

# Covariance Structure Approach to Within-cases

Using SAS proc mixed

This slide show is a free open source document.  
See the last slide for copyright information.

# General mixed model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\epsilon}$$

- $\mathbf{b} \sim N(\mathbf{0}, \boldsymbol{\Sigma}_b)$  is a vector of random effects.
- $\mathbf{Z}$  is another matrix of fixed explanatory variable values.
- $\text{cov}(\boldsymbol{\epsilon})$  need not be diagonal – can accommodate non-independence between observations from the same case.
- We won't even use  $\mathbf{Z}\mathbf{b}$ .
- So we are just scratching the surface of what `proc mixed` can do.

# Advantages

- Straightforward: It's familiar univariate regression.
- Variances of beta-hats are different, because of correlated observations.
- Nicer treatment of missing data (valid if missing at random).
- Can have time-varying covariates.
- Flexible modeling of non-independence within cases.
- Can accommodate more factor levels than cases (with assumptions).

Usual covariance matrix of

$y_1, \dots, y_n$

$$\begin{bmatrix} \sigma^2 & 0 & 0 & \dots & 0 & 0 \\ 0 & \sigma^2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \sigma^2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & 0 \\ 0 & 0 & 0 & \dots & \sigma^2 & 0 \\ 0 & 0 & 0 & \dots & 0 & \sigma^2 \end{bmatrix}$$

# In the covariance structure approach

- There are  $n$  “subjects.”
- There are  $k$  (“repeated”) measurements per subject.
- There are  $nk$  rows in the data file:  $n$  blocks of  $k$  rows.
- Data are multivariate normal (dimension  $nk$ )
- Familiar regression model for the vector of means.
- Special structure for the variance-covariance matrix: not just a diagonal matrix with  $\sigma^2$  on the main diagonal.

# Structure of the variance-covariance matrix

- Covariance matrix of the data has a **block diagonal** structure:  $n \times n$  matrix of little  $k \times k$  variance-covariance matrices (partitioned matrix)
- Off diagonal matrices are all zeros -- no correlation between data from different cases
- Matrices on the main diagonal are all the same (equal variance assumption)

# Block Diagonal Covariance Matrix of $y_1, \dots, y_n$

$$\begin{bmatrix} \Sigma & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \Sigma & \dots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \dots & \Sigma & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \Sigma \end{bmatrix}$$

$\Sigma$  is the matrix of variances and covariances of the data from a single subject.

# $\Sigma$ may have different *structures*

- May be unknown

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{1,2} & \sigma_{1,3} & \sigma_{1,4} \\ \sigma_{2,1} & \sigma_2^2 & \sigma_{2,3} & \sigma_{2,4} \\ \sigma_{3,1} & \sigma_{3,2} & \sigma_3^2 & \sigma_{3,4} \\ \sigma_{4,1} & \sigma_{4,2} & \sigma_{4,3} & \sigma_4^2 \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \sigma_{1,2} & \sigma_{1,3} & \sigma_{1,4} \\ & \sigma_2^2 & \sigma_{2,3} & \sigma_{2,4} \\ & & \sigma_3^2 & \sigma_{3,4} \\ & & & \sigma_4^2 \end{bmatrix}$$

- May be something else



# Available covariance structures include

- Unknown: type=un
- Compound symmetry: type=cs
- Variance components: type=vc
- First-order autoregressive: type=ar(1)
- Spatial autocorrelation: covariance is a function of Euclidian distance
- Factor analysis
- Many others

# Compound Symmetry

- Why are data from the same case correlated?
- Because each case makes its own contribution -- add a (random) quantity that is different for each case.
- So variances of measurements are all equal.
- And correlations are all equal.
- Classical univariate approach implies compound symmetry.

# Compound Symmetry

$$\Sigma = \begin{bmatrix} \sigma^2 + \sigma_1 & \sigma_1 & \sigma_1 & \sigma_1 \\ \sigma_1 & \sigma^2 + \sigma_1 & \sigma_1 & \sigma_1 \\ \sigma_1 & \sigma_1 & \sigma^2 + \sigma_1 & \sigma_1 \\ \sigma_1 & \sigma_1 & \sigma_1 & \sigma^2 + \sigma_1 \end{bmatrix}$$

- Fewer parameters to estimate
- Implied by the random shock model.
- Not always realistic.

# Why not always assume covariance structure unknown?

- No reason why not, if you have enough data.
- Multivariate approach assumes  $\Sigma$  is completely unknown.
- When number of unknown parameters is large relative to sample size, variances of estimators are large  $\Rightarrow$  confidence intervals wide, tests weak.
- In some studies, there can be more treatment conditions than cases, and unique estimates of parameters don't even exist.
- There is always a tradeoff between assumptions and amount of data.

# First-order autoregressive time series

$$\Sigma = \sigma^2 \begin{bmatrix} 1 & \rho & \rho^2 & \rho^3 \\ \rho & 1 & \rho & \rho^2 \\ \rho^2 & \rho & 1 & \rho \\ \rho^3 & \rho^2 & \rho & 1 \end{bmatrix}$$

- Usually much bigger matrix
- Could have a handful of cases measured at hundreds of time points
- Or even just one “case,” say a company

# Eating Norm Study

- Two free meals at the psych lab (on different days)
- One with another student, one alone
- But it's not really another student. It's a "confederate."
- Confederate either eats a lot or a little.
- Dine with the confederate first, or second.
- Response variable is how much you eat. They weigh it.
- Covariates: How long since you ate, and how hungry you are. (Self Report)

# Variables

- Amount subject eats: Response variable
- Amount confederate eats (between)
- Eat alone or with confederate (within)
- Eat with confederate first, or second (between)
- Reported time since ate (covariate)
- Reported hunger (covariate)
  
- Notice these are **time-varying covariates**

# Multivariate approach can't handle time-varying covariates

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_k \end{pmatrix} = \begin{pmatrix} E(y_1 | \mathbf{X}=\mathbf{x}) \\ E(y_2 | \mathbf{X}=\mathbf{x}) \\ \vdots \\ E(y_k | \mathbf{X}=\mathbf{x}) \end{pmatrix} = \begin{pmatrix} \beta_{0,1} + \beta_{1,1}x_1 + \cdots + \beta_{p-1,1}x_{p-1} \\ \beta_{0,2} + \beta_{1,2}x_1 + \cdots + \beta_{p-1,2}x_{p-1} \\ \vdots \\ \beta_{0,k} + \beta_{1,k}x_1 + \cdots + \beta_{p-1,k}x_{p-1} \end{pmatrix}$$

Classical mixed model approach also fails.



# Copyright Information

This slide show was prepared by Jerry Brunner, Department of Statistical Sciences, University of Toronto. It is licensed under a Creative Commons Attribution - ShareAlike 3.0 Unported License. Use any part of it as you like and share the result freely. These Powerpoint slides are available from the course website:

<http://www.utstat.toronto.edu/~brunner/oldclass/441s20>