

Maximum Likelihood solution via EM - random intercepts model

- E-step (expectation - “Expected A Posteriori” or Empirical Bayes)

$$\tilde{v}_i = \rho_{n_i n_i} \left[\frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij} - \mathbf{x}'_{ij} \boldsymbol{\beta} \right]$$

$$\sigma_{v|y_i}^2 = \sigma_v^2 (1 - \rho_{n_i n_i}) \quad \text{where } \rho_{n_i n_i} = \frac{n_i r}{1 + (n_i - 1)r} \quad \text{and} \quad r = \frac{\sigma_v^2}{\sigma_v^2 + \sigma^2}$$

- M-step (maximization - “Maximum Likelihood”)

$$\hat{\boldsymbol{\beta}} = \left(\sum_i^N \mathbf{X}'_i \mathbf{X}_i \right)^{-1} \sum_i^N \mathbf{X}'_i (\mathbf{y}_i - \mathbf{1}_i \tilde{v}_i)$$

$$\hat{\sigma}_v^2 = \frac{1}{N} \sum_i^N \tilde{v}_i^2 + \sigma_{v|y_i}^2$$

$$\hat{\sigma}^2 = \left(\sum_i^N n_i \right)^{-1} \sum_i^N (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}} - \mathbf{1}_i \tilde{v}_i)' (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}} - \mathbf{1}_i \tilde{v}_i) + n_i \sigma_{v|y_i}^2$$

- provide starting values for $\boldsymbol{\beta}$, σ_v^2 , and σ^2
- perform E-step, perform M-step, repeat early and often (until convergence)

Wald z or χ^2 tests

$$z = \frac{\hat{\theta}}{se(\hat{\theta})}$$

is approximately normally distributed under the large sample theory of ML estimates

- for single parameter tests
 - *e.g.*, $H_0 : \beta_1 = 0$
- asymptotic test
- not generally recommended for testing variance terms
- can be generalized for multi-parameter tests

Note:

- $z^2 \sim \chi_1^2$, so either z or χ^2 representation
- since asymptotic, t -tests based on sample size have been proposed (in MIXED, the DDFM=KENWARDROGER option on the MODEL statement performs the degrees-of-freedom calculations detailed by Kenward and Roger, Biometrics, 1997)
- this approach also used to test hypotheses regarding random effects \mathbf{v}_i (however, multiple tests issue)
- asymptotic confidence intervals for θ formed as

$$\hat{\theta} \pm z_{1-\alpha/2} \times se(\hat{\theta})$$

Likelihood-ratio χ^2 test

Suppose Model I is nested within Model II, then

$$2 \times \log(L_{\text{II}} / L_{\text{I}}) = 2 \times (\log L_{\text{II}} - \log L_{\text{I}}) \sim \chi_q^2$$

where q = number of additional terms in Model II

$-2 \log L$ is called the *deviance*, $D_{\text{I}} - D_{\text{II}} \sim \chi_q^2$

- for single or multiparameter tests
e.g., $H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$
or $H_0 : \sigma_{v_0}^2 = \sigma_{v_1}^2 = \sigma_{v_0 v_1} = 0$ (divide p -value by 2)
- good for comparing models *with same sample size*
(be careful with missing values on \mathbf{X} variables)
- With REML estimation - *cannot* use for testing fixed effects
(*i.e.*, $\boldsymbol{\beta}$)

Maximum Likelihood Estimation - R.A. Fisher (1922, Foundations of Theoretical Statistics, *Philos. Trans. Royal Society of London, A*, 309-368)

“The object of statistical methods is the reduction of data. A quantity of data, which usually by its mere bulk is incapable of entering the mind, is to be replaced by relatively few quantities which shall adequately represent the whole, or which, in other words, shall contain as much as possible, ideally the whole, of the relevant information contained in the original data.”

This object is accomplished by:

- constructing a hypothetical infinite population
- data regarded as constituting a random sample (from this population)
- distribution of hypothetical population is specified by relatively few parameters
- information from sample is used to estimate these parameters

The idea of maximum likelihood is to obtain estimates for these parameters which are most likely given the observed data

Sherlock Holmes as an early proponent of ML

We now join Dr. Watson and Sherlock Holmes in their room at No. 221B, Baker Street, in the opening scene of “The Sign of Four” already in progress

(Watson remarks) “ . . . you spoke just now of observation and deduction. Surely the one to some extent implies the other.”

“Why, hardly,” he answered, leaning back luxuriously in his armchair and sending up thick blue wreaths from his pipe. “For example, observation shows me that you have been to the Wigmore Street Post-Office this morning, but deduction lets me know that when there you dispatched a telegram.”

“Right!” said I. “Right on both points! But I confess that I don’t see how you arrived at it. It was a sudden impulse upon my part, and I have mentioned it to no one.”

“It is simplicity itself,” he remarked, chuckling at my surprise,

“so absurdly simple an explanation is superfluous; and yet it may serve to define the limits of observation and of deduction. Observation tells me that you have a little reddish mould adhering to your instep. Just opposite the Wigmore Street Office they have taken up the pavement and thrown up some earth, which lies in such a way that it is difficult to avoid treading in it in entering. The earth is of this peculiar reddish tint which is found, as far as I know, nowhere else in the neighbourhood. So much is observation. The rest is deduction.”

“How, then, did you deduce the telegram?”

“Why, of course, I knew that you had not written a letter, since I sat opposite to you all morning. I see also in your open desk there that you have a sheet of stamps and a thick bundle of postcards. What could you go into the post-office for, then, but to send a wire? *Eliminate all other factors, and the one which remains must be the truth.*”

For Holmes, deducing (estimating) where Watson went that morning (parameter of interest),

observation = data - red mould on shoe, red earth by post-office

deduction = parameter estimate - that Watson went to the post-office

Deducing (estimating) what Watson did there (other parameter of interest),

observation = data - Watson did not write a letter, Watson has stamps and postcards

deduction = parameter estimate - Watson sent a wire



*⇒ Method of Deduction = Maximum Likelihood since
Holmes' recommendation is to choose the option (parameter)
with the highest likelihood, given the data*

Likelihood Function $p(\mathbf{y}, \boldsymbol{\theta})$ or $p(\mathbf{y} | \boldsymbol{\theta})$

- \mathbf{y} represents the observed data
- $\boldsymbol{\theta}$ represents the parameters which characterize hypothetical infinite population, $\boldsymbol{\theta}$ belonging to a set of possible values from Θ
- the likelihood function $p(\mathbf{y}, \boldsymbol{\theta})$ is the probability density at \mathbf{y} when $\boldsymbol{\theta}$ are the true parameters; expresses plausibility of the parameters after having observed the data
- ML estimates are obtained by maximizing the likelihood function
- obtain $\hat{\boldsymbol{\theta}}$ with the highest plausibility given the data

Definition in Silvey (1975): MLE $\hat{\boldsymbol{\theta}}(\mathbf{y})$ is any element of Θ such that $p\{\mathbf{y}, \hat{\boldsymbol{\theta}}(\mathbf{y})\} = \max_{\boldsymbol{\theta} \in \Theta} p(\mathbf{y}, \boldsymbol{\theta})$

Example: Linear Regression Model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i$$

$i = 1 \dots N$ subjects

or in matrix form,

$$\begin{array}{ccccccc} \mathbf{y} & = & \mathbf{X} & \boldsymbol{\beta} & + & \boldsymbol{\varepsilon} & \text{with } \varepsilon_i \sim \mathcal{NID}(0, \sigma^2) \\ N \times 1 & & N \times p & p \times 1 & & N \times 1 & \end{array}$$

$$p(\varepsilon_i | \sigma^2) = (2\pi\sigma^2)^{-\frac{1}{2}} \exp \left[-\frac{1}{2\sigma^2} \varepsilon_i^2 \right]$$

as a result, $y_i \sim \mathcal{NID}(\mathbf{x}_i' \boldsymbol{\beta}, \sigma^2)$ and the density of y_i is

$$p(y_i | \boldsymbol{\beta}, \sigma^2) = (2\pi\sigma^2)^{-\frac{1}{2}} \exp \left[-\frac{1}{2\sigma^2} (y_i - \mathbf{x}_i' \boldsymbol{\beta})^2 \right]$$

Since ε_i are independent

$$p(\boldsymbol{\varepsilon} \mid \sigma^2) = (2\pi\sigma^2)^{-\frac{N}{2}} \exp\left[-\frac{1}{2\sigma^2}\boldsymbol{\varepsilon}'\boldsymbol{\varepsilon}\right]$$

and the density of \mathbf{y} is

$$p(\mathbf{y} \mid \boldsymbol{\beta}, \sigma^2) = (2\pi\sigma^2)^{-\frac{N}{2}} \exp\left[-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right]$$

Need to maximize $p(\mathbf{y} \mid \boldsymbol{\beta}, \sigma^2)$ with respect to $\boldsymbol{\beta}$ and σ^2 to give the most plausible parameter values, given the data

$p(\mathbf{y} \mid \boldsymbol{\theta})$ is at a maximum (for $\boldsymbol{\theta} \in \Theta$) when

- derivative $p'(\mathbf{y} \mid \boldsymbol{\theta}) = 0$
- second derivative $p''(\mathbf{y} \mid \boldsymbol{\theta}) < 0$

The likelihood function is given as:

$$p(\mathbf{y} \mid \boldsymbol{\beta}, \sigma^2) = (2\pi\sigma^2)^{-\frac{N}{2}} \exp \left[-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right]$$

maximizing the log-likelihood gives same result as maximizing the likelihood, however,

- logs are our friends & can be natural
- easier to differentiate the log-likelihood

$$\begin{aligned} \log L &= -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2}\sigma^{-2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &= -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma^2) - \frac{1}{2}\sigma^{-2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &= C - \frac{N}{2} \log(\sigma^2) - \frac{1}{2}\sigma^{-2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \end{aligned}$$

Differentiating with respect to σ^2

$$\begin{aligned}\frac{\partial \log L}{\partial \sigma^2} &= -\frac{N}{2}\sigma^{-2} + \frac{1}{2}\sigma^{-4}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &= \frac{1}{2}\sigma^{-2} \left[\sigma^{-2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) - N \right]\end{aligned}$$

and setting equal to zero,

$$\hat{\sigma}^2 = \frac{1}{N}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

Notice that the above ML estimator yields a biased estimate when sample size is small and/or p is large:

$$\hat{\sigma}_{OLS}^2 = \frac{1}{N-p}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

$$\begin{aligned}
\left(\frac{\partial \log L}{\partial \sigma^2}\right)^2 &= \frac{N}{2}\sigma^{-4} - \sigma^{-6}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\
&= \sigma^{-4} \left[\frac{N}{2} - \sigma^{-2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right]
\end{aligned}$$

and for this to be a maximum,

$$\begin{aligned}
0 &> \frac{N}{2} - \sigma^{-2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\
&> \sigma^2 - \frac{2}{N}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\
&> \sigma^2 - 2\sigma^2 \\
&> -\sigma^2
\end{aligned}$$

Similarly, differentiating

$$\log L = C - \frac{N}{2} \log(\sigma^2) - \frac{1}{2} \sigma^{-2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

with respect to $\boldsymbol{\beta}$, yields

$$\frac{\partial \log L}{\partial \boldsymbol{\beta}} = -\frac{1}{2\sigma^2} \left\{ \underbrace{(-\mathbf{X})'}_{p \times N} \underbrace{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}_{N \times 1} + \underbrace{[(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})']}_{1 \times N} \underbrace{(-\mathbf{X})}_{1 \times p} \right\}$$

and since $(\mathbf{AB})' = \mathbf{B}'\mathbf{A}'$

$$\begin{aligned} &= \frac{1}{2} \sigma^{-2} [\mathbf{X}'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \mathbf{X}'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})] \\ &= \sigma^{-2} \mathbf{X}'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \end{aligned}$$

setting equal to zero,

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$$

Notice

$$\left(\frac{\partial \log L}{\partial \boldsymbol{\beta}} \right)^2 = -(\mathbf{X}'\mathbf{X})$$

which always contains diagonal elements < 0

The mixed-effects regression model for unit i :

$$\mathbf{y}_i = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \mathbf{v}_i + \boldsymbol{\varepsilon}_i$$

$$\begin{matrix} n_i \times 1 & n_i \times p & p \times 1 & n_i \times r & r \times 1 & n_i \times 1 \end{matrix}$$

$$i = 1 \dots N \text{ units ; } \quad j = 1 \dots n_i \text{ observations within unit } i$$

with $\boldsymbol{\varepsilon} \sim \mathcal{NID}(0, \sigma^2 \mathbf{I}_{n_i})$ and $\mathbf{v} \sim \mathcal{NID}(0, \boldsymbol{\Sigma}_v)$

As a result, the observations \mathbf{y} and random coefficients \mathbf{v} have the joint multivariate normal distribution:

$$\begin{bmatrix} \mathbf{y}_i \\ \mathbf{v}_i \end{bmatrix} \sim \mathcal{NID} \left(\begin{bmatrix} \mathbf{X}_i \boldsymbol{\beta} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{Z}_i \boldsymbol{\Sigma}_v \mathbf{Z}_i' + \sigma^2 \mathbf{I}_{n_i} & \mathbf{Z}_i \boldsymbol{\Sigma}_v \\ \boldsymbol{\Sigma}_v \mathbf{Z}_i' & \boldsymbol{\Sigma}_v \end{bmatrix} \right)$$

Parameter estimation

- Empirical Bayes (EB) for random-effects \mathbf{v}_i
- Maximum (marginal) likelihood (ML) for variance parameters, σ^2 and $\boldsymbol{\Sigma}_v$, and covariate effects $\boldsymbol{\beta}$

Articles about Estimation

Bock, R.D. (1989). Measurement of human variation: A two-stage model. In R.D. Bock (Ed.), *Multilevel analysis of educational data*. New York: Academic Press.

Dempster, A. P., Rubin, D. B., & Tsutakawa, R. K. (1981). Estimation in covariance components models. *Journal of the American Statistical Association*, **76**: 341-353.

Laird, N. M. & Ware, J. H. (1982). Random effects models for longitudinal data. *Biometrics*, **38**: 963 - 974.

Empirical Bayes Estimation

Morris, C.N. (1983) Parametric Empirical Bayes Inference: Theory and Applications. *JASA*, **78**:47-55.

EB inference concerns two stochastic processes

- one for the data $f(y_i | \mathbf{v}_i; \boldsymbol{\zeta})$
 - f is the likelihood function (probability density)
 - describes probability of the data associated with the random effects \mathbf{v}_i and parameters $\boldsymbol{\zeta}$ ($\boldsymbol{\beta}$ and σ^2)
- one for the random effects $g(\mathbf{v}; \boldsymbol{\eta})$
 - g is the prior probability density
 - $\boldsymbol{\eta}$ characterizes the distribution ($\boldsymbol{\Sigma}_v$)

All of the information about \mathbf{v}_i that is available in y_i is conveyed by the posterior probability or probability density of \mathbf{v} , given y_i , as expressed in Bayes Theorem:

$$p(\mathbf{v} \mid y_i) = \frac{f(y_i \mid \mathbf{v}_i; \boldsymbol{\zeta}) g(\mathbf{v}; \boldsymbol{\eta})}{h(y_i)}$$

where

$$h(y_i) = \int_{\mathbf{v}} f(y_i \mid \mathbf{v}_i; \boldsymbol{\zeta}) g(\mathbf{v}; \boldsymbol{\eta}) d\mathbf{v}$$

is the marginal probability of the observation

The posterior has all the information about \mathbf{v}_i , but what information do you want from the posterior?

- the mean or expected value $\tilde{\mathbf{v}}_i$

$$\tilde{\mathbf{v}}_i = \int_{\mathbf{v}} \mathbf{v} p(\mathbf{v} | y_i) d\mathbf{v}$$

- the variance-covariance matrix $\Sigma_{\mathbf{v}|y_i}$

$$\Sigma_{\mathbf{v}|y_i} = \int_{\mathbf{v}} (\mathbf{v} - \tilde{\mathbf{v}}_i)(\mathbf{v} - \tilde{\mathbf{v}}_i)' p(\mathbf{v} | y_i) d\mathbf{v}$$

How do we get these quantities?

- posterior is conditional distribution of \mathbf{v} given \mathbf{y}_i
- \mathbf{v} and \mathbf{y}_i jointly multivariate normally distributed
- use property of multivariate normal distribution (Rao, Bock, *etc.*), conditional distribution of \mathbf{a}_2 given \mathbf{a}_1 is $\mathcal{NID}(\boldsymbol{\mu}_2 + \Sigma_{21}\Sigma_{11}^{-1}(\mathbf{a}_1 - \boldsymbol{\mu}_1), \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12})$

Using this property, the posterior is given by

$$\mathbf{v} \mid \mathbf{y}_i \sim \mathcal{NID} \left[\boldsymbol{\Sigma}_v \mathbf{Z}'_i (\mathbf{Z}_i \boldsymbol{\Sigma}_v \mathbf{Z}'_i + \sigma^2 \mathbf{I}_{n_i})^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) , \right. \\ \left. \boldsymbol{\Sigma}_v - \boldsymbol{\Sigma}_v \mathbf{Z}'_i (\mathbf{Z}_i \boldsymbol{\Sigma}_v \mathbf{Z}'_i + \sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i \boldsymbol{\Sigma}_v \right]$$

The posterior mean provides the EB estimate of \mathbf{v}_i

The posterior variance-covariance matrix provides the EB estimate of uncertainty about \mathbf{v}_i

however, the forms above require inversion of $n_i \times n_i$ matrices

Using the really cool property (Dempster *et al.*, 81; Bock, 83)

$$\boldsymbol{\Sigma}_v \mathbf{Z}'_i (\mathbf{Z}_i \boldsymbol{\Sigma}_v \mathbf{Z}'_i + \sigma^2 \mathbf{I}_{n_i})^{-1} = (\mathbf{Z}'_i (\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i + \boldsymbol{\Sigma}_v^{-1})^{-1} \mathbf{Z}'_i (\sigma^2 \mathbf{I}_{n_i})^{-1}$$

which only requires inversion of $r \times r$ matrices

Now, the mean of the posterior distribution of \mathbf{v} , given \mathbf{y}_i , is

$$\begin{aligned} \tilde{\mathbf{v}}_i &= (\mathbf{Z}'_i (\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i + \boldsymbol{\Sigma}_v^{-1})^{-1} \mathbf{Z}'_i (\sigma^2 \mathbf{I}_{n_i})^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \\ &= (\mathbf{Z}'_i \mathbf{Z}_i + \sigma^2 \boldsymbol{\Sigma}_v^{-1})^{-1} \mathbf{Z}'_i (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \end{aligned}$$

with the variance-covariance matrix as

$$\boldsymbol{\Sigma}_{v|y_i} = (\mathbf{Z}'_i (\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i + \boldsymbol{\Sigma}_v^{-1})^{-1}$$

These are the Empirical Bayes or EAP estimates (“Expected A Posteriori”) of the random effects with corresponding posterior variance terms (sometimes PSDs are given instead - positive square roots of the diagonal elements of $\Sigma_{v|y_i}$)

From RIESBYM2.SAS, add **S** (or **SOLUTION**) on **RANDOM** statement:

```
PROC MIXED METHOD=ML;  
  CLASS id;  
  MODEL hamd = week / SOLUTION ;  
  RANDOM INT week / SUBJECT=id TYPE=UN G GCORR S;
```

Solution for Random Effects

Obs	Effect	id	Estimate	StdErr	Pred	DF	tValue	Probt
1	Intercept	101	1.0290	2.0432	243	0.50	0.6150	
2	week	101	-2.1035	0.6995	243	-3.01	0.0029	
3	Intercept	103	3.6392	2.0432	243	1.78	0.0761	
4	week	103	-0.4743	0.6995	243	-0.68	0.4983	
5	Intercept	104	2.6386	2.0432	243	1.29	0.1978	
6	week	104	-1.4891	0.6995	243	-2.13	0.0343	
.	

- **Estimate** column provides $\tilde{\mathbf{v}}_i$ (posterior means)
- **StdErr** **Pred** column provides square root of diagonal elements of $\Sigma_{v|y_i}$ (posterior std devs)

Suppose a random-intercepts models

$$\begin{aligned}\tilde{v}_i &= (\mathbf{1}'_i \mathbf{1}_i + \sigma^2 / \sigma_v^2)^{-1} \mathbf{1}'_i (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \\ &= (n_i + \sigma^2 / \sigma_v^2)^{-1} \mathbf{1}'_i (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \\ &= \left(\frac{n_i \sigma_v^2 + \sigma^2}{\sigma_v^2} \right)^{-1} \mathbf{1}'_i (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \\ &= \left(\frac{\sigma_v^2}{\sigma_v^2 + \sigma^2 / n_i} \right) \frac{1}{n_i} \mathbf{1}'_i (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \\ &= \left(\frac{n_i \sigma_v^2}{n_i \sigma_v^2 + \sigma^2} \right) \frac{1}{n_i} \mathbf{1}'_i (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})\end{aligned}$$

$$\begin{aligned}
&= \left(\frac{n_i \sigma_v^2}{\sigma_v^2 + \sigma^2} \right) \left(\frac{1}{1 + (n_i - 1) \frac{\sigma_v^2}{\sigma_v^2 + \sigma^2}} \right) \frac{1}{n_i} \mathbf{1}'_i (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \\
&= \left(\frac{n_i r}{1 + (n_i - 1) r} \right) \frac{1}{n_i} \mathbf{1}'_i (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \\
&= \rho_{n_i n_i} \frac{1}{n_i} \mathbf{1}'_i (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) = \rho_{n_i n_i} \frac{1}{n_i} \sum_j^{n_i} (y_{ij} - \mathbf{x}'_{ij} \boldsymbol{\beta}) \\
&= \rho_{n_i n_i} [\bar{\mathbf{y}}_{i.} - \bar{\mathbf{x}}'_{i.} \boldsymbol{\beta}] = \rho_{n_i n_i} \hat{v}_{i OLS}
\end{aligned}$$

where $\rho_{n_i n_i} = n_i r / [1 + (n_i - 1) r]$ is Spearman-Brown reliability

suppose $r = .5$

$$\rho_{n_i n_i} = \frac{.5n_i}{.5[2 + (n_i - 1)]} = n_i / (n_i + 1)$$

if $n_i = 1 \Rightarrow \rho_{n_i n_i} = .5$ while if $n_i = 5 \Rightarrow \rho_{n_i n_i} = 5/6$

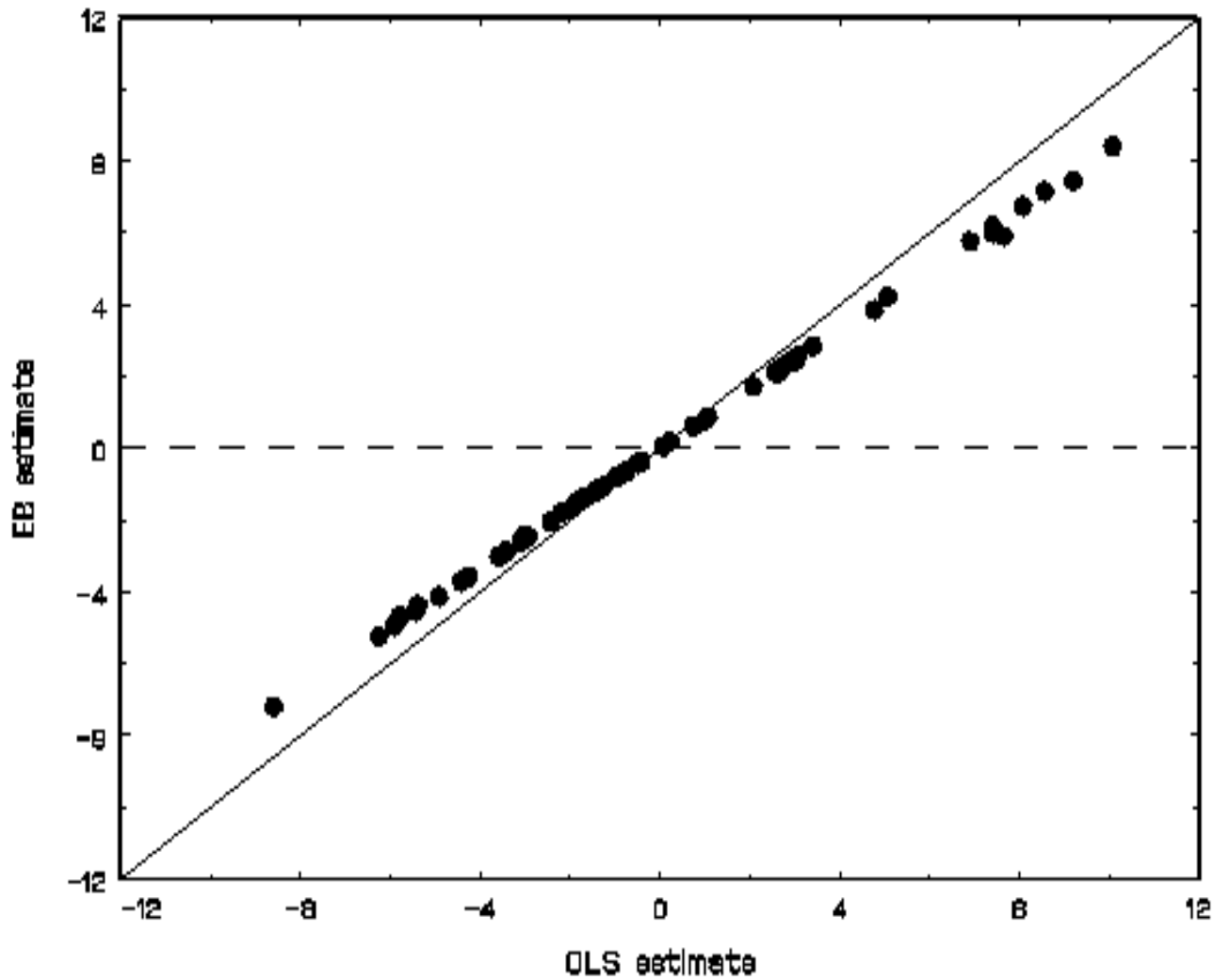
as $n_i \uparrow$ then $\rho_{n_i n_i} \approx 1$

suppose $r = .75$

$$\rho_{n_i n_i} = \frac{n_i}{4/3 + (n_i - 1)} = n_i / (n_i + 1/3)$$

if $n_i = 1 \Rightarrow \rho_{n_i n_i} = 3/4$ while if $n_i = 5 \Rightarrow \rho_{n_i n_i} = 5/5.33$

as $r \uparrow$ then $\rho_{n_i n_i}$ also goes to 1



Plot of the EB estimates versus their OLS counterparts for a random-intercept model of the Reisby data

Similarly,

$$\begin{aligned}\sigma_{v|y_i}^2 &= \left(\mathbf{1}'_i (\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{1}_i + \sigma_v^{-2} \right)^{-1} \\ &= \left(n_i \sigma^{-2} + \sigma_v^{-2} \right)^{-1} \\ &= \left(\frac{n_i}{\sigma^2} + \frac{1}{\sigma_v^2} \right)^{-1} \\ &= \left(\frac{n_i \sigma_v^2 + \sigma^2}{\sigma^2 \sigma_v^2} \right)^{-1} \\ &= \sigma_v^2 \left(\frac{\sigma^2}{n_i \sigma_v^2 + \sigma^2} \right)\end{aligned}$$

$$\begin{aligned}
&= \sigma_v^2 \left(\frac{\sigma^2/n_i}{\sigma_v^2 + \sigma^2/n_i} \right) \\
&= \sigma_v^2 \left(\frac{\sigma_v^2 + \sigma^2/n_i}{\sigma_v^2 + \sigma^2/n_i} - \frac{\sigma_v^2}{\sigma_v^2 + \sigma^2/n_i} \right) \\
&= \sigma_v^2 (1 - \rho_{n_i n_i})
\end{aligned}$$

where

$$\rho_{n_i n_i} = \frac{\sigma_v^2}{\sigma_v^2 + \sigma^2/n_i} = \frac{n_i r}{1 + (n_i - 1)r}$$

EB estimates - multiple random effects

$$\tilde{\mathbf{v}}_i = \left[\mathbf{Z}'_i(\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i + \boldsymbol{\Sigma}_v^{-1} \right]^{-1} \mathbf{Z}'_i(\sigma^2 \mathbf{I}_{n_i})^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})$$

note, using the property $(\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}$

$$\begin{aligned} & \left[\mathbf{Z}'_i(\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i + \boldsymbol{\Sigma}_v^{-1} \right]^{-1} \\ &= \left\{ \left[\mathbf{Z}'_i(\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i \right] \left[\mathbf{I}_r + \left(\mathbf{Z}'_i(\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i \right)^{-1} \boldsymbol{\Sigma}_v^{-1} \right] \right\}^{-1} \\ &= \left\{ \left[\boldsymbol{\Sigma}_v + \left(\mathbf{Z}'_i(\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i \right)^{-1} \right] \boldsymbol{\Sigma}_v^{-1} \right\}^{-1} \left[\mathbf{Z}'_i(\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i \right]^{-1} \\ &= \boldsymbol{\Sigma}_v \left[\boldsymbol{\Sigma}_v + \left(\mathbf{Z}'_i(\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i \right)^{-1} \right]^{-1} \left[\mathbf{Z}'_i(\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i \right]^{-1} \\ &= \mathbf{R} \left[\mathbf{Z}'_i(\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i \right]^{-1} \end{aligned}$$

where $\mathbf{R} = \boldsymbol{\Sigma}_v \left[\boldsymbol{\Sigma}_v + (\mathbf{Z}'_i(\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i)^{-1} \right]^{-1}$ is the “multivariate analogue of reliability” (Bock, 1966). Thus,

$$\begin{aligned}
 \tilde{\mathbf{v}}_i &= \mathbf{R} \left[\mathbf{Z}'_i(\sigma^2 \mathbf{I}_{n_i})^{-1} \mathbf{Z}_i \right]^{-1} \mathbf{Z}'_i(\sigma^2 \mathbf{I}_{n_i})^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \\
 &= \mathbf{R} \left[\mathbf{Z}'_i \mathbf{Z}_i \right]^{-1} \mathbf{Z}'_i (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \\
 &= \mathbf{R} \hat{\mathbf{v}}_i \quad \text{where } \hat{\mathbf{v}}_i = \text{OLS estimator}
 \end{aligned}$$

Similarly,

$$\begin{aligned}\Sigma_{v|y_i} &= (\mathbf{Z}'_i(\sigma^2\mathbf{I}_{n_i})^{-1}\mathbf{Z}_i + \Sigma_v^{-1})^{-1} \\ &= \mathbf{R} \left[\mathbf{Z}'_i(\sigma^2\mathbf{I}_{n_i})^{-1}\mathbf{Z}_i \right]^{-1} \\ &= \mathbf{R} \left\{ \left[\Sigma_v + (\mathbf{Z}'_i(\sigma^2\mathbf{I}_{n_i})^{-1}\mathbf{Z}_i)^{-1} \right] \Sigma_v^{-1}\Sigma_v - \Sigma_v \right\} \\ &= \mathbf{R} (\mathbf{R}^{-1}\Sigma_v - \Sigma_v) \\ &= (\mathbf{I} - \mathbf{R})\Sigma_v\end{aligned}$$

Properties of EB estimates

- \tilde{v}_i “shrunk” towards assumed mean of v_i (0)
- degree of “shrunkness” depends on $\rho_{n_i n_i}$
- number of obs. from individual i , n_i , and intra-unit correlation r influence \tilde{v}_i
 - as $n_i \uparrow$ and/or $r \uparrow$, then $\rho_{n_i n_i} \rightarrow 1$ and $\Rightarrow \tilde{v}_i$ approaches average deviation of the individual (LS estimate)
 - as $n_i \downarrow$ and/or $r \downarrow$, then $\rho_{n_i n_i} \rightarrow 0$ and $\Rightarrow \tilde{v}_i$ approaches assumed mean of v_i , namely 0
- n_i and r have similar effect on the uncertainty about the random effect, $\sigma_{v|y_i}^2$

As a result,

- calculated variance of $\tilde{\mathbf{v}}_i$ estimates will be less than
 - the calculated variance based on the average individual deviations
 - the estimated population variance $\hat{\sigma}_v^2$

Maximum (marginal) Likelihood Solution

The ML solution uses maximum likelihood estimation in the distribution obtained by integrating over the distribution of \mathbf{v} , that is, the marginal distribution:

$$h(\mathbf{y}_i) = \int_{\mathbf{v}} f(\mathbf{y}_i | \mathbf{v}; \boldsymbol{\beta}, \sigma^2) g(\mathbf{v}; \boldsymbol{\Sigma}_v) d\mathbf{v}$$

where

$$g(\mathbf{v}; \boldsymbol{\Sigma}_v) = (2\pi)^{-r/2} |\boldsymbol{\Sigma}_v|^{-1/2} \exp \left[-1/2 \mathbf{v}' \boldsymbol{\Sigma}_v^{-1} \mathbf{v} \right]$$

$$f(\mathbf{y}_i | \mathbf{v}; \boldsymbol{\beta}, \sigma^2) = (2\pi)^{-n_i/2} |\sigma^2 \mathbf{I}_{n_i}|^{-1/2} \exp \left[-1/2 \right. \\ \left. (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{v})' (\sigma^2 \mathbf{I}_{n_i})^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{v}) \right]$$

Need to sum marginal log-likelihoods over the sample:

$$\log L = \sum_{i=1}^N \log h(\mathbf{y}_i)$$

and find out when $\boldsymbol{\beta}$, σ^2 , and $\boldsymbol{\Sigma}_v$ maximize this function

For this, denote the posterior density p_i , likelihood function f_i , prior density g , and marginal log-likelihood h_i

$$\begin{aligned} \log L &= \sum_{i=1}^N \log h_i \\ &= \sum_{i=1}^N \log \int_{\boldsymbol{v}} f_i \cdot g \, d\boldsymbol{v} \end{aligned}$$

and $p_i = f_i \cdot g / h_i$

For the coefficient vector $\boldsymbol{\beta}$ of the fixed covariates we obtain:

$$\begin{aligned}
\frac{\partial \log L}{\partial \boldsymbol{\beta}} &= \sum_{i=1}^N \frac{\partial \log h_i}{\partial \boldsymbol{\beta}} \\
&= \sum_{i=1}^N \frac{1}{h_i} \frac{\partial [\int_{\mathbf{v}} f_i \cdot g \, d\mathbf{v}]}{\partial \boldsymbol{\beta}} \\
&= \sum_{i=1}^N \frac{1}{h_i} \int_{\mathbf{v}} \frac{\partial f_i}{\partial \boldsymbol{\beta}} g \, d\mathbf{v} \\
&= \sum_{i=1}^N \int_{\mathbf{v}} \frac{f_i \cdot g}{h_i} \frac{\partial \log f_i}{\partial \boldsymbol{\beta}} \, d\mathbf{v} \quad (\text{since } f_i = \exp \log f_i) \\
&= \sum_{i=1}^N \int_{\mathbf{v}} p_i \mathbf{X}'_i (\sigma^2 \mathbf{I}_{n_i})^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{v}) \, d\mathbf{v} \\
&= \sigma^{-2} \sum_{i=1}^N \mathbf{X}'_i (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \tilde{\mathbf{v}}_i) \quad (\text{since } \tilde{\mathbf{v}}_i = \int_{\mathbf{v}} \mathbf{v} p_i \, d\mathbf{v})
\end{aligned}$$

Equating to zero yields,

$$\sum_{i=1}^N \mathbf{X}'_i \mathbf{X}_i \boldsymbol{\beta} = \sum_{i=1}^N \mathbf{X}'_i (\mathbf{y}_i - \mathbf{Z}_i \tilde{\mathbf{v}}_i)$$

Thus,

$$\hat{\boldsymbol{\beta}} = \left[\sum_{i=1}^N \mathbf{X}'_i \mathbf{X}_i \right]^{-1} \left[\sum_{i=1}^N \mathbf{X}'_i (\mathbf{y}_i - \mathbf{Z}_i \tilde{\mathbf{v}}_i) \right]$$

as $\tilde{\mathbf{v}}_i \rightarrow 0$, then $\hat{\boldsymbol{\beta}} \rightarrow$ OLS estimates ignoring subject effects

For the residual variance, with $\mathbf{e}_i = \mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{Z}_i\tilde{\mathbf{v}}_i$

$$\begin{aligned}
\frac{\partial \log L}{\partial \sigma^2} &= \sum_{i=1}^N \frac{1}{h_i} \int_{\mathbf{v}} \frac{\partial f_i}{\partial \sigma^2} g \, d\mathbf{v} \\
&= \sum_{i=1}^N \int_{\mathbf{v}} \frac{f_i \cdot g}{h_i} \frac{\partial \log f_i}{\partial \sigma^2} \, d\mathbf{v} \\
&= \sum_{i=1}^N \int_{\mathbf{v}} p_i \left[-\frac{n_i}{2} \sigma^{-2} + \frac{1}{2} \sigma^{-4} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{Z}_i\mathbf{v})' \right. \\
&\quad \left. \times (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{Z}_i\mathbf{v}) \right] \, d\mathbf{v} \\
&= \frac{1}{2} \sigma^{-4} \sum_{i=1}^N \left(-n_i \sigma^2 + \mathbf{e}_i' \mathbf{e}_i + \text{tr} \left[\mathbf{Z}_i \boldsymbol{\Sigma}_{\mathbf{v}|y_i} \mathbf{Z}_i' \right] \right)
\end{aligned}$$

The mysterious last step is because of:

$$\begin{aligned}
& \int_{\mathbf{v}} p_i (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{v})' (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{v}) d\mathbf{v} \\
&= \int_{\mathbf{v}} p_i [(\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \tilde{\mathbf{v}}_i) - \mathbf{Z}_i (\mathbf{v} - \tilde{\mathbf{v}}_i)]' \\
&\quad \times [(\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \tilde{\mathbf{v}}_i) - \mathbf{Z}_i (\mathbf{v} - \tilde{\mathbf{v}}_i)] d\mathbf{v} \\
&= (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \tilde{\mathbf{v}}_i)' (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \tilde{\mathbf{v}}_i) + \text{tr} \mathbf{Z}_i \boldsymbol{\Sigma}_{\mathbf{v}|y_i} \mathbf{Z}_i'
\end{aligned}$$

using

$$\boldsymbol{\Sigma}_{\mathbf{v}|y_i} = \int_{\mathbf{v}} (\mathbf{v} - \tilde{\mathbf{v}}_i)(\mathbf{v} - \tilde{\mathbf{v}}_i)' p_i d\mathbf{v}$$

Equating to zero yields (with $\hat{\mathbf{e}}_i = \mathbf{y}_i - \mathbf{X}_i\hat{\boldsymbol{\beta}} - \mathbf{Z}_i\tilde{\mathbf{v}}_i$),

$$\hat{\sigma}^2 = \left(\sum_i^N n_i \right)^{-1} \sum_i^N \hat{\mathbf{e}}_i' \hat{\mathbf{e}}_i + \text{tr} \left[\mathbf{Z}_i \boldsymbol{\Sigma}_{v|y_i} \mathbf{Z}_i' \right]$$

Notice similarity of above to ML estimator in multiple regression model:

$$\hat{\sigma}^2 = \left(\sum_i^N n_i \right)^{-1} \sum_i^N \hat{\mathbf{e}}_i' \hat{\mathbf{e}}_i$$

as $\tilde{\mathbf{v}}_i \rightarrow 0$, then $\boldsymbol{\Sigma}_{v|y_i} \rightarrow 0$, then $\hat{\sigma}^2$ goes to ML error estimate ignoring subjects (which is approximately equal to OLS error estimate, *i.e.*, MSE)

Now things get really ugly!

For the solution for Σ_v we will use the “vec” and “vech” notation (Henderson & Searle, 1979) and

$$\frac{\partial \log |\Sigma_v|}{\partial \text{vec} \Sigma_v} = \text{vec} \Sigma_v^{-1}$$

$$\begin{aligned} \frac{\partial \mathbf{v}' \Sigma_v^{-1} \mathbf{v}}{\partial \text{vec} \Sigma_v} &= \frac{\partial \text{tr} \Sigma_v^{-1} \mathbf{v} \mathbf{v}'}{\partial \text{vec} \Sigma_v} \\ &= -\text{vec} \Sigma_v^{-1} \mathbf{v} \mathbf{v}' \Sigma_v^{-1} \end{aligned}$$

the “vec” of an $n \times n$ matrix \mathbf{A} is the column vector $[\mathbf{a}'_1, \mathbf{a}'_2, \dots, \mathbf{a}'_n]'$, or in words, the “vec” operator stacks the column vectors of a matrix one on top of each other to form one large column vector

the “vech” of an $n \times n$ symmetric matrix \mathbf{A} is the $n(n + 1)/2$ column vector $[a_{11}, a_{21}, a_{22}, a_{31}, a_{32}, a_{33}, \dots, a_{nn}]'$, or in words, the “vech” operator stacks the on or above diagonal elements of a square matrix on top of each other to form one column vector

We also need:

$$\begin{aligned} \int_{\mathbf{v}} p_i \mathbf{v} \mathbf{v}' d\mathbf{v} &= \int_{\mathbf{v}} p_i [(\mathbf{v} - \tilde{\mathbf{v}}_i) + \tilde{\mathbf{v}}_i] [(\mathbf{v} - \tilde{\mathbf{v}}_i) + \tilde{\mathbf{v}}_i]' d\mathbf{v} \\ &= \Sigma_{v|y_i} + \tilde{\mathbf{v}}_i \tilde{\mathbf{v}}_i' \end{aligned}$$

Derivatives with respect to the $r(r + 1)/2$ unique elements of $\boldsymbol{\Sigma}_v$:

$$\begin{aligned}
\frac{\partial \log L}{\partial \text{vech} \boldsymbol{\Sigma}_v} &= \sum_{i=1}^N \frac{1}{h_i} \int \boldsymbol{v} \frac{\partial g}{\partial \text{vech} \boldsymbol{\Sigma}_v} f_i d\boldsymbol{v} \\
&= \sum_{i=1}^N \int \boldsymbol{v} \frac{f_i \cdot g}{h_i} \frac{\partial \log g}{\partial \text{vech} \boldsymbol{\Sigma}_v} d\boldsymbol{v} \\
&= \mathbf{G}' \text{vec} \sum_i^N \int \boldsymbol{v} p_i \left[-\frac{1}{2} \boldsymbol{\Sigma}_v^{-1} + \frac{1}{2} \boldsymbol{\Sigma}_v^{-1} \boldsymbol{v} \boldsymbol{v}' \boldsymbol{\Sigma}_v^{-1} \right] \\
&= \frac{1}{2} \mathbf{G}' \sum_i^N \text{vec} \boldsymbol{\Sigma}_v^{-1} \left(-\boldsymbol{\Sigma}_v + \boldsymbol{\Sigma}_{v|y_i} + \tilde{\boldsymbol{v}}_i \tilde{\boldsymbol{v}}_i' \right) \boldsymbol{\Sigma}_v^{-1}
\end{aligned}$$

where the matrix \mathbf{G} transforms the vech of a square matrix into the vec (McCulloch, 1982), that is, $\text{vec} \mathbf{A} = \mathbf{G} \text{vech} \mathbf{A}$ for a square matrix \mathbf{A}

setting this monster equal to zero, we find

$$\hat{\Sigma}_v = \frac{1}{N} \sum_i^N (\tilde{\mathbf{v}}_i \tilde{\mathbf{v}}_i' + \Sigma_{v|y_i})$$

e.g., Subject variance for random intercepts model

$$\hat{\sigma}_v^2 = \frac{1}{N} \sum_i^N \tilde{v}_i^2 + \sigma_{v|y_i}^2$$

- Average of EB individual effect squared (\approx sample estimate of variance of EB effect) plus the EB estimate of the uncertainty about the individual effect
- $\sigma_v^2, \tilde{v}_i, \sigma_{v|y_i}^2$ are one happy family (they move together!)
 - as \tilde{v}_i and $\sigma_{v|y_i}^2$ go to 0, so does σ_v^2

- as covariates included in model $\Rightarrow \tilde{v}_i$ gets smaller, thus so does σ_v^2
- $\hat{\sigma}_v^2$ is larger than the calculated variance of \tilde{v}_i
 - since calculated variance equals

$$\frac{1}{N} \sum_i (\tilde{v}_i - \bar{x}_{\tilde{v}_i})^2 / (N - 1)$$

and $\bar{x}_{\tilde{v}_i} \approx 0$

note: this difference is diminished as all n_i get large, then $\sigma_{v|y_i}^2 \rightarrow 0$ and $\hat{\sigma}_v^2$ is approximately equal to the calculated variance of \tilde{v}_i

Within-subjects and between-subjects components

random intercepts model

Within-subjects model - level 1

$$y_{ij} = b_{0i} + \boldsymbol{\beta}'_{(1)} \mathbf{x}_{(1)ij} + \varepsilon_{ij}$$

i = $1 \dots N$ subjects
 j = $1 \dots n_i$ observations for subject i

Between-subjects models - level 2

$$b_{0i} = \beta_0 + \boldsymbol{\beta}'_{(2)} \mathbf{x}_{(2)i} + v_i$$

Total model

$$y_{ij} = \boldsymbol{\beta}' \mathbf{x}_{ij} + v_i + \varepsilon_{ij}$$

where

$$\mathbf{x}'_{ij} = [\mathbf{1} \quad \vdots \quad \mathbf{x}_{(1)ij} \quad \vdots \quad \mathbf{x}_{(2)i}]$$

and

$$\boldsymbol{\beta}' = [\beta_0 \quad \vdots \quad \boldsymbol{\beta}'_{(1)} \quad \vdots \quad \boldsymbol{\beta}'_{(2)}]$$

$$\varepsilon_{ij} \sim \mathcal{NID}(0, \sigma^2) \quad v_i \sim \mathcal{NID}(0, \sigma_v^2)$$

Adding in $\mathbf{x}_{(1)ij}$ can reduce σ^2 and σ_v^2

Adding in $\mathbf{x}_{(2)i}$ can reduce σ_v^2

Estimation via EM algorithm: opposite process of “I do cocaine so I can work more, so I can do more cocaine, so I can work more, *etc.*,”

<u>Cocaine</u>		<u>Work</u>	<u>Health</u>
do cocaine	→	work more	declines
do more cocaine	→	work more	declines more
do even more cocaine	→	work even more	declines even more
...
do a ton of cocaine	→	always working	death

<u>M-Step (ML)</u>		<u>E-Step (EB)</u>	<u>Estimation</u>
starting values $\beta, \sigma^2, \Sigma_v$	→	estimate $\tilde{\mathbf{v}}_i \Sigma_{v y_i}$	improves
re-estimate $\beta, \sigma^2, \Sigma_v$	→	re-estimate $\tilde{\mathbf{v}}_i \Sigma_{v y_i}$	improves more
re-re-estimate $\beta, \sigma^2, \Sigma_v$	→	re-re-estimate $\tilde{\mathbf{v}}_i \Sigma_{v y_i}$	improves even more
...
RE-estimate $\beta, \sigma^2, \Sigma_v$	→	RE-estimate $\tilde{\mathbf{v}}_i \Sigma_{v y_i}$	convergence

⇒ EM is better than cocaine since EM leads to convergence and not death

EM solution - shifts between

- estimation of individual parameters $\tilde{\mathbf{v}}_i$ and $\Sigma_{v|y_i}$
- estimation of the structural β and σ^2 and population parameters Σ_v
- convergence by EM can be slow and torturous
- need solution which utilizes second derivatives
 - Newton-Raphson procedure
 - Fisher-scoring procedure
 - * replaces second derivatives with their expectations (negative of the information matrix)
 - * multiplying the vector of first derivatives by the inverse of the information matrix provides the corrections at each iteration

For a parameter vector Θ , the estimated values for iteration $\iota + 1$ are obtained by

Newton-Raphson

$$\Theta_{\iota+1} = \Theta_{\iota} - \left[\frac{\partial^2 \log L}{\partial \Theta_{\iota} \partial \Theta'_{\iota}} \right]^{-1} \frac{\partial \log L}{\partial \Theta_{\iota}}$$

Fisher scoring

$$\Theta_{\iota+1} = \Theta_{\iota} - \varepsilon \left[\frac{\partial^2 \log L}{\partial \Theta_{\iota} \partial \Theta'_{\iota}} \right]^{-1} \frac{\partial \log L}{\partial \Theta_{\iota}}$$

where the information matrix has the form:

$$I = -\varepsilon \left[\frac{\partial^2 \log L}{\partial \Theta_{\iota} \partial \Theta'_{\iota}} \right] = \begin{bmatrix} I(\beta) \\ 0 & I(\Sigma_v) \\ 0 & I(\sigma^2, \Sigma_v) & I(\sigma^2) \end{bmatrix}$$

The information matrix - random intercepts model

$$I = \begin{bmatrix} I(\boldsymbol{\beta}) & 0 & 0 \\ 0 & I(\sigma_v^2) & I(\sigma^2, \sigma_v^2) \\ 0 & I(\sigma^2, \sigma_v^2) & I(\sigma^2) \end{bmatrix}$$

where

$$I(\boldsymbol{\beta}) = \sigma^{-4} \sum_{i=1}^N \mathbf{X}'_i (\sigma^2 \mathbf{I}_{n_i} - \sigma_{v|y_i}^2 \mathbf{1}_i \mathbf{1}'_i) \mathbf{X}_i$$

$$I(\sigma_v^2) = \frac{1}{2} \sigma_v^{-8} \sum_{i=1}^N (\sigma_v^2 - \sigma_{v|y_i}^2)^2$$

$$I(\sigma^2) = \frac{1}{2} \sigma^{-8} \sum_{i=1}^N n_i (\sigma^2 - \sigma_{v|y_i}^2)^2$$

$$I(\sigma^2, \sigma_v^2) = \frac{1}{2} \sigma^{-4} \sigma_v^{-4} \sum_{i=1}^N n_i \sigma_{v|y_i}^4$$

- At convergence, I^{-1} provides asymptotic variances and covariances of ML estimates
- For standard errors, take the square root of the diagonal elements of I^{-1}
- single parameter hypothesis test $H_0 : \theta = 0$

$$z = \frac{\hat{\theta}}{\hat{se}(\hat{\theta})}$$

Reparameterization of the Variance terms

- Fisher scoring estimation of Σ_v is problematic when diagonal elements get small (near-zero variances)
- reparameterize Σ_v in terms of the Gaussian factorization of a symmetric matrix
- use exponential transformation for diagonal matrix \mathbf{D} corresponding to variance parameters

Specifically, consider $\Sigma_v = \mathbf{L}\mathbf{D}\mathbf{L}' = \mathbf{L} \exp(\mathbf{\Pi}) \mathbf{L}'$

e.g., , if $r = 3$,

$$\Sigma_v = \begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & 1 \end{bmatrix} \begin{bmatrix} e^{\pi_1} & 0 & 0 \\ 0 & e^{\pi_2} & 0 \\ 0 & 0 & e^{\pi_3} \end{bmatrix} \begin{bmatrix} 1 & l_{21} & l_{31} \\ 0 & 1 & l_{32} \\ 0 & 0 & 1 \end{bmatrix}$$

\Rightarrow estimate \mathbf{L} and $\mathbf{\Pi}$

Derivatives - Magnus (1988); Magnus & Neudecker (1988)

diagonal elements in $\mathbf{\Pi}$

$w(\mathbf{A})$ denotes vector which contains only the diagonal elements of a square matrix \mathbf{A}

$$\frac{\partial \log L}{\partial(w(\mathbf{\Pi}))} = \frac{1}{2} \mathbf{D}^{-1} \mathbf{H}_r(\mathbf{L}^{-1} \otimes \mathbf{L}^{-1}) \left\{ \left(\sum_i^N \text{vec} [\boldsymbol{\Sigma}_v|y_i + \tilde{\mathbf{v}}_i \tilde{\mathbf{v}}_i'] \right) - N \text{vec} \boldsymbol{\Sigma}_v \right\}$$

where \otimes denotes the Kronecker product, and where the $r \times r^2$ matrix \mathbf{H}_r eliminates the off-diagonal elements from the vec of a $r \times r$ matrix

elements of \mathbf{L} below the main diagonal

$\tilde{\mathbf{v}}(\mathbf{A})$ denotes vector which contains only the elements below the main diagonal of a square matrix \mathbf{A}

$$\frac{\partial \log L}{\partial (\tilde{\mathbf{v}}(\mathbf{L}))} = \tilde{\mathbf{J}}_r(\mathbf{L}^{-1} \otimes \boldsymbol{\Sigma}_v^{-1}) \left(\sum_i^N \text{vec} \left[\boldsymbol{\Sigma}_{v|y_i} + \tilde{\mathbf{v}}_i \tilde{\mathbf{v}}_i' \right] \right)$$

where the $r \times r^2$ matrix $\tilde{\mathbf{J}}_r$ eliminates the elements of the vec which are on or above the main diagonal

Restricted Maximum Likelihood (REML)

- introduced by Patterson & Thompson (1971) *Biometrika*
- object: to remove bias in ML estimates of variance terms
- *e.g.*, in linear regression

$$\hat{\sigma}_{ML}^2 = SSE/N \quad \hat{\sigma}_{REML}^2 = SSE/(N - p)$$

⇒ important as N is small and/or p is large
(here, p includes intercept)

- idea: in ML estimation of variance terms treats β as known, however β are estimated
- difference occurs as p gets relatively large ⇒ MLEs of variances and their standard errors are downward biased
- for fixed p , methods asymptotically equivalent as either or both of N and n_i tend to infinity

- relatively simple modification of ML algorithm (DLZ p. 64-68)
- REML estimators are preferable, however
- **the REML deviance ($-2 \log L$) can be used only for comparison of models with identical regression parts**
- in SAS, use METHOD=REML (the default)