# Random Effects vs. Marginal Models: Different Approaches to Analyzing Repeated Measures / Longitudinal Data

Presented at
Midwest SAS Users Group

Kathy Welch
CSCAR, The University of Michigan
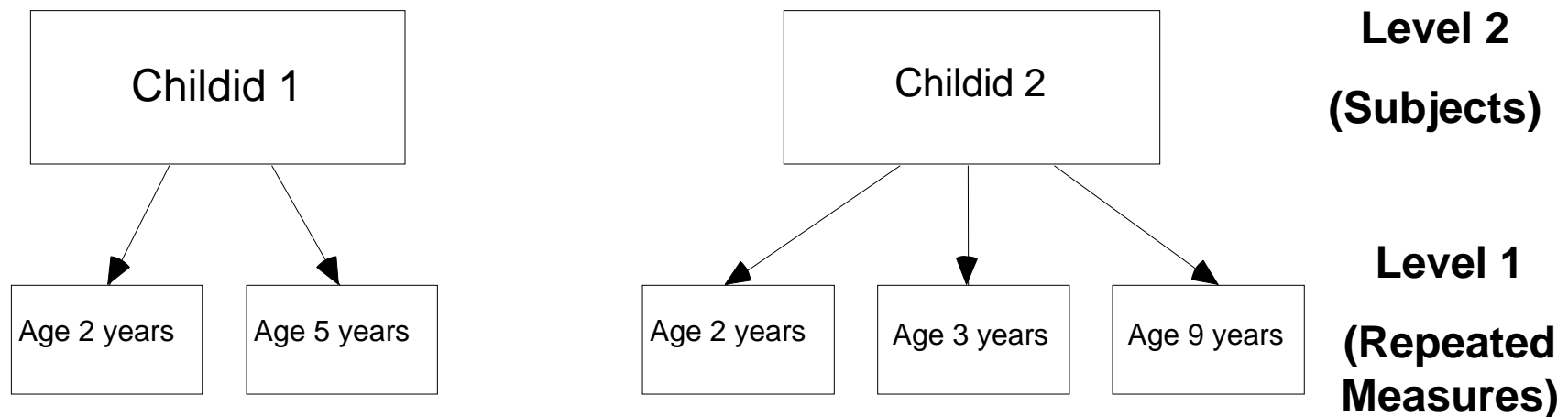February 5, 2009

# Background:
# What is a Linear Mixed Model (LMM)?

- A parametric linear model for
  - Clustered data
  - Repeated Measures / Longitudinal data
- Continuous response
- Predictors may be
  - Fixed
  - Random
- This presentation will focus on an analysis of a longitudinal data set.

# Repeated Measures / Longitudinal Data

- **Longitudinal Data:**
  - Dependent variable measured multiple times for each unit of analysis, basically a type of repeated measures data
  - Repeated measures factor is time
  - Time may be over an extended period (e.g. years)
- Example
  - Autistic children measured at different ages
- Dropout may be a problem
- Missing data at some time points may be a problem (not a problem if MAR)

# Example of Repeated Measures/ Longitudinal Data Structure

| Childid 1 | | Childid 2 | | | **Level 2**<br>**(Subjects)** |
|---|---|---|---|---|---|

| Age 2 years | Age 5 years | Age 2 years | Age 3 years | Age 9 years | **Level 1**<br>**(Repeated Measures)** |
|---|---|---|---|---|---|

- Each subject measured more than once

- Number of measurements does not need to be equal for all subjects

- Spacing of intervals does not have to be equal for all  measurement times

4

# Fixed Factor

- Fixed Factor: A categorical/classification variable
  - all levels of interest are included
    - Treatment level
    - Gender

- Levels of fixed factors can be defined to represent contrasts of interest
  - Female vs. Male
  - High Dose vs. Control, Medium Dose vs. Control

# Random Factor

- Random Factor: A classification variable
  - Levels can be thought of as being randomly sampled from a population
    - Classroom
    - Subject
- Variation in the dependent variable across levels of the random factor can be estimated and assessed
- Usually, random factors do not represent conditions chosen to meet the needs of the study
- Results can be generalized to a greater population

# Fixed **Effects**

- Also called regression coefficients or fixed-effect parameters
  - Describe the relationship between the dependent variable and predictor variables for an entire population

- Represented as unknown **fixed** quantities $(\beta)$ in a LMM
  - The value of a given $\beta$ does not vary across subjects

- $\beta$ is estimated based on data

# Random **Effects**

- Random values associated with levels of a random factor
- Represented as random variables ($u_i$ for the $i$th subject) in a LMM
  - Specific to a given level of a random factor
  - Vary across subjects
    - Classroom-specific intercepts in a clustered design
    - Subject-specific intercepts in a repeated measures design
- Usually describe random deviations in the relationships described by fixed effects
- Can be for categorical or continuous variables
  - Random intercepts
  - Random slopes

# General Specification of an LMM for the $i$th Subject:

$$Y_i = \underbrace{X_i \beta}_{\text{fixed}} + \underbrace{Z_i u_i + \varepsilon_i}_{\text{random}}$$

$$u_i \sim N(0, D)$$

$$\varepsilon_i \sim N(0, R_i)$$

Called a **mixed model** because it has a mix of fixed (*β)* and random (*$u_i$)* effects.

Both *D* and *$R_i$* are variance-covariance matrices, and as such, are required to be positive-definite

# The $D$ Matrix

Variance-covariance matrix for the $q$ random effects ($u_i$) for the ith subject. SAS calls this the $G$ matrix and defines it for all subjects, rather than for individuals.

$$D = Var(u_i) = \begin{pmatrix} Var(u_{1i}) & cov(u_{1i}, u_{2i}) & \cdots & cov(u_{1i}, u_{qi}) \\ cov(u_{1i}, u_{2i}) & Var(u_{2i}) & \cdots & cov(u_{2i}, u_{qi}) \\ \vdots & \vdots & \ddots & \vdots \\ cov(u_{1i}, u_{qi}) & cov(u_{2i}, u_{qi}) & \cdots & Var(u_{qi}) \end{pmatrix}$$

For Example: If there were only one random effect per subject (e.g., a random intercept), then $D$ would be a 1 X 1 matrix.
If there were two random effects per subject, e.g., a random intercept and a random slope, then $D$ would be 2 X 2.

# Two Common Structures for $D$

Many different structures for $D$ are possible:

**Variance components**

**type=vc**

$$D = Var(u_i) = \begin{pmatrix} \sigma^2_{u1} & 0 \\ 0 & \sigma^2_{u2} \end{pmatrix}$$

**Unstructured**

**type=un**

$$D = \begin{pmatrix} \sigma^2_{u1} & \sigma_{u1,u2} \\ \sigma_{u1,u2} & \sigma^2_{u2} \end{pmatrix}$$

Note: In these examples, we have two random effects defined for each subject. The diagonal elements represent variances of the random effects; the off-diagonal elements represent covariances between the random effects

# The $R$ Matrix

Variance-covariance matrix for the $n_i$ **residuals** ($\boldsymbol{\varepsilon}_i$) for the ith subject

$$\boldsymbol{R}_i = Var(\boldsymbol{\varepsilon}_i) = \begin{pmatrix} Var(\varepsilon_{1i}) & cov(\varepsilon_{1i}, \varepsilon_{2i}) & \cdots & cov(\varepsilon_{1i}, \varepsilon_{n_i i}) \\ cov(\varepsilon_{1i}, \varepsilon_{2i}) & Var(\varepsilon_{2i}) & \cdots & cov(\varepsilon_{2i}, \varepsilon_{n_i i}) \\ \vdots & \vdots & \ddots & \vdots \\ cov(\varepsilon_{1i}, \varepsilon_{n_i i}) & cov(\varepsilon_{2i}, \varepsilon_{n_i i}) & \cdots & Var(\varepsilon_{n_i i}) \end{pmatrix}$$

Note: The dimension of $\boldsymbol{R}_i$ depends on the number of observations ($n_i$) for subject $i$. For a subject with 5 repeated measures, the $\boldsymbol{R}_i$ matrix would be 5 X 5.

# Some Commonly Used Structures for $R$

**Unstructured**
**type = UN**

$$\begin{pmatrix} \sigma_1^2 & \sigma_{12} & \sigma_{13} \\ \sigma_{12} & \sigma_2^2 & \sigma_{23} \\ \sigma_{13} & \sigma_{23} & \sigma_3^2 \end{pmatrix}$$

**Variance**
**Components**
**type=VC**

$$\begin{pmatrix} \sigma_1^2 & 0 & 0 \\ 0 & \sigma_2^2 & 0 \\ 0 & 0 & \sigma_3^2 \end{pmatrix}$$

**Compound Symmetry**
**type = CS**

$$\begin{pmatrix} \sigma^2 + \sigma_1 & \sigma_1 & \sigma_1 \\ \sigma_1 & \sigma^2 + \sigma_1 & \sigma_1 \\ \sigma_1 & \sigma_1 & \sigma^2 + \sigma_1 \end{pmatrix}$$

**Banded**
**type = UN(2)**

$$\begin{pmatrix} \sigma_1^2 & \sigma_{12} & 0 \\ \sigma_{12} & \sigma_2^2 & \sigma_{23} \\ 0 & \sigma_{23} & \sigma_3^2 \end{pmatrix}$$

# More structures for $R$

**First-order Autoregressive**

**type = AR(1)**

$$\begin{pmatrix} \sigma^2 & \rho\sigma^2 & \rho^2\sigma^2 \\ \rho\sigma^2 & \sigma^2 & \rho\sigma^2 \\ \rho^2\sigma^2 & \rho\sigma^2 & \sigma^2 \end{pmatrix}$$

**Toeplitz**

**type = Toep**

$$\begin{pmatrix} \sigma^2 & \sigma_1 & \sigma_2 \\ \sigma_1 & \sigma^2 & \sigma_1 \\ \sigma_2 & \sigma_1 & \sigma^2 \end{pmatrix}$$

**Toeplitz (2)**

**type = Toep(2)**

$$\begin{pmatrix} \sigma^2 & \sigma_1 & 0 \\ \sigma_1 & \sigma^2 & \sigma_1 \\ 0 & \sigma_1 & \sigma^2 \end{pmatrix}$$

**Heterogeneous**
  **Compound Symmetry**

**type = CSH**

$$\begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 & \rho\sigma_1\sigma_3 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 & \rho\sigma_2\sigma_3 \\ \rho\sigma_1\sigma_3 & \rho\sigma_2\sigma_3 & \sigma_3^2 \end{pmatrix}$$

**Heterogeneous 1$^{st}$-order**
  **Autoregressive**

**type = ARH(1)**

$$\begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 & \rho^2\sigma_1\sigma_3 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 & \rho\sigma_2\sigma_3 \\ \rho^2\sigma_1\sigma_3 & \rho\sigma_2\sigma_3 & \sigma_3^2 \end{pmatrix}$$

**Heterogeneous Toeplitz**

**type = Toeph**

$$\begin{pmatrix} \sigma_1^2 & \rho_1\sigma_1\sigma_2 & \rho_2\sigma_1\sigma_3 \\ \rho_1\sigma_1\sigma_2 & \sigma_2^2 & \rho_1\sigma_2\sigma_3 \\ \rho_2\sigma_1\sigma_3 & \rho_1\sigma_2\sigma_3 & \sigma_3^2 \end{pmatrix}$$

# Covariance Parameters

- We estimate a set of covariance **parameters** for the variance-covariance matrices, $D$ and $R$.
  - For $D$ we estimate $\theta_D$
  - For $R$ we estimate $\theta_R$
- The number of covariance parameters that we estimate depends on the structure we specify for $D$ and $R$.

# Marginal Model vs. LMM

- **LMM** uses random effects explicitly to explain between-subject variance
  - Subject-specific model
- **Marginal model** does not use random effects in its specification at all
  - Population-averaged model
- **Implied marginal model**
  - Marginal model that results from fitting a LMM

# A Strictly Marginal Model
# With no random effects

$$Y_i = X_i \beta + \varepsilon_i^*$$

$$\varepsilon_i^* \sim N(0, V_i)$$

$$V_i = R_i$$

$V_i$ is the marginal variance-covariance matrix for $Y_i$

In this marginal model, we do not specify any random effects.

There is no $G$ matrix in this model.

Covariances, and hence correlations, among residuals
are specified directly through the $R_i$ matrix

# Implied Marginal Distribution of $Y_i$ Based on a LMM

$$Y_i \sim N(X_i \beta, \ Z_i D Z_i' + R_i)$$

$$E(Y_i) = X_i \beta$$

$$Var(Y_i) = V_i = Z_i D Z_i' + R_i.$$

In the **implied marginal model**, $V_i$ is formed from $D$ and $R_i$, but while $V_i$ is required to be positive-definite, $D$ and $R_i$ are not.

# Model Fit:
# Akaike Information Criteria (AIC)

- SAS calculates the AIC based on the (ML or REML) log likelihood, as shown below:

$$AIC = -2 \times l(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\theta}}) + 2p$$

- The penalty is 2p, where p represents the total number of parameters being estimated for both the fixed and random effects.

- Can be used to compare two models fit for the same observations, models need not be nested.

- Smaller is better.

# Model Fit:
# Bayes Information Criterion (BIC)

- BIC applies a greater penalty for models with more parameters than does AIC.

$$BIC = -2 \times l(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\theta}}) + p \times \ln(n)$$

- The penalty to the likelihood is number of parameters, p, times ln(n), where n is the total number of observations in the data set.
- Can be used to compare two models for the same observations, need not be nested.
- Smaller is better.

# Repeated Measures / Longitudinal Data Setup

- Data are in Long Form, one row for each repeated measurement on each subject

- Each row contains:
  - Information on the repeated measurements
    - Dependent variable
    - Time-varying covariates to be included in the model
  - Plus information on the subject / unit of analysis
    - Unit / subject ID
    - Time-invariant covariates to be included in the model
    - These are repeated for each row of data for a subject

# Proc Mixed Syntax

- **Model** statement specifies the fixed factors and covariates in the model

- **Random** statement specifies the random effects to be included in the model, and specifies the structure of the $D$ matrix of variances and covariances for the random effects (called $G$ matrix by SAS)

- **Repeated** statement specifies the structure of the residual covariance matrix, $R$

# The Autism Data Set

- **autism.csv**   This data set was derived from a study of 158 children with Autism Spectrum Disorder (Oti, Anderson, Lord, 2006).

- Measurements were made at five basic ages for each child: 2, 3, 5, 9, and 13 years. Not all children were measured at all time points.

- We will analyze VSAE, a measure of socialization, for these children as a function of their expressive skills (SICDEGP) measured at baseline (time invariant), and their current age (time-varying).

# Structure of Autism.csv data set

age,vsae,sicdegp,childid
2,6,3,1
3,7,3,1
5,18,3,1
9,25,3,1
13,27,3,1
2,17,3,3
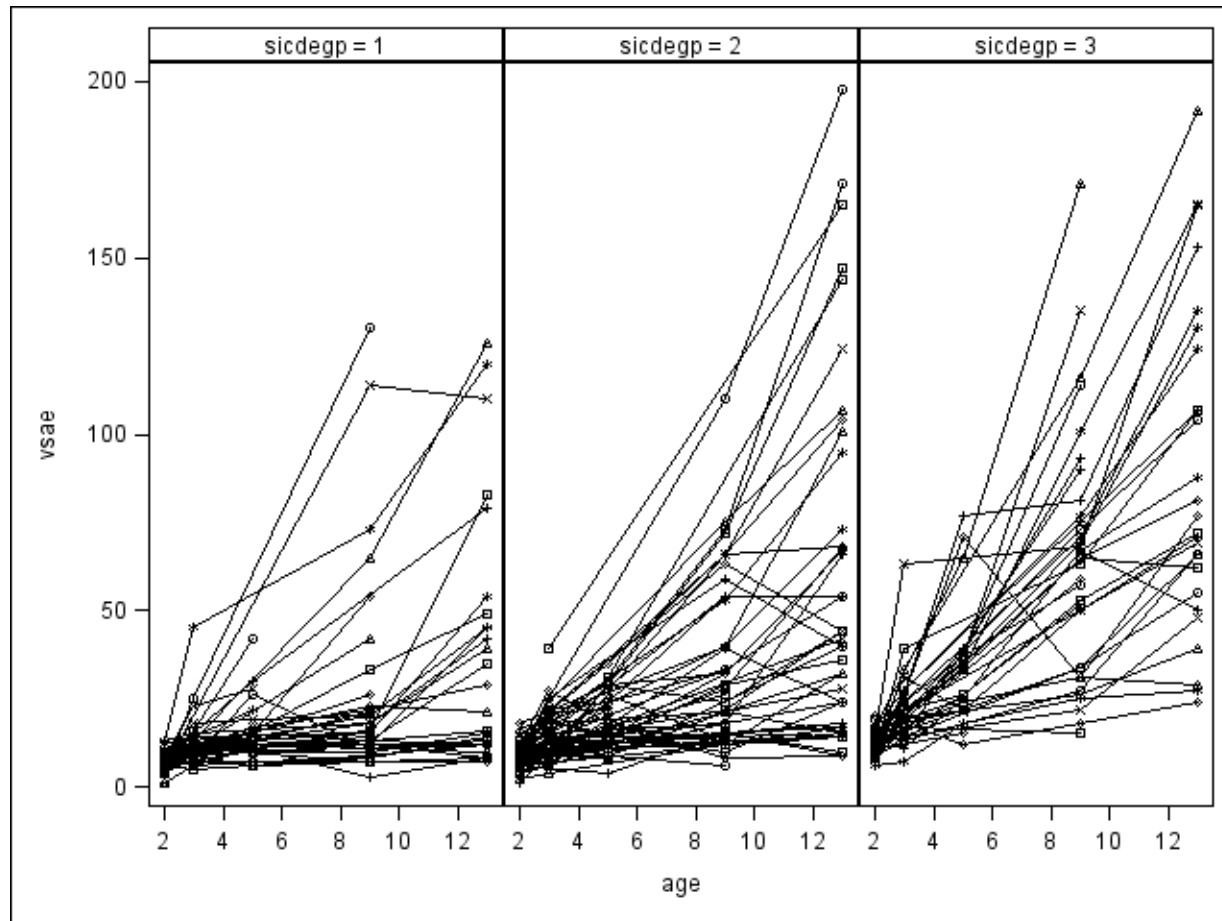3,18,3,3
5,12,3,3
9,18,3,3
13,24,3,3
2,12,3,4
3,14,3,4
5,38,3,4
9,114,3,4

| Obs | childid | age | sicdegp |
|-----|---------|-----|---------|
| 1   | 1       | 2   | 3       |
| 2   | 1       | 3   | 3       |
| 3   | 1       | 5   | 3       |
| 4   | 1       | 9   | 3       |
| 5   | 1       | 13  | 3       |
| 6   | 3       | 2   | 3       |
| 7   | 3       | 3   | 3       |
| 8   | 3       | 5   | 3       |
| 9   | 3       | 9   | 3       |
| 10  | 3       | 13  | 3       |
| 11  | 4       | 2   | 3       |
| 12  | 4       | 3   | 3       |
| 13  | 4       | 5   | 3       |
| 14  | 4       | 9   | 3       |

# Plots of VSAE Over Time for Each Child by Baseline Expressive Language Group
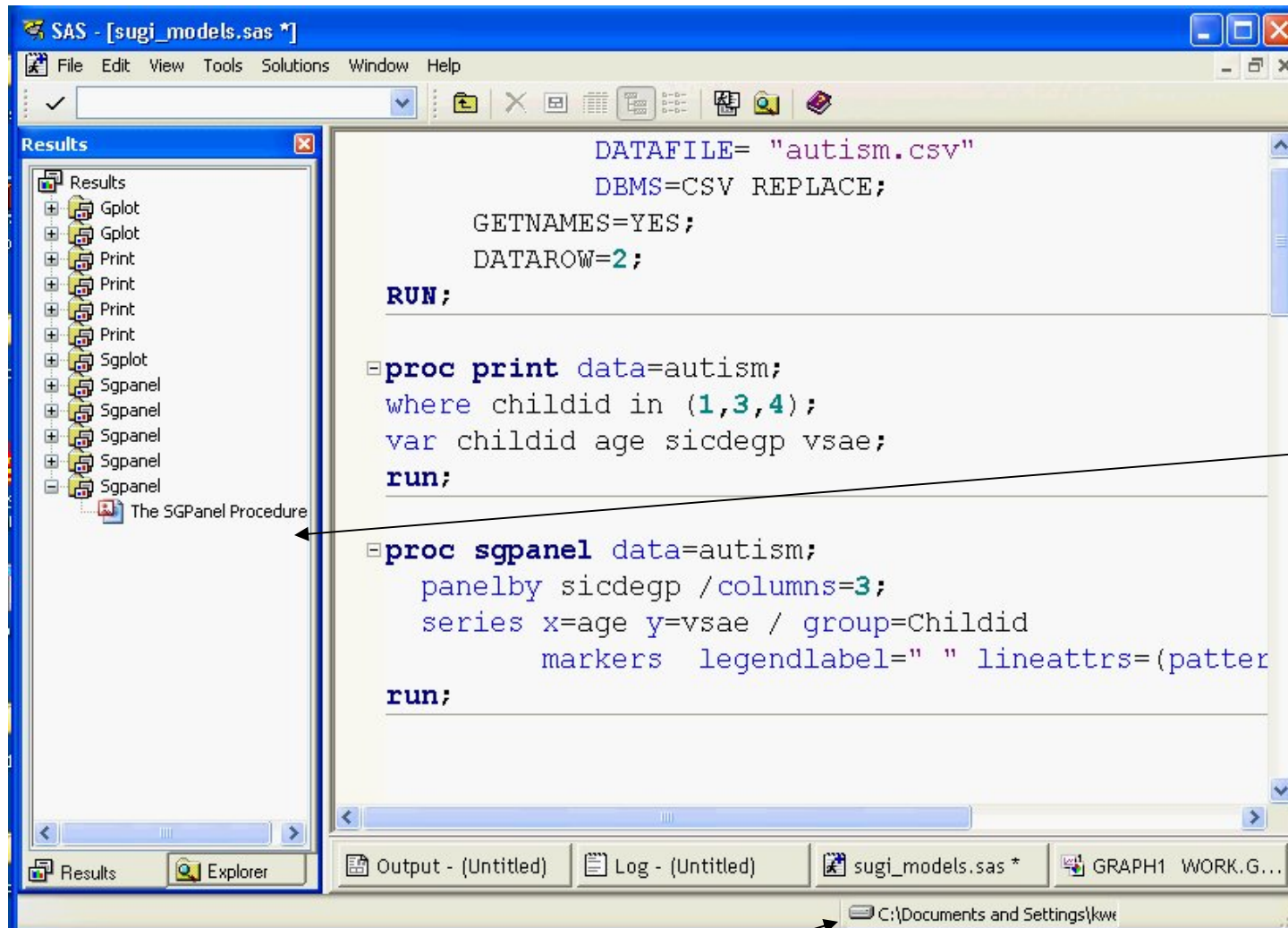
# Proc Sgpanel code for Individual line Graphs (SAS 9.2)

```
proc sgpanel data=autism;
  panelby sicdegp /columns=3;
  series x=age y=vsae / group=Childid
       markers  legendlabel=" " lineattrs=(pattern=1 color=black);
run;
```

The statistical graphics in SAS 9.2 are terrific. I'm still experimenting.

To get help, go to "SAS Help and Documentation"….SAS Products…SAS/Graph…Statistical Graphics Procedures.
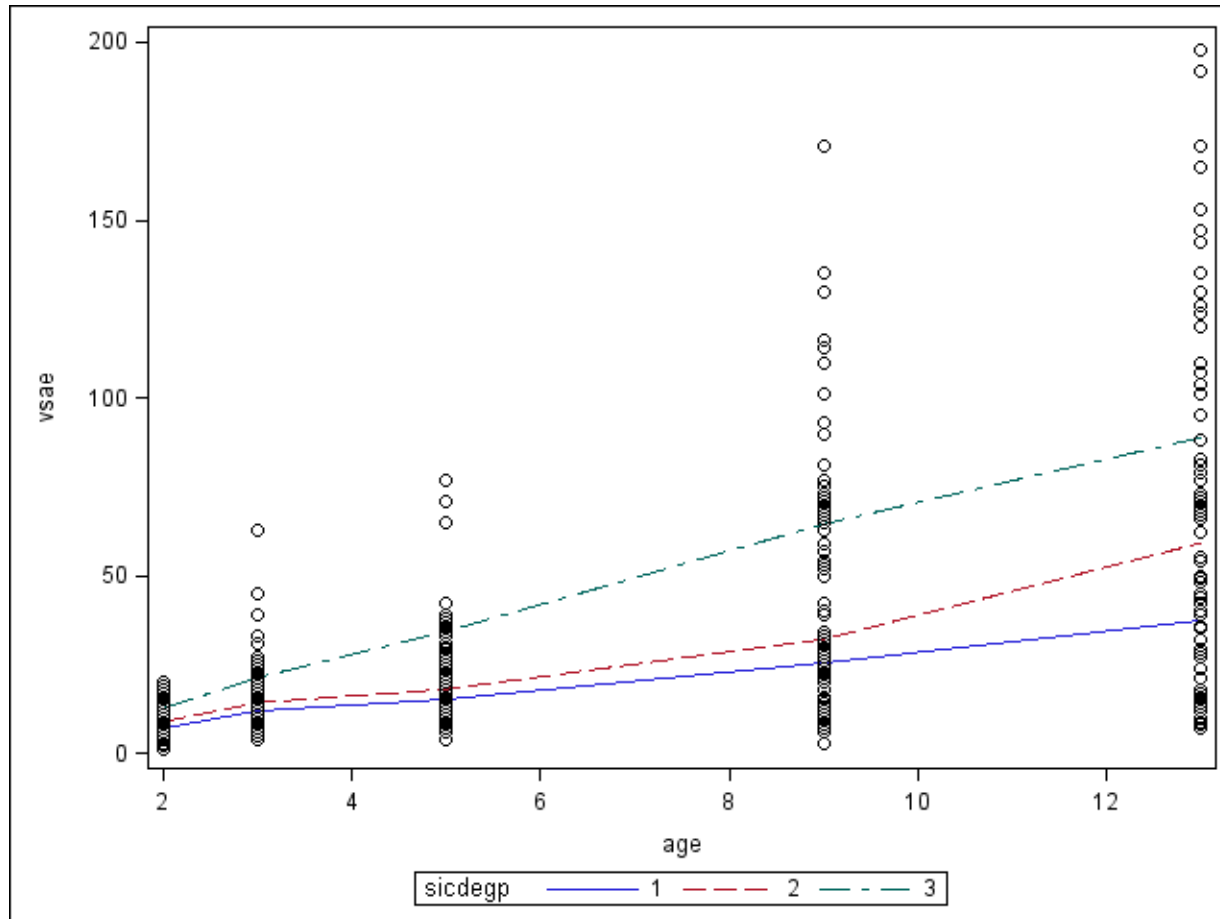
# Location of .png file from Sgpanel



Graph is in results window, not Graph window.

Graph will be saved in default folder as a .png file

# Mean Profiles by SICD Group

# SAS Sgplot Code for Mean Plots by SICD Group

```
proc sort data=autism;
by sicdegp age;
run;
proc means data=autism noprint;
  by sicdegp age;
  output out=meandat mean(VSAE)=mean_VSAE;
run;
data autism2;
  merge autism meandat(drop=_type_ _freq_);
  by sicdegp age;
run;
proc sgplot data=autism2;
  series x=age y=mean_VSAE / group=SICDEGP;
  scatter x=age y=VSAE ;
run;
```

# Discussion of Plots

- There is substantial variation in VSAE scores between children, and this variation gets larger over time.

- Although some children's scores do not seem to increase, there is a generally increasing trend in the means of VSAE over time in all three SICD groups.

- It looks like there may be a quadratic trend in mean VSAE scores, especially in group two.

# Modeling Strategy

- We first attempt to fit a LMM with 3 random effects for each subject: a random intercept, random slope for AGE, and random quadratic effect of AGE.

  - This is known as a random coefficients, growth-curve, or Laird-Ware Model

- We then fit an implied marginal model, in which we relax constraints on the variance-covariance matrices, $D$ and $R_i$

- Finally, we fit a strictly marginal model.

# LMM with Random Child-Specific Intercepts, Slopes and Quadratic Effects

**LMM:**

$$\text{VSAE}_{ti} = \beta_0 + \beta_1 \times \text{AGE\_2}_{ti} + \beta_2 \times \text{AGE\_2SQ}_{ti} + \beta_3 \times \text{SICDEGP1}_i$$

$$+ \beta_4 \times \text{SICDEGP2}_i + \beta_5 \times \text{AGE\_2}_{ti} \times \text{SICDEGP1}_i$$

$$+ \beta_6 \times \text{AGE\_2}_{ti} \times \text{SICDEGP2}_i + \beta_7 \times \text{AGE\_2SQ}_{ti} \times \text{SICDEGP1}_i$$

$$+ \beta_8 \times \text{AGE\_2SQ}_{ti} \times \text{SICDEGP2}_i + $$

$$\left. u_{0i} + u_{1i} \times \text{AGE\_2}_{ti} + u_{2i} \times \text{AGE\_2SQ}_{ti} + \varepsilon_{ti} \right\} \ \textbf{random}$$

**fixed**

We include the **fixed effects** of AGE, AGE-Squared, SICDEGP, and interactions between AGE, AGE-Squared and SICDEGP.

We also include three **random effects** for each child: the intercept ($u_{0i}$), the linear slope of AGE ($u_{1i}$), and the quadratic effect of AGE ($\boldsymbol{u}_{2i}$), to capture between-child variability.

32

# SAS Code for LMM

```
proc mixed data=autism2;
   class childid sicdegp;
   model vsae = age_2 age_2sq sicdegp age_2*sicdegp
         age_2sq*sicdegp / solution ddfm=sat;
   random int age_2 age_2sq
                     / subject=childid type=un g v ;
run;
```

# Distribution of Random Effects for the LMM

$$\boldsymbol{u}_i = \begin{pmatrix} u_{0i} \\ u_{1i} \\ u_{2i} \end{pmatrix} \sim N(\boldsymbol{0}, \boldsymbol{D}).$$

$$\boldsymbol{D} = \begin{pmatrix} \sigma^2_{int} & \sigma_{int,age} & \sigma_{int,age\text{-}squared} \\ \sigma_{int,age} & \sigma^2_{age} & \sigma_{age,age\text{-}squared} \\ \sigma_{int,age\text{-}squared} & \sigma_{age,age\text{-}squared} & \sigma^2_{age\text{-}squared} \end{pmatrix}, \qquad \varepsilon_{ti} \sim N(0, \sigma^2)$$

We specify an unstructured **D** matrix for the random effects. There are 6 covariance parameters in the **D** matrix for the 3 random effects.

Note: There is no **R** matrix specified for this model, so **R** is assumed to be $\sigma^2\boldsymbol{I}$

# LMM: Problem with *G* Matrix

SAS reports problems fitting Model 6.1. We see the following note in the SAS log:

**NOTE: Convergence criteria met.**

**NOTE: Estimated G matrix is not positive definite.**

**NOTE: Asymptotic variance matrix of covariance parameter estimates has been found to be singular and a generalized inverse was used. Covariance parameters with zero variance do not contribute to degrees of freedom computed by DDFM=SATTERTH.**

**We have a problem with the *G* matrix (referred to as the *D* matrix in this presentation). We need to investigate this problem.**

# LMM: SAS Output for *G* Matrix

We see that the value in the G matrix corresponding to the variance of the random intercepts is blank here.

| | | | Estimated G Matrix | | |
|---|---|---|---|---|---|
| Row | Effect | childid | Col1 | Col2 | Col3 |
| 1 | Intercept | 1 | | 0.6171 | 0.5669 |
| 2 | age_2 | 1 | 0.6171 | 14.0300 | -0.6353 |
| 3 | age_2sq | 1 | 0.5669 | -0.6353 | 0.1664 |

# Fit the Implied Marginal Model

- We now refit the model

- Use the **nobound** option, to get the implied marginal model,

- Positive definiteness constraints on *G* and *R* are relaxed.

proc mixed data=autism2  **nobound**;

# Look at Unconstrained *G* Matrix

| Estimated G Matrix | | | | | |
|---|---|---|---|---|---|
| **Row** | **Effect** | **childid** | **Col1** | **Col2** | **Col3** |
| **1** | Intercept | 1 | -10.5406 | 4.2760 | 0.1423 |
| **2** | age_2 | 1 | 4.2760 | 11.9673 | -0.4038 |
| **3** | age_2sq | 1 | 0.1423 | -0.4038 | 0.1383 |

The new estimate of the variance of the random intercepts is **negative**!

Clearly, this LMM is not working, as a negative variance is impossible!

38

# Revised LMM: Remove the Random Intercept

```
proc mixed data = autism2;
   class childid sicdegp;
   model vsae = sicdegp age_2 age_2sq age_2*sicdegp
            age_2sq*sicdegp /
            solution ddfm=sat influence;
   random age_2 age_2sq /
                           subject = childid solution g v vcorr type = un;
run;
```

**Note: the only change in the model is that "int" has been deleted from the random statement.**

# G Matrix for Revised LMM without Random Intercept

There are no error messages in the log.

The 2x2 G matrix is positive-definite.

| Estimated G Matrix | | | | |
|---|---|---|---|---|
| Row | Effect | childid | Col1 | Col2 |
| 1 | age_2 | 1 | 14.6674 | -0.4401 |
| 2 | age_2sq | 1 | -0.4401 | 0.1315 |

| Covariance Parameter Estimates | | |
|---|---|---|
| Cov Parm | Subject | Estimate |
| UN(1,1) | childid | 14.6674 |
| UN(2,1) | childid | -0.4401 |
| UN(2,2) | childid | 0.1315 |
| Residual | | 38.4988 |

# *V* Matrix for the Revised LMM

The only variability in the intercept is the estimated residual variance.

There is no covariance between the Y values at baseline with any other ages

| Estimated V Matrix for childid 1 | | | | | |
|---|---|---|---|---|---|
| Row | Col1 | Col2 | Col3 | Col4 | Col5 |
| 1 | 38.4988 | | | | |
| 2 | | **52.4175** | 39.9043 | 84.4687 | 119.16 |
| 3 | | 39.9043 | **157.39** | 273.57 | 423.87 |
| 4 | | 84.4687 | 273.57 | **770.95** | 1298.89 |
| 5 | | 119.16 | 423.87 | 1298.89 | **2566.53** |

41

# Alternative Marginal Model

Proc Mixed syntax for the marginal model:

```
proc mixed data=autism2 noclprint;
   class childid sicdegp age;
   model vsae = sicdegp age_2 age_2sq age_2*sicdegp
      age_2sq*sicdegp
      / solution ddfm=bw;
   repeated age / subject=childid type=un r rcorr;
run;
```

Note: there is no random statement in this model, and hence, the **G** matrix = **0**.

The repeated statement specifies that the **R** matrix should be unstructured.

# Marginal Model: *R* Matrix

Note: This *R* matrix shows positive covariance, and hence, positive correlation, between residuals at baseline and later ages.

| Estimated R Matrix for childid 1 | | | | | |
|---|---|---|---|---|---|
| Row | Col1 | Col2 | Col3 | Col4 | Col5 |
| 1 | **10.8223** | **8.8584** | **8.3251** | **24.6243** | **43.3342** |
| 2 | **8.8584** | **54.5604** | 55.5313 | 94.3265 | 186.88 |
| 3 | **8.3251** | 55.5313 | **141.39** | 194.08 | 368.95 |
| 4 | **24.6243** | 94.3265 | 194.08 | **810.49** | 1153.08 |
| 5 | **43.3342** | 186.88 | 368.95 | 1153.08 | **2181.67** |

Again, the variance increases at each time point, as was apparent in the initial graphs, and in the LMM.

# Model Fit Comparison

| | Full LMM | Revised LMM minus random intercept | Marginal Model |
|---|---|---|---|
| AIC | 4616.7 | 4623.3 | 4459.5 |
| BIC | 4635.1 | 4635.5 | 4505.4 |

From this comparison, the marginal model is preferable to the two LMMs for this data set. It has the smallest AIC and BIC.

# Summing Up

- The ability to fit a wide array of different covariance structures gives Proc Mixed a lot of flexibility

- By examining the *G* matrix, *R* matrix and *V* matrix we can see how different structures affect the model covariance parameters.

- A LMM may not always be the best choice.

- At times, a marginal model may give a better fit.

- Research goals can help in the choice of the "right" modeling approach.

# References

- Brady West, Kathleen Welch and Andrzej Galecki, "Linear Mixed Models: A Practical Guide Using Statistical Software", Chapman & Hall/CRC, 1986, 353 pp.
- Oti, R., Anderson, D., and Lord, C. Social trajectories among individuals with autism spectrum disorders, Journal of Developmental Psychopathology.