

Dominic Adduci

Homework 6 BIOS 6643

Question 1

Below are two spaghetti plots of these individuals using two different x-axes. Assess the between and within subject variation in 3 sentences or less.

Left Plot

There is a noticeable amount of variation between subjects, with subjects covering a wide range at specific timepoints. There also seems to be a decent amount of variation within subjects. As time progresses subjects have a decrease along the y-axis.

Right Plot

There is a lot of variation between subjects, with each timepoint having a wide range between subjects. There is less variation within a subject, with not much change along the y-axis for a subject as time changes.

Question 2

You decide to model these data using a linear mixed effects model with a linear trend for observation year and use a marginal compound symmetry model to incorporate correlation between the cognitive measures on the same person. Write out the individual (subject) level mixed effects model for this model. Define all matrix dimensions and indices. Include the distribution of model error and any random effects components including the structure of their variance terms.

$$Y_i = \overset{\text{Fixed}}{X_i \beta} + \epsilon_i$$

Because we are using a compound symmetry model we are not including a random effect. This model will be equivalent to a LMM with a random intercept.

Age of 0 would be the reference.

$Y_i = (n_i \times 1)$; n_i is the number of observation years.

$X_i = (n_i \times p)$; where n_i is the number of observation years and $p=2$ because of the fixed intercept and the effect of linear time. For subjects with less than 8, n_i will be < 8 .

$\beta = (p \times 1)$; There is only 1 covariate, age, so $p=2$.

$\epsilon_i = (n_i \times 1)$; n_i is the number of observation years for a subject.

$i = 1, 2, \dots, n$, where n is the number of subjects.

$$\epsilon_i \sim \text{MVP}(0, R_i),$$

$$\text{where } R_i = \begin{bmatrix} \sigma^2 + \sigma_{12} & \dots & \sigma_{12} \\ \vdots & \ddots & \vdots \\ \sigma_{12} & \dots & \sigma^2 + \sigma_{12} \end{bmatrix}$$

Question 3

Show that the minimum of the GEE objective function $\sum_{i=1}^N [y_i - \mu_i(\beta)]^T V_i^{-1} [y_i - \mu_i(\beta)]$ is $\sum_{i=1}^N D_i^T V_i^{-1} [y_i - \mu_i(\beta)] = 0$

$$\sum_{i=1}^N (y_i - \mu_i(\beta))^T V_i^{-1} (y_i - \mu_i(\beta)) = \sum_{i=1}^N (y_i^T - \mu_i(\beta)^T) V_i^{-1} (y_i - \mu_i(\beta))$$

Moving V_i^{-1} into first set of parenthesis

$$\sum_{i=1}^N (y_i^T V_i^{-1} - \mu_i^T(\beta) V_i^{-1}) (y_i - \mu_i(\beta))$$

Then multiplying the two parenthesis together

$$\sum_{i=1}^N (y_i^T V_i^{-1} y_i - \mu_i^T(\beta) V_i^{-1} y_i - y_i^T V_i^{-1} \mu_i(\beta) + \mu_i^T(\beta) V_i^{-1} \mu_i(\beta))$$

Now evaluating all the terms:

$$- y_i^T V_i^{-1} \mu_i(\beta) \rightarrow (1 \times N)(N \times N)(N \times 1) = -1$$

$$- \mu_i^T(\beta) V_i^{-1} y_i \rightarrow (1 \times N)(N \times N)(N \times 1) = -1$$

$$(\mu_i^T(\beta) V_i^{-1} y_i)^T = y_i^T (V_i^{-1})^T \mu_i(\beta) = y_i^T V_i^{-1} \mu_i(\beta)$$

This works because V_i^{-1} is a symmetric matrix, meaning $(V_i^{-1})^T = V_i^{-1}$

This then evaluates to:

$$\sum_{i=1}^N (y_i^T V_i^{-1} y_i - 2 \mu_i^T(\beta) V_i^{-1} y_i + \mu_i^T(\beta) V_i^{-1} \mu_i(\beta))$$

Taking the first derivative and setting to 0 will find the minimum.

$$\frac{\partial}{\partial \beta} \sum_{i=1}^N (y_i^T V_i^{-1} y_i - 2 \mu_i^T(\beta) V_i^{-1} y_i + \mu_i^T(\beta) V_i^{-1} \mu_i(\beta)) = 0$$

$$\sum_{i=1}^N (0 - 2 [\mu_i^T(\beta)]' V_i^{-1} y_i + 2 [\mu_i(\beta)]^T V_i^{-1} \mu_i(\beta)) = 0$$

$$\sum_{i=1}^N ([\mu_i(\beta)]^T V_i^{-1} [y_i - \mu_i(\beta)]) = 0$$

$$\Rightarrow \sum_{i=1}^N D_i^T V_i^{-1} (y_i - \mu_i(\beta)) = 0$$

$$\text{where } D_i = \mu_i(\beta)' = \frac{\partial}{\partial \beta} \mu_i(\beta)$$

Dominic Adducci Homework 6 BIOS 6643

Question 4

Example 1: We received electronic health records data on 90,000 patients who had a positive test result or were administered monoclonal antibodies (mAbs) for COVID treatment in the CU Health system. All patients were eligible to receive mAbs treatment. The investigator was interested in whether mAbs treatment reduced 28-day hospitalization after adjusting for important precision and confounding variables of age, insurance status, and number of comorbid conditions.

A) Approach

I interpreted the investigator's question as wanting to know if treatment was less than 28 days if mAbs treatment was used. The precision and confounding variables of age, insurance status, and number of comorbid conditions are not longitudinal. mAbs is also not longitudinal, as patients either received mAbs or they did not. A binomial generalized linear model will be used to determine if the probability of 28-day hospitalization is reduced with mAbs treatment. A logit link will be used and the previously noted covariates will be controlled for.

B) Outcome

The outcome for this is the odds ratio of 28 day hospitalization.

C) Distributional assumption for this model

The distributional assumption is a Bernoulli distribution for the outcome.

D) Write out the distribution of the outcome with correct subscripts and definitions of indices. Write out the expected value and variances for this distribution.

$Y_i \sim \text{Bernoulli}(\mu_i)$; i refers to the subject. μ_i is probability of < 28 day hospitalization.

$$E[Y_i] = \mu_i ; \text{Var}[Y_i] = \mu_i(1 - \mu_i)$$

E) What is the linear model for these data (i.e. the systematic component)? Be specific using the covariates you were given.

$$\log\left(\frac{\mu_i}{1 - \mu_i}\right) = \eta_i = \overbrace{\beta_0 + \beta_1 \text{ mAb} + \beta_2 \text{ Age} + \beta_3 \text{ Insurance Status} + \beta_4 \text{ Number of Comorbidities}}^{\text{systematic component}}$$

F) What is the link function?

$$\eta_i = \underbrace{\log\left(\frac{\mu_i}{1 - \mu_i}\right)}_{\text{Link Function (logit)}}$$

G) How will you interpret the coefficient on treatment for your model? Write a formal statistical interpretation.

The mAb coefficient is the log odds of less than 28 days of hospitalization if treatment is received. The exponentiated coefficient will give the odds of hospitalization within 28 days. 95% CI will give estimates of the bounds and a p-value will assess significance. Will you use?

H) How will this model be estimated?

The coefficients will be estimated using maximum likelihood.

Domini Adduci Homework 6 BIOS 6643

Question 4 [Continued]

I) Write pseudo R code for estimating your model.

Less-Than-28 = glm(Days-28 ~ mAbst + Age + Insurance-Status + Number-of-Comorbs,
family = binomial, data = Electronic-Health-Data)

Domestic Address Homework 6 BIOS 6643

Question 5

Example 2: We received electronic health records data on 90,000 patients who had a positive test result or were administered monoclonal antibodies (mAbs) for COVID treatment in the CH health system. All patients were eligible to receive mAbs treatment. The investigator was interested in whether mAbs treatment reduced the number of emergency department (ED) visits after adjusting for important precision and confounding variables of age, insurance status, and number of comorbid conditions.

A) Approach

In this analysis we want to determine if the number of ED visits (counts) were reduced if mAbs treatment was utilized. The precision and confounding variables of age, insurance status, and number of comorbid conditions are not longitudinal. mAbs is also not longitudinal, as patients either received mAbs or they did not. A Quasi-Poisson generalized linear model will be determined if ED visits will be reduced with mAbs treatment. The number of ED visits will be the outcome, and the previously mentioned covariates will be controlled for.

B) Outcome

The number of ED visits will be the outcome.

C) Distributional assumption for this model

The distributional assumption is Poisson. A Quasi-Poisson model will account for overdispersion.

D) Write out the distribution of the outcome with correct subscripts and definitions of indices. Write out the expected value and variance for this distribution.

$Y_i \sim \text{Poisson}(\lambda_i)$; i refers to the subject, λ_i is the rate of hospitalization.

$E[Y_i] = \lambda_i$; $\text{var}[Y_i] = \phi \lambda_i$; ϕ accounts for overdispersion.

E) What is the linear model for these data (i.e. the systematic component)?

$$Y_i = \beta_0 + \beta_1 \text{mAbs}_i + \beta_2 \text{Age}_i + \beta_3 \text{Insurance Status}_i + \beta_4 \text{Number of comorbid cond.}_i$$

F) What is the link function?

$$\eta_i = \log(\lambda_i)$$

G) How will you interpret the coefficient on treatment for your model?

The mAbs coefficient is the log of rate for someone who received mAbs treatment versus someone who did not receive the treatment. Less than 1 will mean a lower number of ED visits. The exponentiated coefficient will show the difference rate of ED visits between the treatment and control group. 95% CI will give the bounds, and a p-value will assess significance.

H) How will this model be estimated?

With a generalized linear model, we can estimate the parameters of the distribution.

The coefficients will be estimated using maximum likelihood.

Log-likelihood function for the generalized linear model, and the maximum likelihood estimates will be used to estimate the parameters.

Domènec Albués Homework 6 BIOS 6643

Question 3 [continued]

I.) write pseudo R code for estimating your model.

ED-visits_treatment = glm(ED-visits ~ mAbs + Age + Insurance-Status + Number-of-Comorbs,
family = quasi poisson, data = Electronic-Health-Data)

Dominic Alucci Homework 6 BIOS 6643

Question 6

Example 3: We received electronic health records data on 90,000 patients who had a positive test result or were administered monoclonal antibodies (mAbs) for COVID treatment in the CU Health system who were followed every 3 months for a year (so baseline (first 28 days), 3 mo, 6 mo, 9 mo, and 12 mo). All patients were eligible to receive mAbs treatment. The investigator was interested in whether mAbs treatment modifies the number and pattern of post COVID doctor visits over the year after COVID. We need to adjust for important precision and confounding variables of age, insurance status, and number of comorbid conditions at baseline.

A) Approach

This analysis has a longitudinal component where subjects are assessed every 3 months. Each subject will have multiple observations, where the primary variable will be mAbs treatment and confounding variables will be age, insurance status, time, and number of comorbid conditions. A marginal GEE will be used to determine if the subjects who receive mAbs treatment have improvements in required doctor visits over time. The mAbs variable will assess the treatment effect, and the time variable will assess the pattern of visits. The GEE will use a log link.

B) Outcome

The outcome will be the number of doctor visits (count).

C) Distributional assumption for this model

A GEE does not require a distributional assumption, but in this case a log link function will be used which will model the results as a Poisson.

D) Write out the distribution of the outcome with correct subscripts and definitions of indices. Write out the expected value and variance for this distribution and any cov/corr struct.

$E[Y_{ij}] = \lambda_{ij}$; where λ_{ij} is the rate parameter, i is the subject, and j is the timepoint.

$\text{Var}[Y_{ij}] = \phi \lambda_{ij}$, ϕ accounts for overdispersion.

$\text{Corr}(Y_{ij}, Y_{ik}) = \alpha_{jk}$; Compound symmetry like relationship.

E) What is the linear model for these data (i.e. systematic component)?

$$\eta_{ij} = \beta_0 + \beta_1 \text{mAbs}_{ij} + \beta_2 \text{Time}_{ij} + \beta_3 \text{Age}_{ij} + \beta_4 \text{Insurance Status}_{ij} + \beta_5 \text{Number of comorbidities}_{ij}$$

F) What is the link function?

$$\eta_{ij} = \log(\lambda_{ij})$$

G) How will you interpret the coefficient of treatment on your model?

The coefficient of treatment will be the log rate ratio of someone who received the mAbs treatment versus someone who did not receive the treatment. The rate in this case will be the rate of doctors visits. A negative coefficient will mean mAbs treatment reduces doctors visits.

H) How will this model be estimated?

The exponentiated treatment coefficient will assess the effect of mAbs on the number of doctors visits. A rate ratio < 1 will mean mAbs reduces the number of doctors visits. A 95% CI will assess the range of the effect and a p-value will assess significance. The exponentiated time coefficient will assess the change in doctors visits over time, where 95% CI will show the range of effect and the p-value will determine if the effect is significant.

Domènec Altaba: Homework 6 BIOS 4643

Question 6 [Continued]

I) Write pseudo R code for estimating your model.

library(gee)

```
mftbs_treat = gee(Doctors_visits ~ mftbs + Time + Age + Insurance_status +  
  Number_of_comorbidities, Family = quasipoisson,  
  Corstr = unstructured, data = Electronic_Health_Data)
```



Dominic Ajjuci Homework 7 BIOS 6643

Question 7

Example 4: We received brain volume metrics and biomarker data from a neuropsychologist who is interested in how the biomarker influences the trajectory of brain volumes over 5 years in those with mild cognitive impairment. Patients received an MRI each year and different brain volumes were computed. For this project the investigator is interested in the hippocampus (one brain region implicated in Alzheimer's). The biomarker data was only measured at baseline. Hippocampal volume can be considered a continuous measure with an infinite number of values possible, meaning it isn't count or binary data. All analysis will be adjusted for baseline age, sex, and SES.

A) Pseudo code (Provided)

```
lmm.un.fit <- gls(hippoc ~ visit + visit^2 + biomarker + visit * biomarker + visit^2 * biomarker +  
age + sex + SES, data = braindata, correlation = corSymm(form = 1 | ID),  
weights = varIdent(form = 1 | visit))
```

```
lmm.un.summary <- summary(lmm.un.fit)
```

```
lmm.un.fit2 <- gls(hippoc ~ as.factor(visit) + biomarker + as.factor(visit) * biomarker + age +  
sex + SES, data = braindata, correlation = corSymm(form = 1 | ID),  
weights = varIdent(form = 1 | as.factor(visit)))
```

B) Approach

Both models fit a linear regression using the generalized least squares method. For the first model visit is considered a continuous variable and there is a square visit term to account for a non-linear relationship between visit and hippocampus volume. An interaction between the biomarker variable and the visit and visit squared term was included in the model. An unstructured correlation for each subject was assumed and different variance were assumed for different visits. The second model has many of the same attributes except that the visit variable is factored. There is still an interaction between factored visit and the biomarker and the weights parameter assigned different variances for different visits.

C) What is the outcome for these analyses?

The outcome in both models is hippocampus volume.

D) What is the distributional assumption you will make for this outcome?

The GLS model makes no assumptions on the distribution of the outcome, but does assume the residuals are distributed normally with constant variance for each visit group.

E) How is time model in each of these models? How do you determine this?

For the first model time is a continuous variable. I determined this by the fact that there is a squared time term. This model includes a visit term, a visit squared term, as well as an interaction between both of these terms and the biomarker variable. For the second model time is Factorial making it discrete. This was simple to determine, as a Factorial variable is inherently discrete. There is also an interaction in this model between factored time and the biomarker.

F) Write out the distribution of the outcome with correct subscripts and definitions of indices. Write out the expected value and variance-covariance for this distribution for each of these models

Model 1 expected value:

$$E[Y_{ij}] = \beta_0 + \beta_1 \text{visit}_{ij} + \beta_2 \text{visit}_{ij}^2 + \beta_3 \text{biomarker}_{ij} + \beta_4 \text{visit}_{ij} * \text{biomarker}_{ij} + \beta_5 \text{visit}_{ij}^2 * \text{biomarker}_{ij} + \beta_6 \text{age}_{ij} + \beta_7 \text{sex}_{ij} + \beta_8 \text{SES}_{ij}$$

where i is the subject; $i = 1, \dots, n$ n = number of subjects

where j is the observation; $j = 1, \dots, n_i$ n_i = number of observations
for a subject

Domènec Aladell Homework 6 BIOS 6643

Question 7 [Continued]

F (Continued)

Model 2 expected value (Visit 1 reference)

$$E[Y_{ij}] = \beta_0 + \beta_1 \text{Visit}2_{ij} + \beta_2 \text{Visit}3_{ij} + \beta_3 \text{Visit}4_{ij} + \beta_4 \text{Visit}5_{ij} + \beta_5 \text{Visit}6_{ij} + \beta_6 \text{biomarker}_{ij} + \beta_7 \text{Visit}2 * \text{biomarker}_{ij} + \beta_8 \text{Visit}3 * \text{biomarker}_{ij} + \beta_9 \text{Visit}4 * \text{biomarker}_{ij} + \beta_{10} \text{Visit}5 * \text{biomarker}_{ij} + \beta_{11} \text{age}_{ij} + \beta_{12} \text{Sex}_{ij} + \beta_{13} \text{SES}_{ij}$$

$$\text{Variance-Covariance} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \sigma_{13} & \sigma_{14} & \sigma_{15} \\ \sigma_{21} & \sigma_2^2 & \sigma_{23} & \sigma_{24} & \sigma_{25} \\ \sigma_{31} & \sigma_{32} & \sigma_3^2 & \sigma_{34} & \sigma_{35} \\ \sigma_{41} & \sigma_{42} & \sigma_{43} & \sigma_4^2 & \sigma_{45} \\ \sigma_{51} & \sigma_{52} & \sigma_{53} & \sigma_{54} & \sigma_5^2 \end{bmatrix} = R_i \quad \text{where } \sigma_{12} = \sigma_{21}; \sigma_{13} = \sigma_{31}, \text{etc.}$$

Symmetric var-cov structure

R_i is the same for both models

c) What is the linear model for these data in observation level, subject level, and complete data notation for the first model.

Observation level:

$$Y_{ij} = \beta_0 + \beta_1 \text{Visit}_{ij} + \beta_2 \text{Visit}2_{ij} + \beta_3 \text{biomarker}_{ij} + \beta_4 \text{Visit} * \text{biomarker}_{ij} + \beta_5 \text{Visit}2 * \text{biomarker}_{ij} + \beta_6 \text{age}_{ij} + \beta_7 \text{Sex}_{ij} + \beta_8 \text{SES}_{ij} + \epsilon_{ij}$$

where $\epsilon_{ij} \sim N(0, R_i)$; R_i defined above

Subject level:

$$\begin{bmatrix} y_{i1} \\ y_{i2} \\ y_{i3} \\ y_{i4} \\ y_{i5} \end{bmatrix} = \begin{bmatrix} 1 & \text{visit}_{i1} & \text{visit}_{sq_{i1}} & \text{biomarker}_{i1} & \text{visit} * \text{biomarker}_{i1} & \text{visit}_{sq} * \text{biomarker}_{i1} \\ 1 & \text{visit}_{i2} & \text{visit}_{sq_{i2}} & \text{biomarker}_{i2} & \text{visit} * \text{biomarker}_{i2} & \text{visit}_{sq} * \text{biomarker}_{i2} \\ 1 & \text{visit}_{i3} & \text{visit}_{sq_{i3}} & \text{biomarker}_{i3} & \text{visit} * \text{biomarker}_{i3} & \text{visit}_{sq} * \text{biomarker}_{i3} \\ 1 & \text{visit}_{i4} & \text{visit}_{sq_{i4}} & \text{biomarker}_{i4} & \text{visit} * \text{biomarker}_{i4} & \text{visit}_{sq} * \text{biomarker}_{i4} \\ 1 & \text{visit}_{i5} & \text{visit}_{sq_{i5}} & \text{biomarker}_{i5} & \text{visit} * \text{biomarker}_{i5} & \text{visit}_{sq} * \text{biomarker}_{i5} \end{bmatrix}$$

$$\begin{bmatrix} \text{age}_{i1} & \text{sex}_{i1} & \text{SES}_{i1} \\ \text{age}_{i2} & \text{sex}_{i2} & \text{SES}_{i2} \\ \text{age}_{i3} & \text{sex}_{i3} & \text{SES}_{i3} \\ \text{age}_{i4} & \text{sex}_{i4} & \text{SES}_{i4} \\ \text{age}_{i5} & \text{sex}_{i5} & \text{SES}_{i5} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \\ \beta_6 \\ \beta_7 \\ \beta_8 \\ \beta_9 \end{bmatrix} + \begin{bmatrix} \epsilon_{i1} \\ \epsilon_{i2} \\ \epsilon_{i3} \\ \epsilon_{i4} \\ \epsilon_{i5} \end{bmatrix} \quad \epsilon_{ij} \sim N(0, \sigma_{ij})$$

Complete data notation

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \\ \beta_6 \\ \beta_7 \\ \beta_8 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix} \quad \begin{matrix} \epsilon_i \sim N(0, \begin{bmatrix} \sigma_{i1} & 0 \\ 0 & \ddots \\ 0 & 0 & \sigma_{in} \end{bmatrix}) \\ \downarrow \\ (n \times 5) \times (n \times 5) \end{matrix}$$

Dominic Adjucci Homework 6 BIOS 6643

Question 7 [Continued]

4.) How will you interpret the coefficient on the interaction with visit in the first model?

If there is a positive interaction coefficient this means that as the visit number increases then the biomarker has an increase in Hippocampus volume. If there is a negative coefficient this means as visits increase in number the biomarker has a decrease in Hippocampus volume.

I.) How would you interpret the coefficients on the interaction with visit in the second model?

	visit 2	visit 3	visit 4	visit 5
Biomarker	visit 2 * biomarker	visit 3 * biomarker	visit 4 * biomarker	visit 5 * biomarker

Because visit is Factorial each visit will have its own interaction term with biomarker. For visit 2 * biomarker this coefficient would account for the interaction between visit 2 and the biomarker on the outcome of Hippocampus volume. If positive this would mean the biomarker would have had a positive (increase) effect on Hippocampus volume. If negative this would mean the biomarker would have decreased Hippocampus volume by visit 2.

5.) How will this model be estimated?

The generalized least squares method will be used to estimate the model.



Question 8

Example 5: We received brain volume metrics and biomarker data from a neuropsychologist who is interested in how the biomarker influences the trajectory of brain volumes over 5 years in those with mild cognitive impairment. Patients received an MRI each year and different brain volumes were computed. For this project the investigator is interested in the hippocampus (one brain region implicated in Alzheimer's). The biomarker data was only measured at baseline. Hippocampal volume can be considered a continuous measure with an infinite number of values possible, meaning it isn't count or binary data.

A) Pseudo code

```
Imm.AR.Fit <- gls(hippoc ~ visit + biomarker + visit * biomarker + age + sex + SES, data = brainData,  
correlation = CorAR1(form = 1 | ID), method = "ML")
```

```
Imm.AR.Summary <- summary(Imm.AR.Fit)
```

```
Imm.AR.Fit2 <- gls(hippoc ~ visit + biomarker + visit * biomarker + age + sex + SES, data = brainData,  
correlation = CorAR1(form = 1 | ID), method = "ML", weights = varIdent(form = 1 | visit))
```

```
Imm.AR2.Summary <- summary(Imm.AR.Fit2)
```

B) Approach

The first model uses the generalized least squares method to estimate the beta coefficients. The relationship between visit and the biomarker in changing hippocampal volume will be assessed using the visit * biomarker interaction term.

The model will also ^{control for} age, sex, and SES. An autoregressive correlation structure will be used and maximum likelihood will be used (as opposed to PGM).

The second model uses the generalized least squares method to estimate the beta coefficients. The relationship between visit and the biomarker in changing hippocampal volume will be assessed using the visit * biomarker interaction term. The model will also control for age, sex, and SES. An autoregressive correlation structure and maximum likelihood will be used. The weights of the variances will be adjusted between visits.

c) Distributional Assumptions

The gls method assumes a multivariate normal distribution.

D.) Write out the distribution of the outcome with correct subscripts and definition of indices. Write out the expected value and variance-covariance for this distribution for each model.

$$E[Y_{ij}] = \beta_0 + \beta_1 \text{visit}_{ij} + \beta_2 \text{biomarker}_{ij} + \beta_3 \text{visit} * \text{biomarker}_{ij} + \beta_4 \text{age}_{ij} + \beta_5 \text{sex}_{ij} + \beta_6 \text{SES}_{ij}$$

Expected value is the same for each model.

Where i is the subject; $i = 1, \dots, n$ n = number of subjects

Where j is the observation; $j = 1, \dots, n_i$ n_i = number of observations for a subject.

First model variance-covariance.

$$R_i = \begin{bmatrix} \sigma^2 & \sigma^2 p & \sigma^2 p^2 & \sigma^2 p^3 & \sigma^2 p^4 \\ \sigma^2 p & \sigma^2 & \sigma^2 p & \sigma^2 p^2 & \sigma^2 p^3 \\ \sigma^2 p^2 & \sigma^2 p & \sigma^2 & \sigma^2 p & \sigma^2 p^2 \\ \sigma^2 p^3 & \sigma^2 p^2 & \sigma^2 p & \sigma^2 & \sigma^2 p \\ \sigma^2 p^4 & \sigma^2 p^3 & \sigma^2 p^2 & \sigma^2 p & \sigma^2 \end{bmatrix}$$

p represents correlation.

Second model Variance - Covariance

$$R_i = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 & \rho^2\sigma_1\sigma_3 & \rho^3\sigma_1\sigma_4 & \rho^4\sigma_1\sigma_5 \\ \rho\sigma_2\sigma_1 & \sigma_2^2 & \rho\sigma_2\sigma_3 & \rho^2\sigma_2\sigma_4 & \rho^3\sigma_2\sigma_5 \\ \rho^2\sigma_3\sigma_1 & \rho\sigma_3\sigma_2 & \sigma_3^2 & \rho\sigma_3\sigma_4 & \rho^2\sigma_3\sigma_5 \\ \rho^3\sigma_4\sigma_1 & \rho^2\sigma_4\sigma_2 & \rho\sigma_4\sigma_3 & \sigma_4^2 & \rho\sigma_4\sigma_5 \\ \rho^4\sigma_5\sigma_1 & \rho^3\sigma_5\sigma_2 & \rho^2\sigma_5\sigma_3 & \rho\sigma_5\sigma_4 & \sigma_5^2 \end{bmatrix}$$

The subscripts on variances and standard deviations represents time points

Domènec Adamec Homework 6 BIOS 6643 Q8 [Continued]

E) What differed between these two model choices?

The variance-covariance matrix was different. For the second model each timepoint was allowed to have its own variance.

F) What is the linear model for these data in subject level and complete data notation for the first model?

Subject-level

$$\begin{bmatrix} Y_{i1} \\ Y_{i2} \\ \vdots \\ Y_{in_i} \end{bmatrix} = \begin{bmatrix} 1 & \text{visit}_{i1} & \text{biomarker}_{i1} & \text{visit} * \text{biomarker}_{i1} & \text{age}_{i1} \\ 1 & \text{visit}_{i2} & \text{biomarker}_{i2} & \text{visit} * \text{biomarker}_{i2} & \text{age}_{i2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & \text{visit}_{in_i} & \text{biomarker}_{in_i} & \text{visit} * \text{biomarker}_{in_i} & \text{age}_{in_i} \end{bmatrix}$$

$$\begin{bmatrix} \text{age}_{i1} & \text{sex}_{i1} & \text{SES}_{i1} \\ \text{age}_{i2} & \text{sex}_{i2} & \text{SES}_{i2} \\ \vdots & \vdots & \vdots \\ \text{age}_{in_i} & \text{sex}_{in_i} & \text{SES}_{in_i} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \\ \beta_6 \end{bmatrix} + \underbrace{\begin{bmatrix} E_{i1} \\ E_{i2} \\ \vdots \\ E_{in_i} \end{bmatrix}}_{= E_i}$$

Where i is the subject; $i = 1, \dots, n$ n = number of subjects

Where j is the observation; $j = 1, \dots, n_i$ n_i is the number of observations for a subject.

$$E_i \sim N(0, R_i) ; R_i \text{ defined previously.}$$

Complete Data

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \\ \beta_6 \end{bmatrix} + \begin{bmatrix} E_1 \\ E_2 \\ \vdots \\ E_n \end{bmatrix}$$

Y_i defined in subject level model

X_i defined in subject level model

$$E_i \sim N\left(0, \begin{bmatrix} R_1 & & 0 \\ & \ddots & \\ 0 & & R_{ni} \end{bmatrix}\right)$$

Q) How will this model be estimated?

This model will be estimated using Maximum likelihood.