is the same across the different trials.

The condition is free of the variances!

## Example 1 (continued)

A study comparing two anaesthetics A and B with respect to the recovery times of patients undergoing short surgical procedures.

Log-transformed recovery times available from nine centers.

$\mu_A, \mu_B$ : population mean log-recovery times for anaesthetics A and B, respectively.

To estimate  $\mu_A - \mu_B$ .

Center (Trial)	Anaesthetic A			Anaesthetic B		
	# patients	Mean	SD	# patients	Mean	SD
1	4	1.141	0.967	5	0.277	0.620
2	10	2.165	0.269	10	1.519	0.913
3	17	1.790	0.795	17	1.518	0.849
4	8	2.105	0.387	9	1.189	1.061
5	7	1.324	0.470	10	0.456	0.619
6	11	2.369	0.401	10	1.550	0.558
7	10	1.074	0.670	12	0.265	0.502
8	5	2.583	0.409	4	1.370	0.934
9	14	1.844	0.848	19	2.118	0.749



$\sigma_j^2$ : Variability in the  $j$ th center. Can be estimated by pooling the pair of sample variances for each center.

Based on summary data, the meta-analysis estimator of  $\mu_A - \mu_B$  has value 0.627, with standard error 0.099.

Based on IPD, the meta-analysis estimator of  $\mu_A - \mu_B$  has value 0.679, with standard error 0.0982.

95% confidence intervals for  $\mu_A - \mu_B$ :

The interval (0.433, 0.821) based on summary data.

The interval (0.486, 0.874) based on the IPD.

## Example 2 (continued)

Data on the costs of counseling in primary care available from four studies.

Short term costs for patients treated by counselors, and for those who remained under the care of a general practitioner.

$\mu_1, \mu_2$ : population average cost for patients treated by counselors, and for patients who remained under the care of a general practitioner, respectively.

To estimate  $\mu_1 - \mu_2$ .

Study (Trial)	Treated by counselor			Care by general practitioner		
	# patients	Mean	SD	# patients	Mean	SD
1	58	304	170	57	226	480
2	87	221	157	45	140	97
3	82	283	142	79	171	291
4	53	322	285	49	166	329

Meta-analysis estimator of  $\mu_1 - \mu_2$  based on summary data: 94.09, with standard error 17.468.

95% confidence intervals for  $\mu_1 - \mu_2$ : (59.86, 128.33)

Meta-analysis estimator of  $\mu_1 - \mu_2$  based on IPD: 117.050, with standard error 15.918

95% confidence intervals for  $\mu_1 - \mu_2$ : (85.85, 148.25)

## To estimate several treatment-control differences:

Suppose several treatments are to be compared to the same control, **in the absence of any other effects or covariates**, based on data from  $k$  trials.

$\mu_i$ : mean for the  $i$ th treatment ( $i = 1, 2, \dots, m - 1$ ),

$\mu_m$ : mean for the control.

To estimate the treatment-control differences  $\mu_i - \mu_m$ ,  $i = 1, 2, \dots, m - 1$ .

$(\mu_1, \mu_2, \dots, \mu_m) = (\mu_1 - \mu_m, \mu_2 - \mu_m, \dots, \mu_{m-1} - \mu_m, \mu_m)$ .

$\mu_m$ : A common nuisance parameter across the different studies.

Meta-analysis estimator based on summary data, and the meta-analysis estimator based on IPD coincide if and only if the fraction of observations for the  $i$ th treatment is the same across the different trials;  $i = 1, 2, \dots, m$ .

## A model with fixed treatment effects and fixed trial effects:

$y_{ij\alpha}$ : the  $\alpha$ th response on the  $i$ th treatment from the  $j$ th trial,  
 $\alpha = 1, \dots, n_{ij}$ ,  $i = 1, 2, \dots, m$ ,  $j = 1, 2, \dots, k$ . Assume  $n_{ij} > 0$   
for all  $i$  and  $j$ .

$\tau_i$ :  $i$ th treatment effect

$\alpha_j$ :  $j$ th trial effect.

Assume

$$E(y_{ij\alpha}) = \alpha_j + \tau_i.$$

$\mathbf{y}_{ij}$ .: vector of  $n_{ij}$  responses on the  $i$ th treatment in the  $j$ th trial

$$\mathbf{y}_{\cdot j} = (\mathbf{y}'_{1j}, \dots, \mathbf{y}'_{mj})'$$

Assume  $\text{Cov}(\mathbf{y}_{\cdot j}) = V_j$  (known)

Let the  $m$ th treatment represent a control

To estimate  $\boldsymbol{\theta} = (\tau_1 - \tau_m, \dots, \tau_{m-1} - \tau_m)'$

Write  $\theta_i = \tau_i - \tau_m$ ,  $i = 1, 2, \dots, m - 1$ , so that  $\tau_i = \theta_i + \tau_m$ .

$$E(y_{ij\alpha}) = \alpha_j + \tau_i = \alpha_j + \theta_i + \tau_m = \alpha'_j + \theta_i.$$

To estimate  $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_{m-1})'$ .

The nuisance parameters are  $\alpha'_j$ :  $j = 1, 2, \dots, k$ .

No common nuisance parameters across the studies.

Meta-analysis estimator based on summary data, and the meta-analysis estimator based on IPD always coincide.

Result due to Olkin and Sampson (1998) and Mathew and Nordström (1999).



## A model with fixed treatment effects and random trial effects:

Now the trial effects will contribute to the covariance matrix.

$$E(y_{ij\alpha}) = \mu_i, \quad i = 1, 2, \dots, m.$$

To estimate  $\theta = (\mu_1 - \mu_m, \dots, \mu_{m-1} - \mu_m)'$

$$(\mu_1, \mu_2, \dots, \mu_m)' = (\mu_1 - \mu_m, \dots, \mu_{m-1} - \mu_m, \mu_m)'$$

$\mu_m$  is a common nuisance parameter across the trials.

The two estimators of  $\theta$  coincide if and only if the fraction of observations for the  $i$ th treatment is the same across the different trials;  $i = 1, 2, \dots, m$

## What if there are covariates?

Consider a simple model with one covariate, and without treatment-covariate interaction.

The problem is that of estimating a treatment–control difference.

$y_{ij}$ : outcome for the  $i$ th patient in the  $j$ th trial.

Consider the model

$$y_{ij} = \beta_{0j} + \beta_1 x_{1ij} + \beta_2 x_{2ij} + \epsilon_{ij},$$

$i = 1, \dots, n_j$  and  $j = 1, \dots, k$  (Higgins et al., 2001, Section 3.1)

$\beta_{0j}$ : trial effects, assumed to be fixed

$x_{1ij}$ : a dummy variable that indicates treatment group (treatment or control)

$x_{2ij}$ : a single covariate

$\beta_1$ : treatment–control difference

$\beta_2$ : regression coefficient.

$\epsilon_{ij}$ 's are independent and identically distributed random variables with mean zero and variance  $\sigma^2$ .

To estimate  $\beta_1$ .

Meta-analysis estimator of  $\beta_1$  based on summary data: Estimate  $\beta_1$  and its variance from each trial, and combine the estimators.

Meta-analysis estimator of  $\beta_1$  based on IPD: Estimate  $\beta_1$  using the model for the entire data, consisting of all the  $y_{ij}$ 's.

Define

$$\bar{x}_{1j} = \frac{1}{n_j} \sum_{i=1}^{n_j} x_{1ij}, \quad \bar{x}_{2j} = \frac{1}{n_j} \sum_{i=1}^{n_j} x_{2ij}.$$

The two estimators coincide if the quantities

$$\sum_{i=1}^{n_j} (x_{1ij} - \bar{x}_{1j})(x_{2ij} - \bar{x}_{2j}) / \sum_{i=1}^{n_j} (x_{2ij} - \bar{x}_{2j})^2$$

are all equal, for all  $j = 1, \dots, k$ .

Requires a certain level of homogeneity with respect to the treatment allocation and the covariates across trials.

Patient-level homogeneity of the covariate within each trial implies the required condition.

Unrealistic in practice.

Thus equality of the two estimators is unlikely to hold when covariates are present.

## References

Mathew, T. and Nordström, K. (2010). Comparison of one-step and two-step meta-analysis models using individual patient data. *Biometrical Journal* **52**, 271-287.

Lin, D. Y. and Zeng, D. (2010). On the relative efficiency of using summary statistics versus individual-level data in meta-analysis. *Biometrika* **97**, 321-332.

Angelillo, I. F. and Villari, P. (2003). Meta-analysis of published studies or meta-analysis of individual data? Caesarean section in HIV-positive women as a case study. *Public Health* **117**, 323328.

Bower, P., Byford, S., Barber, J., Beecham, J., Simpson, S., Friedli, K., Corney, R., King, M. and Harvey, I. (2003). Meta-analysis of data on costs from trials of counseling in primary care: using individual patient data to overcome sample size limitations in economic analyses. *British Medical Journal* **326**, 1247–1250.

Higgins, J. P. T., Whitehead, A., Turner, R. M., Omar, R. Z. and Thompson, S. G. (2001). Meta-analysis of continuous outcome data from individual patients. *Statistics in Medicine* **20**, 2219–2241.

Marshall, A. W. and Olkin, I. (1979). *Inequalities: Theory of Majorization and Its Applications*. Academic Press, New York.

Mathew, T. and Nordström, K. (1999). On the equivalence of meta-analysis using literature and using individual patient data. *Biometrics* **55**, 1221–1223.

Olkin, I. and Sampson, A. (1998). Comparison of meta-analysis versus analysis of variance of individual patient data. *Biometrics* **54**, 317–322.

Whitehead, A. (2003). *Meta-Analysis of Controlled Clinical Trials*. Wiley, New York.