



JOHNS HOPKINS
BLOOMBERG SCHOOL
of PUBLIC HEALTH

Lecture 12

Finish Implementation of WLS/robust variance in R

Introduction to Linear Mixed Models

Weighted least squares review

Assume a longitudinal design with (Y_{ij}, X_{ij}) for $i = 1, \dots, m$ and $j = 1, \dots, n_i$

The model for subject i can be expressed as $Y_i = X_i\beta + \varepsilon_i, \varepsilon_i \sim MVV(0, V_i)$

Then, we can stack the subject models together as $Y = X\beta + \varepsilon, \varepsilon \sim MVN(0, \Sigma)$

$$\text{where } Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_m \end{bmatrix}, X = \begin{bmatrix} X_1 \\ \vdots \\ X_m \end{bmatrix}, \Sigma = \begin{bmatrix} V_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & V_m \end{bmatrix}$$

We have shown how to extend OLS to WLS to account for Σ (instead of OLS assumption of $\sigma^2 I$).

The score equations for the WLS solution is: $X'\Sigma^{-1}(Y - X\beta) = 0$

Yielding: $\hat{\beta}_{wls} = (X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1}Y$ and $Var(\hat{\beta}_{wls}) = (X'\Sigma^{-1}X)^{-1}$



Implementation in R

- ▶ In Lecture 11, we walked through the required exploratory analysis for longitudinal data. In order to fit a WLS model, we have to understand
 - ▶ the mean model
 - ▶ the correlation structure
 - ▶ the patterns of variance
- ▶ Using the simulated NEPAL1 dataset, we settled on:
 - ▶ $Y_{ij} = \beta_0 + \beta_1 age_{ij} + \beta_2 (age_{ij} - 6)^+ + \varepsilon_{ij}, \varepsilon_{ij} \sim Normal$
 - ▶ $Corr(\varepsilon_{ij}, \varepsilon_{ik}) = \rho^{|j-k|}$
 - ▶ $Var(\varepsilon_{ij}) = f(age_{ij})$
- ▶ Fit this model using gls, but also considered two other correlation models:
 - ▶ $Corr(\varepsilon_{ij}, \varepsilon_{ik}) = \rho$
 - ▶ $Corr(\varepsilon_{ij}, \varepsilon_{ik}) = \rho_{jk}$



Implementation in R

- ▶ AR1 model, constant variance
 - ▶ `mod.gls.exch.het = gls(wt ~ age + age_sp6, data = nepal1, correlation = corAR1(form = ~num | id))`
- ▶ Exchangeable / compound symmetry, constant variance
 - ▶ `mod.gls.exch.het = gls(wt ~ age + age_sp6, data = nepal1, correlation = corCompSymm(form = ~1 | id))`
- ▶ Unstructured, constant variance
 - ▶ `mod.gls.exch.het = gls(wt ~ age + age_sp6, data = nepal1, correlation = corSymm(form = ~num | id))`
- ▶ Allowing for variance to depend on age
 - ▶ `mod.gls.exch.het = gls(wt ~ age + age_sp6, data = nepal1, correlation = corAR1(form = ~num | id), weights = varFunc(~age))`



Generalized Estimating Equations

Weighted least squares is a special case of a general method called Generalized Estimating Equations (GEE).

In the case of $Y_i \sim MVN(X_i\beta, V_i)$, the WLS/GEE method finds the values of β that equates the score equations (i.e. estimating equations) to 0. In the case of independent Y_i , for $i = 1, \dots, m$, the $\hat{\beta}_{wls}$ solves:

$$S(\beta, \theta) = \sum_{i=1}^m \frac{\partial X_i \beta}{\partial \beta} V_i(\theta)^{-1} (Y_i - X_i \beta) = 0$$

The estimation procedure is iterative, same as WLS.

- ▶ We will discuss GEE again in 654.
- ▶ One advantage of GEE is that you don't have to specify a multivariate distribution for Y_i , so long as we can specify $E(Y_{ij})$, $Var(Y_{ij})$, $Corr(Y_{ij}, Y_{ik})$, then we can solve for β and make inferences.
 - ▶ This is nice because multivariate Bernoulli or Poisson distributions are quite complicated.



Generalized Estimating Equations

- ▶ Why do we care about GEE?
- ▶ Historical: Kung-Yee Liang and Scott Zeger derived the method; motivated by a longitudinal design with binary outcome
- ▶ The *gls* function in R is limiting in that it does not directly compute robust variance estimates; so you only have access to standard error estimates based on the model you specify.
 - ▶ See *clubSandwich* which should work on *gls* objects and produce robust variance estimates; I have had trouble with this function
- ▶ The typical implementation of GEE is to provide both model based (similar to *gls*) and robust variance estimates.
- ▶ Note: in *gee* function in R, you specify a model for the within subject/cluster correlation structure (R_i) and the model assumes a constant variance



Robust Variance Estimation

- ▶ We refit several of the models we considered before (plus a few additional models) and obtained robust variance estimates for some of the models
 - ▶ OLS with and without robust variance estimate
 - ▶ Exchangeable, constant variance with and without robust variance estimate
 - ▶ Exchangeable, variance depends on age
 - ▶ AR1, constant variance with and without robust variance estimate
 - ▶ AR1, variance depends on age

	OLS	OLS-RV	Exch	Exch-Het	Exch-RV
Intercept	5.074 (0.289)	5.074 (0.157)	4.915 (0.199)	5.116 (0.082)	4.916 (0.192)
Age	0.486 (0.06)	0.486 (0.026)	0.512 (0.029)	0.468 (0.025)	0.511 (0.032)
Age_SP1	-0.344 (0.071)	-0.344 (0.028)	-0.366 (0.034)	-0.314 (0.03)	-0.366 (0.034)

	OLS	OLS-RV	AR1	AR1-Het	AR1-RV
Intercept	5.074 (0.289)	5.074 (0.157)	4.98 (0.19)	5.106 (0.073)	4.99 (0.165)
Age	0.486 (0.06)	0.486 (0.026)	0.495 (0.022)	0.467 (0.024)	0.494 (0.023)
Age_SP1	-0.344 (0.071)	-0.344 (0.028)	-0.348 (0.025)	-0.313 (0.028)	-0.347 (0.023)

Which model is “best”?

- ▶ You can use an information criteria statistic, which combines information about the fit (i.e. sums of squares residuals) and the complexity of the model (i.e. number of parameters in the model, including the parameters for variance/covariance)
- ▶ Akaike’s Information Criteria: $-2 \log\text{-likelihood} + 2 \times p$, where p is the number of parameters in the model
- ▶ Models with smaller AIC values are “better”

##	df	AIC
## mod.gls.exch.fit	5	729.3473
## mod.gls.exch.het.fit	5	827.3266
## mod.gls.ar1.fit	5	589.6508
## mod.gls.ar1.het.fit	5	731.3010

Two approaches for modeling longitudinal data

- ▶ Descriptive: Marginal model, goal is to describe and make inference for the mean model.
 - ▶ Have to account for the variance/correlation structure to get valid inferences
 - ▶ But we don't necessary care about describing that structure.
- ▶ Etiologic: Conditional models: we are specifically interested in describing where the correlation comes from.
 - ▶ E.g. the current observation may depend on the prior observation (transition model)
 - ▶ E.g. each subject may be distinguished by latent variables/random effects which separate their data from other subjects data.
 - ▶ The goal is to describe the population level patterns (similar to marginal models) but also quantify heterogeneity across subjects in features of the data that are very important for public health researchers, e.g. variation in child specific growth rates.



Transition Models

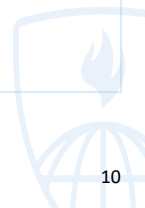
Here past observations of the outcome cause future values of the outcomes. Namely, a transition model where the current value of Y_{ij} depends on the p past observations can be expressed as:

$$E(Y_{ij}|Y_{ij-1}, \dots, Y_{ij-p}, X_{ij}) = X_{ij}'\beta^c + \sum_{k=1}^p \alpha_k Y_{ij-k}$$

The special case of the AR-1 model is where $p = 1$.

$$E(Y_{ij}|Y_{ij-1}, X_{ij}) = X_{ij}'\beta^c + \alpha Y_{ij-1}$$

Note that the models above make a strong assumption: the relationship between the mean of Y_{ij} and X_{ij} is the same regardless of the past values of Y . This assumption can be made flexible by including interaction terms of components of X_{ij} and past values of Y .



Subject specific or random effects models

- ▶ Consider the data generating structure within the NEPAL1 and NEPAL2 simulated datasets:
 - ▶ Children are enrolled between 1 and 5 months of age
 - ▶ Children are followed over time and growth in weight is recorded every 4 months for a total of 5 assessments (enrollment + 4 follow-ups)
- ▶ For each child, we can think of the child's growth:

$$Y_{ij} = \beta_{0i} + \beta_{1i}age_{ij} + \beta_{2i}(age_{ij} - 6)^+ + e_{ij}$$



Subject specific or random effects models

- ▶ The β describe characteristics of the specific children and we assume that these characteristics can vary from child to child, specifically,

$$\begin{bmatrix} \beta_{0i} \\ \beta_{1i} \\ \beta_{2i} \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} + \begin{bmatrix} b_{0i} \\ b_{1i} \\ b_{2i} \end{bmatrix}$$

$$\beta_i = \beta + b_i, b_i \sim MVN(0, D), D = \begin{bmatrix} \tau_0^2 & \tau_{01} & \tau_{02} \\ \tau_{01} & \tau_1^2 & \tau_{12} \\ \tau_{02} & \tau_{12} & \tau_2^2 \end{bmatrix}$$

Visualization



General Model

We can rewrite the model above as:

$$Y_{ij} = (\beta_0 + b_{0i}) + (\beta_1 + b_{1i})age_{ij} + (\beta_2 + b_{2i})(age_{ij} - 6)^+ + e_{ij}$$

In vector notation,

$$Y_{ij} = \begin{bmatrix} 1 \\ age_{ij} \\ (age_{ij} - 6)^+ \end{bmatrix}' \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} + \begin{bmatrix} 1 \\ age_{ij} \\ (age_{ij} - 6)^+ \end{bmatrix}' \begin{bmatrix} b_{0i} \\ b_{1i} \\ b_{2i} \end{bmatrix} + e_{ij}$$

Even more generally,

$$Y_{ij} = X_{ij}'\beta + Z_{ij}'b_i + e_{ij}$$

where $b_i \sim MVN(0, D)$, e_{ij} iid $N(0, \sigma^2)$ and b_i and e_{ij} are independent!

Means and Variances

- In the random effects model, we express the mean function for an individual subject as:

$$E(Y_{ij}|X_{ij}, b_i) = X_{ij}\beta + Z_{ij}b_i$$

- We can express the population mean (i.e. the average over all subjects) as:

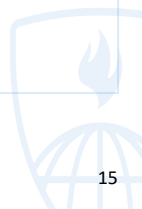
$$E(Y_{ij}|X_{ij}) = E[E(Y_{ij}|X_{ij}, b_i)] = E[X_{ij}\beta + Z_{ij}b_i] = X_{ij}\beta$$

- We can derive the variance of Y_{ij} as

$$Var(Y_{ij}|X_{ij}) = E_{b_i}[Var(Y_{ij}|X_{ij}, b_i)] + Var_{b_i}[E(Y_{ij}|X_{ij}, b_i)]$$

$$Var(Y_{ij}|X_{ij}) = E_{b_i}[\sigma^2] + Var_{b_i}[X_{ij}'\beta + Z_{ij}'b_i]$$

$$Var(Y_{ij}|X_{ij}) = \sigma^2 + Z_{ij}'DZ_{ij}$$



Correlation

- ▶ Assume a random intercept only model:

- ▶ $Y_{ij} = \beta_{0i} + \beta_1 age_{ij} + \beta_2 (age_{ij} - 6)^+ + e_{ij}, \beta_{0i} \sim N(\beta_0, \tau_0^2), e_{ij} \sim N(0, \sigma^2), Cov(\beta_{0i}, e_{ij}) = 0$
- ▶ $Y_{ij} = \beta_0 + b_{0i} + \beta_1 age_{ij} + \beta_2 (age_{ij} - 6)^+ + e_{ij}, b_{0i} \sim N(0, \tau_0^2), e_{ij} \sim N(0, \sigma^2), Cov(b_{0i}, e_{ij}) = 0$

- ▶ What is $Cov(Y_{ij}, Y_{ik})$?



Correlation

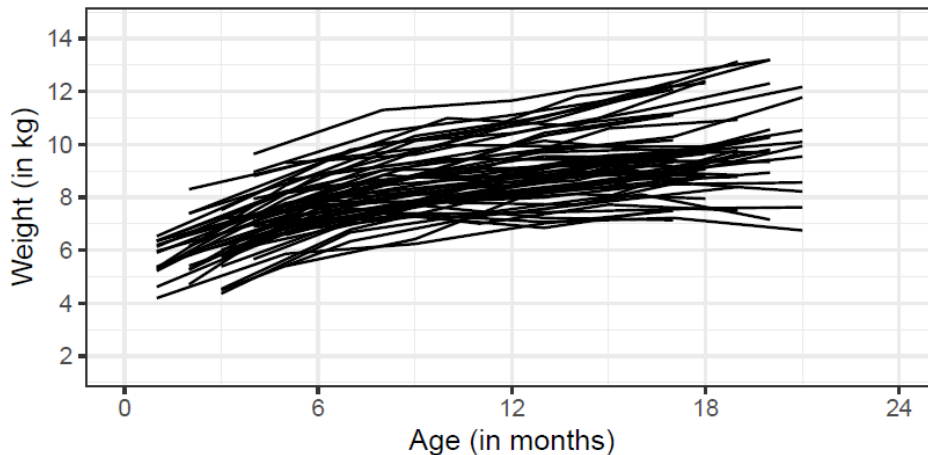
- ▶ Assume a random intercept and random slope for age model:

$$Y_{ij} = \beta_0 + b_{0i} + (\beta_1 + b_{1i})age_{ij} + \beta_2(age_{ij} - 6)^+ + e_{ij}, \text{ where}$$

$$b_{0i} \sim N(0, \tau_0^2), b_{1i} \sim N(0, \tau_1^2), Cov(b_{0i}, b_{1i}) = \tau_{01}, e_{ij} \sim N(0, \sigma^2), Cov(b_{0i}, e_{ij}) = 0, Cov(b_{1i}, e_{ij}) = 0$$



Example: NEPAL1 simulated data



- Fit the following model to the NEPAL1 simulated dataset

$$Y_{ij} = \beta_0 + b_{0i} + (\beta_1 + b_{1i})age_{ij} + \beta_2(age_{ij} - 6)^+ + e_{ij}, \text{ where}$$

$$b_{0i} \sim N(0, \tau_0^2), b_{1i} \sim N(0, \tau_1^2), Cov(b_{0i}, b_{1i}) = \tau_{01}, e_{ij} \sim N(0, \sigma^2), Cov(b_{0i}, e_{ij}) = 0, Cov(b_{1i}, e_{ij}) = 0$$

```
lmer(wt~age+age_sp6+(1+age|id),data=nepal2, control = lmerControl(optimizer="Nelder_Mead"))
```

Example: NepalI simulated data

- ▶ What is the estimate of the population mean birth weight?

```
##               Estimate Std. Error   t value
## (Intercept)  4.9777731 0.15426618  32.26743
## age         0.4984283 0.01867078  26.69563
## age_sp6     -0.3497761 0.01802296 -19.40725
```

```
summary(fit)$varcor
```

- ▶ What is the estimate of the population mean growth rate in the first 6 months of life?

```
## Groups      Name      Std.Dev. Corr
## id          (Intercept) 1.045047
##            age          0.082252 -0.345
## Residual                                0.281274
```

```
est = fixef(fit)
```

- ▶ What is the estimate of the difference in the population mean growth rate after 6 months compared to during the first 6 months of life?

Example: Nepali simulated data

- ▶ For a given child at a specific age, how much do the observed weights differ (+/-) on average from the child's average weight at that age?

```
##               Estimate Std. Error   t value
## (Intercept)  4.9777731 0.15426618  32.26743
## age          0.4984283 0.01867078  26.69563
## age_sp6      -0.3497761 0.01802296 -19.40725
```

```
summary(fit)$varcor
```

```
## Groups   Name      Std.Dev. Corr
## id       (Intercept) 1.045047
##          age          0.082252 -0.345
## Residual              0.281274
```

```
est = fixef(fit)
```

- ▶ Construct an interval that contains 95% of birthweights for Nepali children.

Example: Nepali simulated data

- Construct an interval that contains 95% of growth rates for Nepali children under 6 months of age

```
##               Estimate Std. Error   t value
## (Intercept)  4.9777731  0.15426618  32.26743
## age          0.4984283  0.01867078  26.69563
## age_sp6      -0.3497761  0.01802296 -19.40725
```

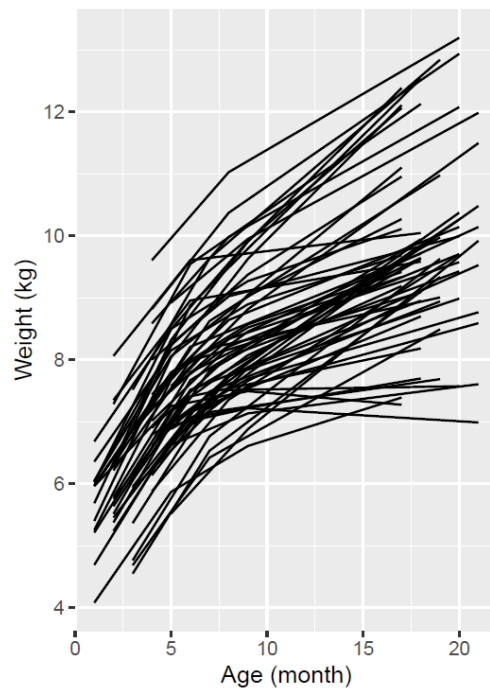
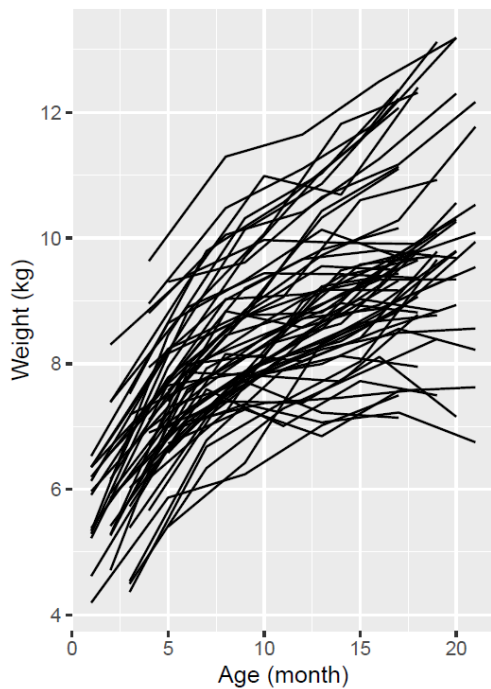
```
summary(fit)$varcor
```

```
## Groups   Name          Std.Dev. Corr
## id       (Intercept)  1.045047
##          age          0.082252 -0.345
## Residual                    0.281274
```

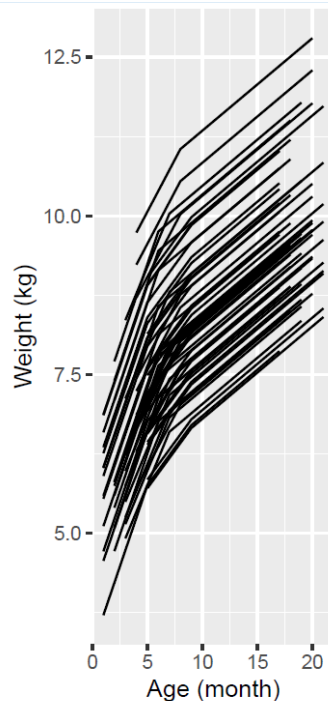
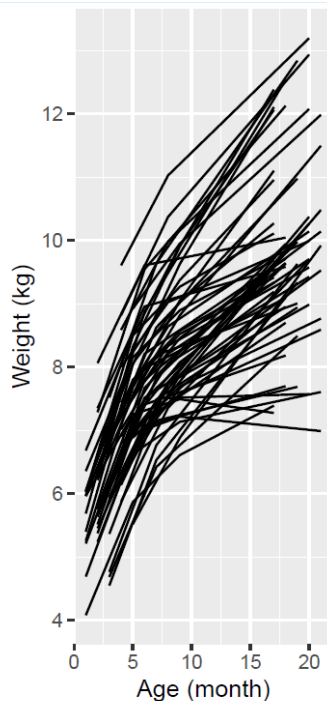
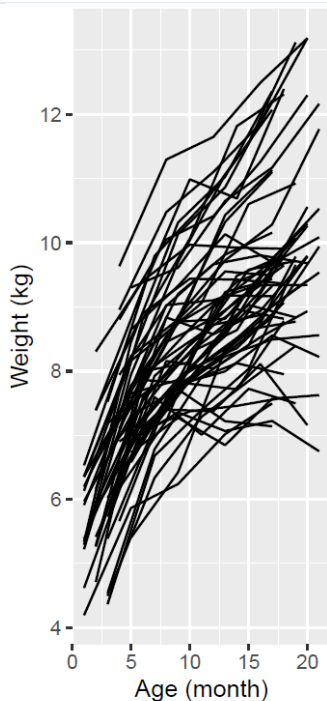
```
est = fixef(fit)
```

- Construct an interval that contains 95% of growth rates for Nepali children over 6 months of age

Example: NepaliI simulated data



Example: NepaliI simulated data



Information criterion comparison

10. Compare the fits from the *gls* models and random intercept and slope models, using AIC.

```
AIC(mod.gls.exch.fit,mod.gls.exch.het.fit,mod.gls.ar1.fit,mod.gls.ar1.het.fit)
```

##		df	AIC
##	mod.gls.exch.fit	5	729.3473
##	mod.gls.exch.het.fit	5	827.3266
##	mod.gls.ar1.fit	5	589.6508
##	mod.gls.ar1.het.fit	5	731.3010

```
AIC(fit,fit.int)
```

##		df	AIC
##	fit	7	526.8088
##	fit.int	5	729.3473

NOTE: The data was generated under the random intercept and random slope for age model!



Next time....

- ▶ On Tuesday March 9th, we will review
 - ▶ linear mixed models
 - ▶ analyze NEPAL2 together

