

# Statistical Methods for Correlated Data

## General Regression Models

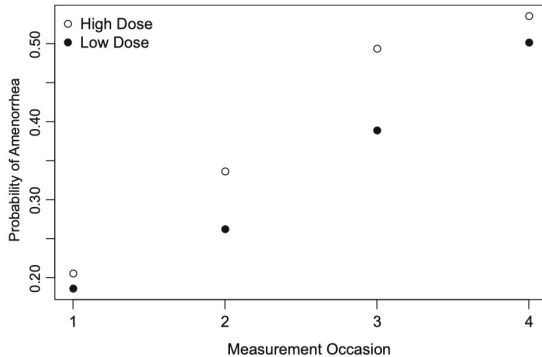
Michele Guindani

Department of Biostatistics  
UCLA

- ▶ GLMs can be extended to incorporate dependencies in obs on the same unit
- ▶ Two ways to do this:
  - ▶ **conditional models:** the dependences are introduced through modeling of unit-specific random effects:
    - ▶ Generalized linear mixed models (GLMMs)
    - ▶ Bayesian methods
  - ▶ **marginal models:** first- and second- moment assumptions only
    - ▶ Generalized estimating equations

## Ex: Contraception Data

- ▶ A randomized longitudinal contraception trial
- ▶ 1,151 women received a low and a high dose of DepoProvera, a drug used for contraceptive purposes, on the day of randomization and three additional injections at 90-day intervals. There was a final follow-up 3 months after the last injection (a year after the initial injection).
- ▶ **Outcome:** Amenorrhea yes/no, the absence of menstrual bleeding for a specified number of days, during each of the four 3-months intervals.
- ▶ Different # of measurements at different occasions
- ▶ Individual spaghetti-plot not informative for binary data; better to plot averages over time (approx probability of amenorrhea over time for the two treatment groups)



**Fig. 9.1** Probability of amenorrhea over time in low- and high-dose groups, in the contraception data

Increasing probabilities of amenorrhea in both groups, with the probabilities in the high dose group being greater than in the low dose group.

## Ex: Contraception Data

- ▶ With binary data, no clear measure of dependence

**Table 9.1** Empirical variances (on the *diagonal*) and correlations (on the *upper diagonal*), between measurements on the same woman at different observation occasions (1–4), in the low- (*left*) and high- (*right*) dose groups of the contraception data

	1	2	3	4		1	2	3	4
1	0.15	0.40	0.28	0.27	1	0.16	0.31	0.25	0.29
2		0.19	0.45	0.35	2		0.22	0.43	0.43
3			0.24	0.13	3			0.25	0.47
4				0.25	4				0.25

- Correlation between observations on the same woman, with a suggestion that the correlations decrease on measurements taken further apart
- Multivariate binary data are needed to take into account within-subject correlation

# Seizure data

- ▶ Thall and Vail (1990) describe data on epileptic seizures in 59 individuals.
- ▶ For each patient, the number of epileptic seizures was recorded during a baseline period of 8 weeks, after which patients were randomized to one of two groups: treatment with either the antiepileptic drug progabide or with placebo.
- ▶ The number of seizures was recorded in four consecutive 2-week periods
- ▶  $m_{placebo} = 28$  and  $m_{progabide} = 31$
- ▶ Let  $Y_{ij}$  represent the number of counts for patient  $i, i = 1, \dots, 59$  at occasion  $j$ , with  $j = 0$  the baseline period and  $j = 1, \dots, 4$  the subsequent set of four 2-week measurement periods.
- ▶ Let  $T_0$  be the baseline period and  $T_j = 2$  for  $j = 1, 2, 3, 4$ .
- ▶ See R

# Generalized Linear Mixed Models (GLMMs)

- ▶ A modeling framework that allows the introduction of random effects in the GLM framework (Breslow and Clayton, 1993)
- ▶ A GLMM is defined by the following two-stage model:
- ▶ **Stage One (conditional):** The distribution of the data is  $Y_{ij}|\theta_{ij}, \phi \sim p(\cdot)$  where  $p(\cdot)$  is a member

$$p(y_{ij}|\theta_{ij}, \phi) = \exp \{ [y_{ij}\theta_{ij} - b(\theta_{ij})] / a(\phi) + c(y_{ij}, \phi) \}$$

for  $i = 1, \dots, m$  units and  $j = 1, \dots, n_i$ , measurements per unit.  
The variance is

$$\text{var}(Y_{ij}|\theta_{ij}, \phi) = \phi v(\mu_{ij})$$

with  $\phi$  a dispersion parameter and  $v(\mu_{ij})$  indicating how the variance is functionally related to the mean response

# Generalized Linear Mixed Models (GLMMs)

- ▶ A modeling framework that allows the introduction of random effects in the GLM framework (Breslow and Clayton, 1993)
- ▶ A GLMM is defined by the following two-stage model:
- ▶ **Stage One (conditional):** The distribution of the data is  $Y_{ij}|\theta_{ij}, \phi \sim p(\cdot)$  where  $p(\cdot)$  is a member

$$p(y_{ij}|\theta_{ij}, \phi) = \exp \{ [y_{ij}\theta_{ij} - b(\theta_{ij})] / a(\phi) + c(y_{ij}, \phi) \}$$

for  $i = 1, \dots, m$  units and  $j = 1, \dots, n_i$ , measurements per unit.  
The variance is

$$\text{var}(Y_{ij}|\theta_{ij}, \phi) = \phi v(\mu_{ij})$$

with  $\phi$  a dispersion parameter and  $v(\mu_{ij})$  indicating how the variance is functionally related to the mean response



# Generalized Linear Mixed Models (GLMMs)

- ▶ Let  $\mu_{ij} = \text{E}[Y_{ij}|\theta_{ij}, \alpha]$  and, for a link function  $g(\cdot)$ , suppose

$$g(\mu_{ij}) = \mathbf{x}_{ij}\boldsymbol{\beta} + \mathbf{z}_{ij}\mathbf{b}_i$$

(systematic component)  $g(\cdot)$  is called the link function

# Generalized Linear Mixed Models (GLMMs)

- ▶ Let  $\mu_{ij} = \text{E}[Y_{ij}|\theta_{ij}, \alpha]$  and, for a link function  $g(\cdot)$ , suppose

$$g(\mu_{ij}) = \mathbf{x}_{ij}\boldsymbol{\beta} + \mathbf{z}_{ij}\mathbf{b}_i$$

(systematic component)  $g(\cdot)$  is called the link function

- ▶ **State two:** The random effects are assigned a normal distribution:

$$\mathbf{b}_i|\mathbf{D} \sim_{iid} \mathbf{N}_{q+1}(\mathbf{0}, \mathbf{D})$$

- ▶ Given  $\mathbf{b}_i$ ,  $Y_{ij} \perp\!\!\!\perp Y_{ik}$  (conditional independence assumption)

## Example: Poisson GLMMs

►  $Y_{ij} | b_i \sim \text{Poisson}(\mu_{ij}) \quad \mu_{ij} = E(y_{ij} | b_i)$

☞  $E(Y_{ij} | b_i) = \mu_{ij}, \text{Var}(Y_{ij} | b_i) = \mu_{ij}$

► log-link

$$g(\mu_i) = \log(\mu_{ij})$$

$$\log(\mu_{ij}) = \mathbf{x}_{ij}\boldsymbol{\beta} + \mathbf{z}_{ij}\mathbf{b}_i$$

## Example: Poisson GLMMs

►  $Y_{ij} | b_i \sim \text{Poisson}(\mu_{ij}) \quad \mu_{ij} = E(y_{ij} | b_i)$

☞  $E(Y_{ij} | b_i) = \mu_{ij}, \text{Var}(Y_{ij} | b_i) = \mu_{ij}$

► log-link

$$g(\mu_i) = \log(\mu_{ij})$$

$$\log(\mu_{ij}) = \mathbf{x}_{ij}\boldsymbol{\beta} + \mathbf{z}_{ij}\mathbf{b}_i$$

► How to interpret the coefficients?

► In GLMMs, the regression parameters  $\boldsymbol{\beta}$  have **subject-specific** interpretations:

☞ They represent the influence of covariates on a subject-specific response

☞ regression coefficients represent the effect of **within-subject** changes in covariates and changes in an individual's transformed mean response

## Example: Poisson GLMMs

►  $\log(\mu_i) = \beta_1 + \beta_2 t_{ij} + b_{1i} \quad b_{1i} \sim N(0, \sigma_b^2)$

👉  $\mu_i = e^{\beta_1 + \beta_2 t_{ij} + b_{1i}}$

►  $e^{\beta_1}$  : mean response for an individual with  $b_{1i} = 0$  (typical individual) ( $\neq$  mean response in the population)

►  $e^{\beta_2}$  : multiplicative change in mean response (rate of occurrence) within an individual for a unit increase in time (e.g.  $\beta_2 = 1.36$  (36% increase);  $\beta_2 = .75$  (25% decrease))

👉  $\beta_2$  represents the change in the log-expected rate for a single unit increase in the predictor (time) within an individual (or with respect to a typical individual, where “typical” typically means  $b_i = 0$ ).

# Binary GLMMs

►  $y_{ij}|b_i \sim \text{Bern}(\pi_{ij})$ , i.e.  $y_{ij} = \begin{cases} 0 & \text{with prob } 1 - \pi_{ij} \\ 1 & \text{with prob } \pi_{ij} \end{cases}$

►  $E(Y_{ij}|b_i) = \pi_{ij}$ ,  $\text{Var}(y_{ij}|b_i) = \pi_{ij}(1 - \pi_{ij})$

►  $g(\mu_{ij}) = \log\left(\frac{\pi_{ij}}{1 - \pi_{ij}}\right) = \text{logit}(\pi_{ij}) = \mathbf{x}_{ij}\boldsymbol{\beta} + \mathbf{z}_{ij}\mathbf{b}_i$

► **Example:**  $\log\left(\frac{\pi_{ij}}{1 - \pi_{ij}}\right) = \beta_1 + \beta_2 t_{ij} + b_{1i} + b_{2i} t_{ij}$

►  $\underbrace{\frac{\pi_{ij}}{1 - \pi_{ij}}}_{\text{odds}} = e^{\beta_1 + \beta_2 t_{ij} + b_{1i} + b_{2i} t_{ij}}$

► **Interpretation:**

►  $e^{\beta_1}$  : odds of  $Y_i = 1$  at baseline for a typical subj ( $b_{ji} = 0$ )

►  $e^{\beta_2}$  : multiplicative change in odds of  $y_{1,j} = 1$  for a “typical subject” for one unit interval in time

- ▶ However in GLMMs it is harder to interpret fixed effects coefficients for covariates that are time-invariant (e.g., treatment/placebo; gender; etc...)
- ▶ In those cases you may want to compare two “typical individuals”, i.e. two individuals with the same random effects (e.g. a typical individual under treatment vs a typical individual under placebo)

- ▶ However in GLMMs it is harder to interpret fixed effects coefficients for covariates that are time-invariant (e.g., treatment/placebo; gender; etc...)
- ▶ In those cases you may want to compare two “typical individuals”, i.e. two individuals with the same random effects (e.g. a typical individual under treatment vs a typical individual under placebo)
- ▶ In GLMMs, covariate effects are generally different at the marginal and individual level:
- ▶ GLMMs are typically parameterized in terms of covariate effects on conditional means
  - 👉 within-subject or subject-specific changes in covariates



- ▶ However in GLMMs it is harder to interpret fixed effects coefficients for covariates that are time-invariant (e.g., treatment/placebo; gender; etc...)
- ▶ In those cases you may want to compare two “typical individuals”, i.e. two individuals with the same random effects (e.g. a typical individual under treatment vs a typical individual under placebo)
- ▶ In GLMMs, covariate effects are generally different at the marginal and individual level:
- ▶ GLMMs are typically parameterized in terms of covariate effects on conditional means
  - 👉 within-subject or subject-specific changes in covariates
- ▶ Differently than in LMEs, marginal effects are not easy to calculate. For example, the marginal mean is:

$$E(\mathbf{Y}_i) = E_{b_i} (g^{-1} (x\boldsymbol{\beta} + \mathbf{z}_i\mathbf{b}_i)) = \dots$$

# Marginal moments in GLMMs

- The marginal variance is

$$\begin{aligned}\text{var}(Y_{ij}) &= \text{E}[\text{var}(Y_{ij}|\mathbf{b}_i)] + \text{var}(\text{E}[Y_{ij}|\mathbf{b}_i]) \\ &= \alpha \text{E}_{b_i} [v\{g^{-1}(\mathbf{x}_{ij}\boldsymbol{\beta} + \mathbf{z}_{ij}\mathbf{b}_i)\}] + \text{var}_{b_i}[g^{-1}(\mathbf{x}_{ij}\boldsymbol{\beta} + \mathbf{z}_{ij}\mathbf{b}_i)]\end{aligned}$$

- The covariances between outcomes on the same unit are

$$\begin{aligned}\text{cov}(Y_{ij}, Y_{ik}) &= \text{E}[\text{cov}(Y_{ij}, Y_{ik}|\mathbf{b}_i)] + \text{cov}[\text{E}(Y_{ij}|\mathbf{b}_i), \text{E}(Y_{ik}|\mathbf{b}_i)] \\ &= \text{cov}_{b_i}[g^{-1}(\mathbf{x}_{ij}\boldsymbol{\beta} + \mathbf{z}_{ij}\mathbf{b}_i), g^{-1}(\mathbf{x}_{ik}\boldsymbol{\beta} + \mathbf{z}_{ik}\mathbf{b}_i)] \\ &\neq 0\end{aligned}$$

for  $j \neq k$  due to shared random effects,

- 👉  $\boldsymbol{\beta}$ 's can only be interpreted with reference to an individual subject profile
- 👉 Interpret the effects on “typical” subjects ( $\mathbf{b}_i = \mathbf{0}$ ) or as effect sizes that only apply to subjects with the same  $\mathbf{b}_i$  values

## Example: Probit model

- ▶  $Y_{ij}|\mu_{ij} \sim \text{Bernoulli}(\mu_{ij})$ , with  $\Phi^{-1}(\mu_{ij}) = \eta_{ij}$  where  $\Phi^{-1}(\cdot)$  is the inverse standard normal cumulative distribution function or probit link.

## Example: Probit model

- ▶  $Y_{ij}|\mu_{ij} \sim \text{Bernoulli}(\mu_{ij})$ , with  $\Phi^{-1}(\mu_{ij}) = \eta_{ij}$  where  $\Phi^{-1}(\cdot)$  is the inverse standard normal cumulative distribution function or probit link.
- ▶ The probit model can be written as:

$$\eta_{ij} = \mathbf{x}_{ij}'\boldsymbol{\beta} + \xi_{ij}, \quad \xi_{ij} = \mathbf{z}_{ij}'\mathbf{b}_i + e_{ij}$$

where  $\xi_{ij}$  is the "total residual" or random part (conditioning on  $\mathbf{x}_{ij}$  and  $\mathbf{z}_{ij}$ ) of the model with variance

$$\sigma_{i,jj} \equiv \mathbf{z}_{ij}'\mathbf{D}\mathbf{z}_{ij} + \sigma_e^2$$

## Example: Probit model

- The marginal cumulative response probabilities become

$$\begin{aligned}\Pr(Y_{ij} = 1|\mathbf{x}_{ij}, \mathbf{z}_{ij}) &= \Pr(\eta_{ij} > 0|\mathbf{x}_{ij}, \mathbf{z}_{ij}) \\ &= \Pr(\mathbf{x}'_{ij}\boldsymbol{\beta} + \xi_{ij} > 0|\mathbf{x}_{ij}, \mathbf{z}_{ij}) \\ &= \Pr(-\xi_{ij} \leq \mathbf{x}'_{ij}\boldsymbol{\beta}|\mathbf{x}_{ij}, \mathbf{z}_{ij}) \\ &= \Pr\left(\frac{\xi_{ij}}{\sqrt{\sigma_{i,jj}}} \leq \frac{\mathbf{x}'_{ij}\boldsymbol{\beta}}{\sqrt{\sigma_{i,jj}}}\middle|\mathbf{x}_{ij}, \mathbf{z}_{ij}\right) \\ &= \Phi\left(\frac{\mathbf{x}'_{ij}\boldsymbol{\beta}}{\sqrt{\sigma_{i,jj}}}\right)\end{aligned}$$

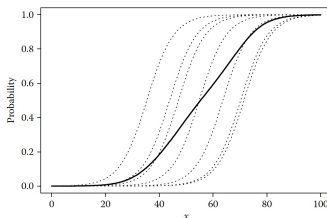
- Here the probit link is preserved for the marginal probabilities but with different regression coefficients; the marginal effects  $\beta/\sqrt{\sigma_{i,jj}}$  are attenuated, or closer to zero, compared to the conditional effects  $\beta$ .

## Remarks:

- For the log link, it can be shown that the link function is also preserved with

$$E(Y_{ij} \mid \mathbf{x}_{ij}, \mathbf{z}_{ij}) = \exp(\mathbf{x}'_{ij}\boldsymbol{\beta} + \mathbf{z}'_{ij}\mathbf{D}\mathbf{z}_{ij}/2)$$

- Apart from the probit and log link, the link function for the marginal model is generally different from that for the conditional model and usually does not have a simple form.
- In logistic models, the marginal regression coefficients are attenuated compared to the conditional ones



Solid: marginal; Dashed: conditional

## Remark: Bad notation

Generalized linear mixed models are often written as non-linear models with an error term. For instance, mixed logit or probit models for binary responses are written as

$$Y_{ij} = \pi_{ij} + \epsilon_{ij} z_{ij}^{(1)}, \quad \pi_{ij} = h(\eta_{ij}), \quad z_{ij}^{(1)} = \sqrt{\pi_{ij}(1 - \pi_{ij})}$$

Constraining the variance of  $\epsilon_{ij}$  to one, we obtain the required Bernoulli variance  $\pi_{ij}(1 - \pi_{ij})$ . However, this formulation is awkward because it requires a very peculiar distribution of  $\epsilon_{ij}$  to produce valid 0 or 1 responses. More importantly, the formulation gives the false impression that the variance of  $\epsilon_{ij}$  could be estimated and should therefore be avoided.