

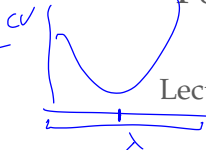
# Linear Regression Models

P8111

Lecture 23

- HW 4 notes

- no shaped CV
- outlier removal
- 



Jeff Goldsmith

April 19, 2016



THE DEPARTMENT OF  
**BIostatISTICS**



Columbia University  
**MAILMAN SCHOOL  
OF PUBLIC HEALTH**

# Today's Lecture

- Multilevel models ✓
  - ▶ Hierarchical / nested models
  - ▶ Crossed designs ✓
- Bayesian methods

# Longitudinal data

- We observe data  $y_{ij}, x_{ij}$  for subjects  $i = 1, \dots, I$  at visits  $j = 1, \dots, J_i$
- Overall, we pose the model

$$\underline{y = X\beta + \epsilon}$$

where  $\text{Var}(\epsilon) = \sigma^2 V$  and

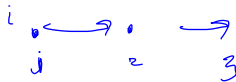
$v(p)$

$$V = \begin{bmatrix} \boxed{V_1} & 0 & \dots & 0 \\ 0 & \boxed{V_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \\ 0 & 0 & & \boxed{V_I} \end{bmatrix}$$

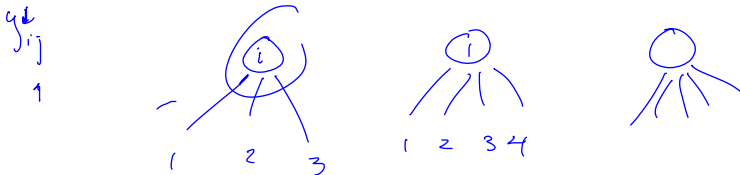
# Longitudinal data

- Extended cross-sectional models to allow repeated subject observations
- Repeated observations had a time element
- One basic approach was random effects

# Multilevel models



- Multilevel models are a 'more general' class of models
- Repeated observations don't necessarily have to be taken in time
- Examples of two-level models include students in a class, members in a family, patients in a hospital, etc



# Two-level model

The repeated observations structure we developed for longitudinal data helps for two-level models. Specifically for repeated observations  $j$  within clusters  $i$ , we could write

$$\underbrace{y_{ij}} = \underbrace{\beta_0} + \underbrace{\beta_1 x_{ij}} + \underbrace{b_i}_{\sim} + \underbrace{\epsilon_{ij}}$$

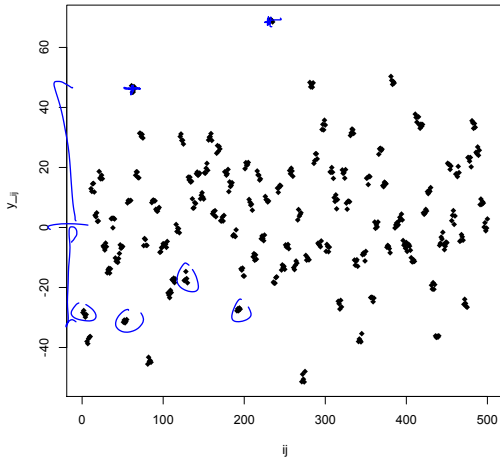
with

- $b_i \sim N[0, \tau^2]$
- $\epsilon \sim N[0, \nu^2]$

Intuition, estimation, induced correlation, interpretation – all of these were established for LDA and transfer here

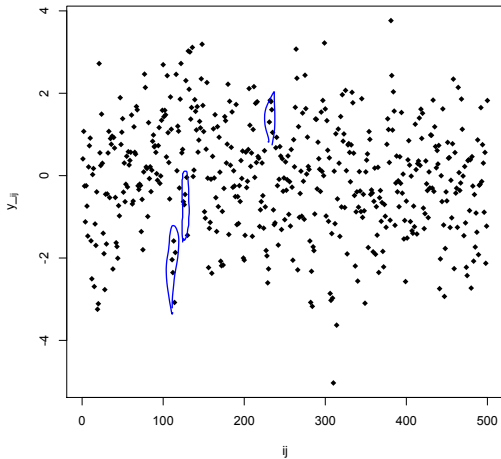
# Example I

$$\tau^2 \gg \sigma^2 \quad / \quad ICC \approx 1$$



## Example II

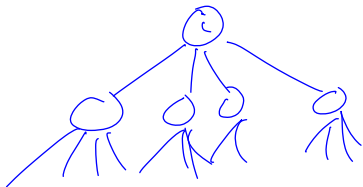
$$\gamma^2 \ll \nu^2 \quad / \quad I(\gamma) \approx 0$$





# Three level model

- Sometimes, the data have a more complex nested structure
- Each cluster is part of a larger cluster
- Examples include students in classes in universities, members in families in towns, patients in hospitals in regions



# Three level model

For a model with three levels (repeated observations  $k$  within clusters  $j$ , within super-clusters  $i$ ), we can write

$$y_{ijk} = \beta_0 + \beta_1 x_{ijk} + \underbrace{b_i}_{\text{super-cluster}} + \underbrace{b_{ij}}_{\text{cluster}} + \epsilon_{ijk}$$

with

- $b_i \sim N[0, \tau_{(1)}^2]$
- $b_{ij} \sim N[0, \tau_{(2)}^2]$
- $\epsilon \sim N[0, \nu^2]$

$$\left[ \begin{array}{ccccc} 1 & 0 & 0^1 & 0^2 & 0^3 \\ 0 & 1 & 0 & 1 & \\ & & & & 1 \end{array} \right]$$

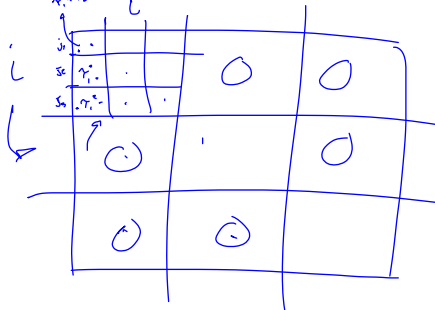
# ICCs

This model gives two levels of correlation (observations within clusters, clusters within super clusters), and therefore a couple of ICCs:

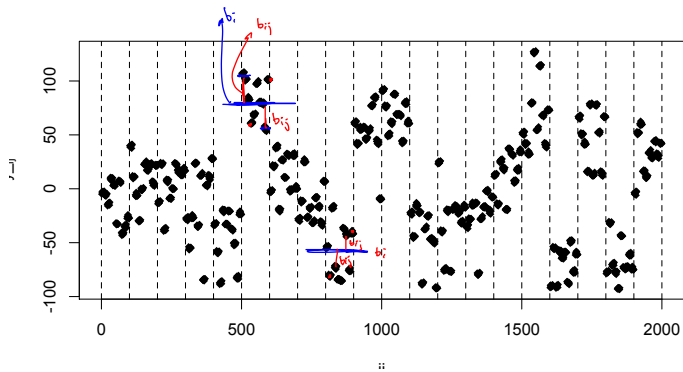
$$\blacksquare \text{cov}(\underline{y}_{ijk}, \underline{y}_{ijk'}) = \text{Cov}(\underbrace{\rho_0 + \rho_1 x_{ijk} + b_i + b_j + \epsilon_{ijk}}_{\text{Predictor}}, \underbrace{\rho_0 + \rho_1 x_{ijk'} + b_i + b_j + \epsilon_{ijk'}}_{\text{Predictor}})$$

$$= \underbrace{\tau_1^2} + \underbrace{\tau_2^2}$$

- $COV(y_{ijk}, y_{ij'k}) = \tau_i^2$



# Example



# Example

## Nested model

$y_{ijk} = \beta_0 + \beta_i + \beta_{ij} + \epsilon_{ijk}$

```
> nested.mod = lmer(yijk ~ (1+x1i | L1) + (1+x2j | L2))  
> summary(nested.mod)  
Linear mixed model fit by REML ['lmerMod']  
Formula: yij ~ (1 | L1) + (1 | L2)
```

REML criterion at convergence: 7464.337

Random effects:

Groups	Name	Variance	Std.Dev.
L2	(Intercept)	527.003	22.957
L1	(Intercept)	2137.453	46.233
Residual		1.004	1.002

Number of obs: 2000, groups: L2, 200; L1, 20

Fixed effects:

	Estimate	Std. Error	t value
(Intercept)	0.7712	10.4646	0.074

$$y_{ijk} = \beta_0 + \beta_i + \beta_{ij} + \epsilon_{ijk}$$

vs

$$(1 + x_1 | L_1) + (1 + x_2 | L_2)$$

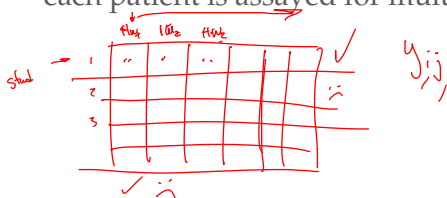
vs

$$(1 | L_1) + (1 + x_1 + x_2 | L_2)$$

# Crossed designs




- Alternatively to nested (hierarchical) models, sometimes there is a crossed design
- Each subject is observed under multiple “treatments”, so there are both subject and treatment effects
- For example, each student is graded in multiple classes; each patient is assayed for multiple genes



# Crossed designs

For a crossed model (with subjects  $i$  and treatments  $j$ ), we can write

$$\underline{y_{ij}} = \beta_0 + \beta_1 x_{ij} + \underline{b_i} + \underline{b_j} + \epsilon_{ij}$$


with

- $b_i \sim N[0, \tau_{(1)}^2]$
- $b_j \sim N[0, \tau_{(2)}^2]$
- $\epsilon \sim N[0, \nu^2]$

Here there is covariance within subjects across treatments, and within treatments across subjects.

# ICCs

Here there is covariance within subjects across treatments, and within treatments across subjects.

- $cov(y_{ij}, y_{ij'}) =$

- $cov(y_{i'j}, y_{ij}) =$



# LDA and MLM

- Estimation works basically the same for these models as for random intercept models
- Intuition is the same as well – you want to borrow strength for one subject from the population of other subjects
- Interpretation of fixed effects is *marginal*; interpretation of random effects is *conditional*  $\hookrightarrow E(y) = E(E(y|b))$
- Using randomness both decreases the number of parameters and induces correlation structures  $= E(y|b=0)$

# Bayesian methods

Longitudinal data analysis and multilevel models are a good place to start “thinking Bayesian”

- Even though they're frequentist, they include randomness at subject levels
- The idea of “shrinking toward a population mean” or “borrowing strength” is a pretty Bayesian concept
- Even writing down random effect distributions is reminiscent of defining prior distributions

# Basic Bayes

LDA and MLM are fairly advanced topics, so we'll start with a simpler example

- Suppose I gather data  $y_i$  and want to learn about  $E(y)$
- Suppose even more I think I already know *something* about  $E(y)$
- I might write down something about what I want to learn and what I think I know

# Basic Bayes

- What do I think I know?
- $y_i | \mu \sim N[\mu, \sigma_y^2]$
  - $\mu \sim N[\mu_0, \sigma_0^2]$
- Handwritten annotations:  $E(y)$  with an arrow pointing to the  $\mu$  in the first equation, and a blue box around the first equation. A blue arrow points from the second equation to the first.

What do I want to learn?

- $\mu | y_i \sim ???$
- Handwritten annotation: a blue underline under the expression  $\mu | y_i$ .

# Basic Bayes

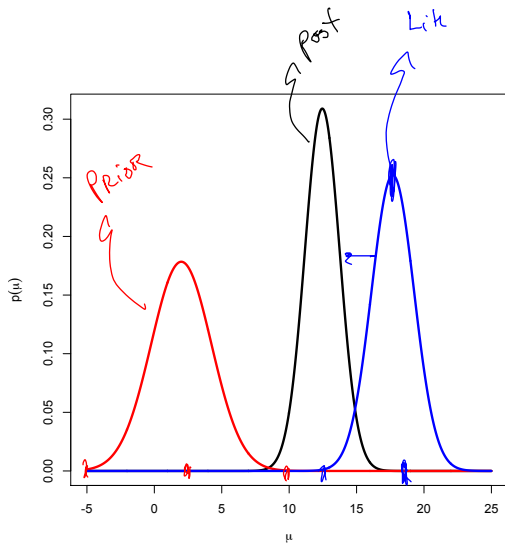
Luckily, this is all related through Bayes' formula:

$$\underbrace{p(\mu|y_i)}_{\text{posterior}} \propto \underbrace{p(y_i|\mu)}_{\text{likelihood}} \underbrace{p(\mu)}_{\text{prior}}$$

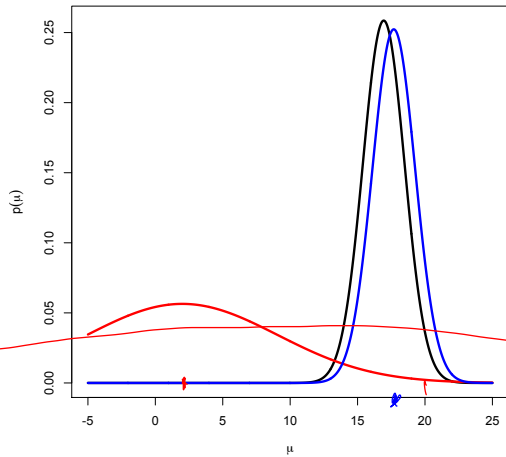
- For the Normal likelihood with a Normal prior for the mean, the posterior is also Normal:

$$\mu|y_i \sim N \left[ \underbrace{\frac{\sigma_\mu^2}{\frac{\sigma_y^2}{n} + \sigma_\mu^2} \bar{y}}_{\text{posterior mean}}, \underbrace{\frac{\frac{\sigma_y^2}{n} \sigma_\mu^2}{\frac{\sigma_y^2}{n} + \sigma_\mu^2}}_{\text{posterior variance}} \right]$$

# Effect of Prior



# Effect of Prior



# Bayesian regression

$$y^{\text{train}}; \mu$$



How can we pose the regression model

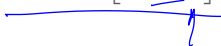
$$\underline{y} = \underline{X}\underline{\beta} + \underline{\epsilon}$$

$$y | \beta \sim \mathcal{N}(X\beta, \sigma^2 I)$$

with  $\epsilon \sim \mathcal{N}[0, I_n]$  in a Bayesian framework?

$$(X^T X)^{-1} X^T y$$

- By making distributional assumptions about the  $\beta$
- Normal priors seemed to work well in the past ...
- Try  $\beta \sim \mathcal{N}[0, \sigma_{\beta}^2 I_p]$  where  $p$  includes the intercept





# Bayesian regression

We want to obtain the posterior

$$p(\beta|\mathbf{y}, \mathbf{X}) \propto \underbrace{p(\mathbf{y}|\beta, \mathbf{X})} \underbrace{p(\beta)}$$

$$\exp \left\{ -\frac{1}{2} \dots \right\}$$

$$\exp \left\{ -\frac{1}{2} \dots \right\}$$

# Bayesian regression

Can show that

$\beta | y, X$

$$[\beta | y, X] \sim N[\mu_p, \Sigma_p]$$

where

$$\Sigma_p = \left( \frac{1}{\sigma_\epsilon^2} \mathbf{X}^T \mathbf{X} + \frac{1}{\sigma_\beta^2} \mathbf{I} \right)^{-1}$$

and

$$\mu_p = \Sigma_p \left( \frac{1}{\sigma_\epsilon^2} \mathbf{X}^t \mathbf{y} \right)$$

$$= \left( \mathbf{X}^T \mathbf{X} + \frac{\sigma_\epsilon^2}{\sigma_\beta^2} \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

So, about the variances

$$\underline{p | y}$$

$$y | p, \sigma_\epsilon^2$$

$$p | \sigma_\beta^2$$

- Throughout all of this we have implicitly conditioned on the variances  $\sigma_\epsilon^2$  and  $\sigma_\beta^2$
- Doesn't affect any of our calculations – the terms involving  $\mu$  don't overlap with terms involving  $\sigma_\epsilon^2$  or  $\sigma_\beta^2$
- $\sigma_\beta^2$  is often treated as fixed;  $\sigma_\epsilon^2$

$$\underline{p, \sigma_\epsilon^2 | y}$$

# The full posterior

- Need  $[\beta, \sigma_\epsilon^2 | \mathbf{y}, \mathbf{X}]$
- “Intractable” problem
- Just as good: sample from the posterior

# Sampling from the posterior

MCMC

- aka where Bayes gets really weird
- You can draw a sample from the posterior even if you can't write down exactly what it is
- That sample is your basis for inference
  - ▶ Posterior sample average is your estimate
  - ▶ Quantiles on the posterior sample define your credible interval
- Sample describes the *joint distribution* of all model parameters

## Some notes on this business

$$\hat{V}_{OLS}(\hat{\beta}_{OLS}) = \frac{\hat{\sigma}^2}{\uparrow} (X^T X)^{-1}$$

- Joint distributions are often worth the trouble
- Bayesian methods were really controversial for a long time, but are at least less controversial now
- The introduction of “prior knowledge” happens even in frequentist methods, although it is often not explicitly acknowledged

# Today's big ideas

- Nested and crossed random effects models
  - Bayesian stuff
-