

Khairil Anwar Notodiputro

Earned BS and Master degrees in statistics from Bogor Agricultural University and Ph.D in statistics from Macquarie University, Australia.

Now he is a **Profesor of Statistics at Bogor Agricultural University** and a **Visiting Professor at Prince of Songkla University, Thailand.**

He is a member of ISI (International Statistical Institute) and a former Chairman of Indonesian Statistical Association. He is actively doing competitive research as well as supervising undergraduate and postgraduate students. His research interests include **mixed models**, **small area estimation** and **time series analysis**.

During the last five years he has taught the following subjects: **Time Series Analysis and Forecasting, Advanced Data Analysis, Categorical Data Analysis** and **Statistical Analysis.**



Mixed Models in Practice

Khairil Anwar Notodiputro Professor of Statistics Bogor Agricultural University

khairilnotodiputro@gmail.com





INVITATION LETTER

March 21, 2016

Professor Khairil Anwar Notodiputro, Ph.D. Department of Statistics Faculty of Mathematics and Sciences Bogor Agricultural University Indonesia

Dear Professor Khairil Anwar Notodiputro,

The Department of Statistics invites you to visit Kasetsart University during May 2016. We expect that your visit will be on May 13th 2016.

We would like to invite you to give a special lecture on the topic of "Mixed Models in Practice" for graduate students. Furthermore, we would like you to collaborate on research projects that will focus on applied sciences in the future.

We would like to thank you in advance for visiting department of statistics at Kasetsart University, Thailand.

Sincerely yours,

honte

Assist.Prof. Boonorm Chomtee Department Head

DEPARTMENT OF STATISTICS, FACULTY OF SCIENCE KASETSART UNIVERSITY P. O. BOX 1086, BANGKOK 10903, THAILAND TEL: 66-0-2502-5444 e-mail: ficiboc@ku.ac.fb

Mixed Models in Practice

Khairil Anwar Notodiputro Department of Statistics, Bogor Agricultural University e-mail: <u>khairilnotodiputro@gmail.com</u>



Abstract. Over the past decade there has been an explosion of developments in mixed effects models and their applications. The models have been applied in many areas such as agriculture, medicine, biology, biostatistics, education, social, economic and management sciences. The mixed models have also been intensively used in small area estimation, an area of research where proper statistical methods are needed to obtain estimates with high precision in a situation where the available samples are very limited. In this presentation, the mixed models have been discussed. We started by discussing the limitation of linear models and then introduced the linear mixed models. The importance of the linear mixed models has become evident if the predictor variables consisted of both fixed effects and random effects. These two types of effects have often been observed in practice. However, the linear mixed models are based on normal distributions of the response variables as well as the random effects whereas in practice these variables may not be normally distributed. Hence, the generalized linear mixed models (GLMM) need to be employed for non-normal distributions of the response variables, especially the exponential family of distributions. Moreover, since the models are very general, then basically the models can be used to handle various distributions of the response variables as well as the random effect distributions. Finally, to demonstrate the practice of GLMM, we have shown our works on mixed models. Several of the works have been related to the Na Thap river project, an IPB - PSU Collaborative Research project, RPM *ID16287*.

Keywords: Bias Reduction, Generalized Linear Mixed Models Hierarchical Models, Linear Models, Mixed Models, Multilocation Eksperiment, Na Thap River



- Introduction
- Types of Model Effects
- Why Mixed Models?
- Estimation and Inference
- Mixed Model in Practice
- Concluding Remarks



Introduction

Introduction



Consider a **linear model** $y = Xb + \varepsilon$, where ε is assumed to have a Gaussian disribution.

What if *ɛ* is non-Gaussian?

Pair	Treatment _0	Treatmen t_1
1	98	94
2	95	36
3	93	85
4	94	88
5	99	91
6	61	82
7	84	43
8	92	71

Imagine a study to compare two treatments, "control" treatment vs "test" treatment. The response variable is the incidence of favourable outcomes. A paired comparison is used by observing 100 individuals on each treatment on each pair. Assuming independent observations, the number of favorable outcomes for each pair-treatment combination has a binomial distribution with n=100and probability π_{ii} for the *i*th treatment and *j*th pair. How should we model the data?

Introduction (cont..)



We can start with the normal approximation and the model can be described as follows:

Pair	Treatment_0	Treatment_1
1	98	94
2	95	36
3	93	85
4	94	88
5	99	91
6	61	82
7	84	43
8	92	71

- The sample proportions are 0.738 and 0.895
- T-test: p-value = 0.1132

- Response variable: Let $p_{ij} = y_{ij}/100$ denotes the sample proportion, where y_{ij} is the number of favorable outcomes out of 100 individuals.
- The model is $p_{ij} = \mu + \tau_i + \rho_j + \varepsilon_{ij}$ where μ denotes an overall mean, τ_i denotes i^{th} treatment effect, ρ_j denotes j^{th} pair effect, and ε_{ij} is the sampling error assumed i.i.d $N(0,\sigma^2)$

Introduction (cont..)



One obvious problem with the previous analysis is as follows

- Normal approximation not-withstanding, the response variable remains binomial, meaning the variance must depend on π.
- Assuming π changes from treatment to treatment (or among pairs) the assumption of equal variance does not hold.
- Transformation is required → the arc sine square root transformation.

At this point in this example, we have reached the limits of the linear model.

Introduction (cont..)



The limitation of the Linear Model (LM) can be overcome by the Generalized Linear Mixed Model (GLMM) which involves 3 components:

- The **distribution** of the observations: $y_{ij} / \rho_j \sim Binomial (100, \pi_{ij}).$
- The linear predictor: $\eta_{ij} = \eta + \tau_i + \rho_j$, where η is the intercept and asume the pair effects ρ_j are i.i.d. $N(0, \sigma_{\rho}^2)$.
- The **link function**: with non-normal data, the canonical parameter of the log likelihood is typically a better candidate for fitting a linear model than the mean itself. For the binomial, the canonical link function is $\eta_{ij} = log[\pi_{ij}/(1-\pi_{ij})]$.



Types of Model Effects

Types of Model Effects



- The types of model effects can be distinguished based on how the levels of an effect are chosen and the scope of the intended inference.
- The model of the previous example:

$$p_{ij} = \mu + \tau_i + \rho_j + \varepsilon_{ij}$$

- Treatment effects r have been fixed with two conditions (either "control" or "test")
- > Pair effects ρ depend on the individuals chosen as a sample.
- > So τ are fixed effects and ρ are random effects

Types of Model Effects (cont..)



In a **fixed-effects model** for an experiment, all the factors in the experiment have a predetermined set of levels and the only inferences are for the levels of the factors actually used in the experiment.

In a **random effects model** for an experiment, the levels of factors used in the experiment are randomly selected from a population of possible levels. The inferences from the data in the experiment are for all levels of the factors in the population from which the levels were selected and not only the levels used in the experiment.

Types of Model Effects (cont..)



In a **mixed-effects model** (or simply **mixed model**) for an experiment, the levels of some of the factors used in the experiment are randomly selected from a population of possible levels, whereas the levels of the other factors in the experiment are predetermined.

The inferences concerning factors with fixed levels are only for the levels of the factors used in the experiment, whereas inferences concerning factors with randomly selected levels are for all levels of the factors in the population from which the levels were selected.



Why Mixed Models?

Why Mixed Models?



Consider a hypothetical data (Stroup, 2013):

Binomial regression data Fav Obs Ν

Multi-batch data, two responses (Y and Fav)

X	Ba	atch 1		B	Batch2		Batch 3			Batch 4		
	Y	Fav	Ν	Y	Fav	Ν	Y	Fav	Ν	Y	Fav	Ν
0	95.6	15	21	96.6	18	21	96.6	16	19	96.6	18	21
3	96.9	13	17	96.7	14	16	96.5	13	19	96.8	18	21
6	98.5	19	23	96.3	14	17	97.7	17	22	96.4	14	20
9	99.0	14	17	96.3	17	20	98.3	23	27	96.6	14	19
12	100.2	18	23	96.5	15	20	99.1	16	21	96.8	14	19
18	101.9	19	27	96.4	14	22	100.5	11	16	96.9	11	20
24	104.1	15	20	95.7	12	22	101.2	13	18	97.1	11	16
36	107.8	14	21	95.2	11	25	103.3	10	17	97.4	10	21
48	111.5	13	18	94.6	5	26	105.9	6	16	97.4	5	19

One observation per level of X

Levels of *X* are observed for multiple batches



	X	Ba	tch 1		E	atch2		Ва	atch 3		В	atch 4	
		Y	Fav	Ν	Y	Fav	Ν	Y	Fav	Ν	Y	Fav	Ν
đ	0	95.6	15	21	96.6	18	21	96.6	16	19	96.6	18	21
5	3	96.9	13	17	96.7	14	16	96.5	13	19	96.8	18	21
at	6	98.5	19	23	96.3	14	17	97.7	17	22	96.4	14	20
9	9	99.0	14	17	96.3	17	20	98.3	23	27	96.6	14	19
ĮĘ,	12	100.2	18	23	96.5	15	20	99.1	16	21	96.8	14	19
Ju .	18	101.9	19	27	96.4	14	22	100.5	11	16	96.9	11	20
<	24	104.1	15	20	95.7	12	22	101.2	13	18	97.1	11	16
	36	107.8	14	21	95.2	11	25	103.3	10	17	97.4	10	21
	48	111.5	13	18	94.6	5	26	105.9	6	16	97.4	5	19

This table shows 9 levels of *X* varying from 0 to 48. For each of four batches, a continuous variable (*Y*) and a binomial variable (*N* = number of Bernoulli trials and *Fav* = number of "successes") are observed at each level of *X*.

How should these data be modeled ?



Look at the response Y. The middle line shows the average linear regression over all batches; the other lines show the regressions for each individual batch. Inspecting the graph suggests that assuming linear regression is reasonable but a common regression for all batches may not be justified.





Separate linear regressions by batch yields the linear predictor $\eta_{ij} = \beta_{0i} + \beta_{1i} X_{ij}$, where β_{0i} and β_{1i} are, respectively, the intercept and slope for the i^{th} batch. $\beta_0 + b_{0i}$

 $(\beta_1 + b_{1i})$

Or
$$\eta_{ij} = \beta_0 + b_{0i} + (\beta_1 + b_{1i})X_{ij}$$
,
 β_0 , β_1 overall intercept and slope
 b_{0i} , b_{1i} batch-specific deviations from β_0 and β_1 .



 Let us assume that the batches represent a sample of a larger population of batches so that b_{0i} and b_{1i} are random variables and, thus, have probability distributions.

Random effect

Hence, the above model, i.e.

 $\eta_{ij} = \beta_0 + b_{0i} + (\beta_1 + b_{1i})X_{ij}$ is a **mixed model**.

Fixed effect

 To summarize, the multi-batch regression example gives rise to the linear predictor:

$$\eta_{ij} = \beta_0 + b_{0i} + (\beta_1 + b_{1i}) X_{ij}$$

 (b_{0i}, b_{1i}) pairs are assumed independent and within each batch, the pairs are bivariate normal, that is,

$$\begin{bmatrix} b_{0i} \\ b_{1i} \end{bmatrix} \sim N \begin{pmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_0^2 & \sigma_{01} \\ \sigma_{01} & \sigma_1^2 \end{bmatrix} \end{pmatrix}$$

where

 σ_0^2 and σ_1^2 are the variances of b_{0i} and b_{1i} , respectively σ_{01} is the covariance between b_{0i} and b_{1i} .

- UT PERANIAN BOGOR
- If the observed data is continuous (Y) and have a Gaussian distribution, then the linear predictor is an estimate of the data's expected value, conditional on the random effects b_{0i} and b_{1i} as follows:
 - 1. Observations: $(y_{ij}/b_{0i}, b_{1i}) \sim N(\mu_{ij}/b_{0i}, b_{1i}, \sigma^2)$
 - 2. Model focus: $E(y_{ij}/b_{0i}, b_{1i}) = \mu_{ij}/b_{0i}, b_{1i}$
 - 3. Linear predictor: $\eta_{ij} = \beta_0 + b_{0i} + (\beta_1 + b_{1i}) X_{ij}$
 - 4. Assumptions about b_{0i} and b_{1i} (if random), for example, as shown previously
 - 5. Link function: identity, that is, μ_{ij}/b_{0i} , b_{1i} modeled by $\beta_0 + b_{0i} + (\beta_1 + b_{1i}) X_{ij}$



- On the other hands, if you model the variable Fav in the table, then you adjust the distribution of the observations and the link function, accordingly.
 - 1. Observations: $(Fav_{ij}/b_{0i}, b_{1i}) \sim Binomial[N_{ij}, (\pi_{ij}/b_{0i}, b_{1i})]$
 - 2. Model focus: $E(Fav_{ij}/b_{0i}, b_{1i}) = \pi_{ij}/b_{0i}, b_{1i}$
 - 3. Linear predictor: $\eta_{ij} = \beta_0 + b_{0i} + (\beta_1 + b_{1i}) X_{ij}$
 - 4. Assumptions about b_{0i} and b_{1i} (if random)
 - 5. Link function: logit for example, $log[\pi_{ij}/(1 \pi_{ij})]$ modeled by $\beta_0 + b_{0i} + (\beta_1 + b_{1i}) X_{ij}$, where π_{ij} is used here as a shorthand for $\pi_{ij}/b_{0i}, b_{1i}$

Back to the linear predictor:

$$\eta_{ij} = \beta_0 + b_{0i} + (\beta_1 + b_{1i}) X_{ij}$$

- This predictor consists of **fixed** effects and **random** effects.
- The fixed effects component of the linear predictor is $(\beta_0 + \beta_1 X_{ij})$ whereas the random effects component is $(b_{0i} + b_{1i} X_{ij})$.
- In matrix form the fixed effects are denoted by *X*β whereas the random effects are denoted by *Zb*.
- Hence, the linear predictor: $\eta = X\beta + Zb$
- Vector **b** ~ N(**0**, **G**)

PERAPN

BOGOR

INSTITUTED



Xß



- Recall that a fully specified GLMM is given by:
 - A linear predictor X\(\mathcal{\mathcal{X}\mathcal{F}} + Z\mathcal{b}\) (or simply X\(\mathcal{\mathcal{A}}\) if there are no random effects)
 - The distribution of *y* | *b* the observations, conditional on the random effects if there are random effects
 - The link function, $\eta = g(\mu | b)$, where $\mu | b = E(y | b)$
 - The random effects distribution, b ~ N(0,G)
- For models with Gaussian data dan random effects:
 - Linear predictor: $\eta = X\beta + Zb$
 - Distributions: y b ~ N(µ b, R); b ~ N(0,G)
 - Link: η = μ | b





 The log-likelihood equations for y b and b, respectively, are

$$\mathcal{E}(\mathbf{y} \mid \mathbf{b}) = -\left(\frac{n}{2}\right)\log(2\pi) - \left(\frac{1}{2}\right)\log\left(|\mathbf{R}|\right) - \left(\frac{1}{2}\right)(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b})'\mathbf{R}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b})$$

and

$$\ell(\mathbf{b}) = -\left(\frac{b}{2}\right)\log(2\pi) - \left(\frac{1}{2}\right)\log\left(|\mathbf{G}|\right) - \left(\frac{1}{2}\right)\mathbf{b'}\mathbf{G}^{-1}\mathbf{b'}$$

Focusing on the "quasi-likelihood" part:

$$\ell(\mathbf{y},\mathbf{b}) = -\left(\frac{1}{2}\right)\left(\mathbf{y} - \mathbf{X}\mathbf{\beta} - \mathbf{Z}\mathbf{b}\right)'\mathbf{R}^{-1}\left(\mathbf{y} - \mathbf{X}\mathbf{\beta} - \mathbf{Z}\mathbf{b}\right) - \left(\frac{1}{2}\right)\mathbf{b}'\mathbf{G}^{-1}\mathbf{b}$$



MLE can be obtained by setting $\partial l(y,b)/\partial \beta'$ and $\partial l(y,b)/\partial b'$ equal to zero and solving the resulting set of equations for β and b:

$$\frac{\partial \left[\ell(\mathbf{y},\mathbf{b})\right]}{\partial \mathbf{\beta}'} = \mathbf{X}'\mathbf{R}^{-1}\mathbf{y} - \mathbf{X}'\mathbf{R}^{-1}\mathbf{X}\mathbf{\beta} - \mathbf{Z}'\mathbf{R}^{-1}\mathbf{X}\mathbf{b}$$
$$\frac{\partial \left[\ell(\mathbf{y},\mathbf{b})\right]}{\partial \mathbf{\beta}'} = \mathbf{X}'\mathbf{R}^{-1}\mathbf{y} - \mathbf{X}'\mathbf{R}^{-1}\mathbf{x}\mathbf{\beta} - \mathbf{Z}'\mathbf{R}^{-1}\mathbf{x}\mathbf{b}$$

$$\frac{\partial \left[\ell(\mathbf{y},\mathbf{b})\right]}{\partial \mathbf{b}'} = \mathbf{Z}'\mathbf{R}^{-1}\mathbf{y} - \mathbf{X}'\mathbf{R}^{-1}\mathbf{Z}\mathbf{\beta} - \mathbf{Z}'\mathbf{R}^{-1}\mathbf{Z}\mathbf{b} - \mathbf{G}^{-1}\mathbf{b}$$

The resulting mixed model equations:

$$\begin{bmatrix} X'R^{-1}X & X'R^{-1}Z \\ Z'R^{-1}X & Z'R^{-1}Z+G^{-1} \end{bmatrix} \begin{bmatrix} \beta \\ b \end{bmatrix} = \begin{bmatrix} X'R^{-1}y \\ Z'R^{-1}y \end{bmatrix}$$

This is called the Linear Mixed Model, i.e. the GLMM with Gaussian response and random effects.

- If the distribution of response variable is not normal then we use the GLMM with the essential features:
 - Linear predictor: $\eta = X\beta + Zb$
 - Distribution: *b* ~ *N(0,G)*
 - Distribution or quasi-likelihood: $E(\mathbf{y} | \mathbf{b}) = \boldsymbol{\mu} | \mathbf{b};$ $Var(\mathbf{y} | \mathbf{b}) = V_{\mu}^{\frac{1}{2}} \mathbf{A} V_{\mu}^{\frac{1}{2}}$, where $V_{\mu}^{\frac{1}{2}} = \text{diag} [(\partial^2 \mathbf{b}(\mathbf{\theta})/\partial \mathbf{\theta}^2)^{\frac{1}{2}}]$ and $\mathbf{A} = diag [1/a(\varphi)]$, and $\mathbf{y} | \mathbf{b}$ either has a distribution that belongs to the exponential family or a quasi-likelihood

• Link: $\eta = g(\mu \mid b)$, or alternatively, inverse link: $X\beta + Zb = h(\eta)$. Typically, $h(\cdot) = g^{-1}(\cdot)$

 The quasi-likelihood of the observation conditional on the random effects is

$$ql(y|b) = y'A\theta - 1'Ab(\theta)$$

The log-likelihood of the random effects is $\ell(b) = -(b/2) \log(2\pi) - (1/2)\log(|G|) - (1/2)b'G^{-1}b$

- The joint log(quasi)-likelihood is thus $\ell(b) + ql(y/b)$
- The marginal (quasi-) likelihood is

$$\iint_{\mathbf{b}} \left[ql(\mathbf{y} \mid \mathbf{b}) + \ell(\mathbf{b}) \right] d\mathbf{b} = ql(\mathbf{y})$$

$$= \iint_{\mathbf{b}} \left[\mathbf{y}' \mathbf{A} \mathbf{\theta} - \mathbf{1}' \mathbf{A} b(\mathbf{\theta}) - \left(\frac{b}{2}\right) \log(2\pi) - \left(\frac{1}{2}\right) \log\left(|\mathbf{G}|\right) - \left(\frac{1}{2}\right) \mathbf{b}' \mathbf{G}^{-1} \mathbf{b} \right] d\mathbf{b}$$

Approximation is required $\rightarrow pseudo-likelihood$ method or Laplace approximation or Gauss Hermite quadrature.

- The idea of *pseudo-likelihood* is to approximate the inverse link function by Taylor series expansion evaluated at $\tilde{\eta}$
- The Taylor series expansion: $h(\eta) \cong h(\tilde{\eta}) + \frac{\partial h(\eta)}{\partial \eta} \Big|_{\eta = \tilde{\eta}} (\eta \tilde{\eta})$
- If *D* = diag [(∂g(µ/b)/∂µ)] then (Jiang, 2007) h(η) ≈h(η̃) + *D̃*⁻¹(*Xβ* + *Zb* - *Xβ̃* - *ZĎ̃*) where *D̃* denotes *D* evaluated at η̃ = *Xβ̃* - *ZĎ̃*.
 Let *y** = η̃+*D̃*⁻¹[*y*-(µ́/*D̃*)] = *Xβ̃* + *ZĎ* + *D̃*⁻¹[*y*-h(η̃)]
 Then E(*y**|*b*) = *D̃*[h(η)-h(η̃)]+*Xβ̃*+*ZĎ̃* ≈ *Xβ*+*Zb* Var(*y**|*b*) = *DV*^{1/2}_µ*AV*^{1/2}_µ*D*



• Let us fit a binomial mixed model to the multibatch data:

X	В	atch 1		Batch2		Batch 3			Batch 4			
	Y	Fav	Ν	Y	Fav	Ν	Y	Fav	Ν	Y	Fav	Ν
0	95.6	15	21	96.6	18	21	96.6	16	19	96.6	18	21
3	96.9	13	17	96.7	14	16	96.5	13	19	96.8	18	21
6	98.5	19	23	96.3	14	17	97.7	17	22	96.4	14	20
9	99.0	14	17	96.3	17	20	98.3	23	27	96.6	14	19
12	100.2	18	23	96.5	15	20	99.1	16	21	96.8	14	19
18	101.9	19	27	96.4	14	22	100.5	11	16	96.9	11	20
24	104.1	15	20	95.7	12	22	101.2	13	18	97.1	11	16
36	107.8	14	21	95.2	11	25	103.3	10	17	97.4	10	21
48	111.5	13	18	94.6	5	26	105.9	6	16	97.4	5	19

1. Observations: $(Fav_{ij}/b_{0i}, b_{1i}) \sim B[N_{ij}, (\pi_{ij}/b_{0i}, b_{1i})]$

2. Model focus:
$$E(Fav_{ij}/b_{0i}, b_{1i}) = \pi_{ij}/b_{0i}, b_{1i}$$

- 3. Linear predictor: $\eta_{ij} = \beta_0 + b_{0i} + (\beta_1 + b_{1i}) X_{ij}$
- 4. Assumptions about b_{0i} and b_{1i} (if random)
- 5. Logit link: $log[\pi_{ij}/(1-\pi_{ij})] = \beta_0 + b_{0i} + (\beta_1 + b_{1i}) X_{ij}$,

Fixed effects: X Random effects : Batch

 $\beta_0 + b_{0i} + (\beta_1 + b_{1i}) X_{ij}$

setwd("d:/Prof.Kan/mixmodel Kan 1D") read.csv("multibatchdata.csv",h=T,as.is=T) -> dt library(lme4) Batch<- as.character(dt\$Batch) gm1 <- glmer(cbind(dt\$Fav, dt\$N- dt\$Fav) ~ dt\$X+ (1 | Batch), family = binomial) summary(gm1)

Generalized linear mixed model fit by maximum likelihood (Laplace Approximation) [glmerMod] Family: binomial (logit) Formula: $cbind(dt\Fav, dt\N - dt\Fav) \sim dt\X + (1 | Batch)$

AIC	BIC	logLik	deviance	df.resid	
151.8	156.5	-72.9	145.8	33	
Scaled re	siduals:				
Min	1Q	Median	3Q	Max	
-1.6315	-0.6042	0.1118	0.4402	3.0004	
Random	effects:				
Groups N	Name	Variance	Std.Dev.		
Batch (I	ntercept)	0.01395	0.1181		
Number of	of obs: 36, gro	oups: Batch,	4		
Fixed effe	ects:				
		Estimate	Std. Error	z value	Pr(> z)
(Intercept	t) 1.584756	0.149651	10.590	<2e-16 ***	
dt\$X		-0.045388	0.005963 -7	7.612	2.7e-14 ***

Binomial reponse





Let us fit a normal mixed model to the multibatch data:

X	Ba	tch 1		B	atch2		Ba	atch 3		Ba	atch 4	
	Y	Fav	Ν	Y	Fav	Ν	Y	Fav	Ν	Y	Fav	Ν
0	95.6	15	21	96.6	18	21	96.6	16	19	96.6	18	21
3	96.9	13	17	96.7	14	16	96.5	13	19	96.8	18	21
6	98.5	19	23	96.3	14	17	97.7	17	22	96.4	14	20
9	99.0	14	17	96.3	17	20	98.3	23	27	96.6	14	19
12	100.2	18	23	96.5	15	20	99.1	16	21	96.8	14	19
18	101.9	19	27	96.4	14	22	100.5	11	16	96.9	11	20
24	104.1	15	20	95.7	12	22	101.2	13	18	97.1	11	16
36	107.8	14	21	95.2	11	25	103.3	10	17	97.4	10	21
48	111.5	13	18	94.6	5	26	105.9	6	16	97.4	5	19

1. Observations: $(y_{ij}/b_{0i}, b_{1i}) \sim N(\mu_{ij}/b_{0i}, b_{1i}, \sigma^2)$

2. Model focus:
$$E(y_{ij}/b_{0i}, b_{1i}) = \mu_{ij}/b_{0i}, b_{1i}$$

- 3. Linear predictor: $\eta_{ij} = \beta_0 + b_{0i} + (\beta_1 + b_{1i}) X_{ij}$
- 4. Assumptions about b_{0i} and $b_{1i} \sim N(0, G)$
- 5. Identity link: $(\mu_{ij}/b_{0i}, b_{1i}) = \beta_0 + b_{0i} + (\beta_1 + b_{1i}) X_{ij}$

Fixed effects: X Random effects : Batch

 $\beta_0 + b_{0i} + (\beta_1 + b_{1i}) X_{ij}$

setwd("d:/Prof.Kan/mixmodel_Kan_1D")
read.csv("multibatchdata.csv",h=T,as.is=T) -> dt
library(lme4)
Batch<- as.character(dt\$Batch)
gm2 <- glmer(dt\$Y ~ dt\$X+ (1 | Batch), family =gaussian)
summary(gm2)</pre>

Linear mixed model fit by REML ['lmerMod'] Formula: dt\$Y ~ dt\$X + (1 | Batch) REML criterion at convergence: 176.2

Scaled residuals:

Min	1Q	Median	3Q	N
-2.41643	-0.47968	0.02534	0.57671	2

Random effects:

Groups Name	Variance	Std.Dev
Batch (Intercept)	6.332	2.516
Residual	5.857	2.420
Number of obs: 36, gro	oups: Batch, 4	1

Fixed effects:

	Estimate	Std. Error	t value
(Intercept)	96.50627	1.40072	68.90
dt\$X	0.13129	0.02867	4.58

lax 43510	Linear model							
Source	DF	Adj SS	Adj MS	F-Value	P-Value			
Х	1	122.838	122.838	264.83	0.000			
Batch	3	0.782	0.261	0.56	0.644			
X*Batch	3	168.579	56.193	121.15	0.000			
Error	28	12.987	0.464					
Total	35	492.943						





Model diagnostics

- Literatures on diagnostics of mixed models involving random effects is not extensive.
- Limited methodology is available, mostly regarding assessing the distribution of the random effects and errors.
- The methods may be classified as diagnostic plots and goodness-of-fit tests.

Diagnostic plots. Lange and Ryan (1989) proposed to examine a Q–Q plot of some standardized linear combinations

$$z_i = \frac{c'\tilde{\alpha}_i}{\{c'\operatorname{Var}(\tilde{\alpha}_i)c\}^{1/2}}, \qquad i = 1, \dots, m,$$

where c is a known vector. They argued that, through appropriate choices of c, the plot can be made sensitive to different types of model departures.



Goodness of fit tests. Few authors have developed tests for checking distributional assumptions involved in linear mixed models. Consider a mixed ANOVA model

$$y = X\beta + Z\alpha + \varepsilon$$

where for $1 \le i \le s$, $\alpha_i = (\alpha_{ij})_{1 \le j \le mi}$, where the $\alpha_{ij}s$ are *i.i.d.* ~ $(0, \sigma^2)$, and continuous distribution $Fi = Fi(\cdot /\sigma_i)$; and $\varepsilon = (\varepsilon_{\theta})_{1 \le j \le N'}$ where the $e_j s$ are *i.i.d.* ~ $(0, \tau^2)$, and continuous distribution $G = G(\cdot |\tau)$; and $\alpha_1, \ldots, \alpha_s, \varepsilon$ are independent.

We are interested in testing the following hypothesis,

$$H_0: F_i(\cdot | \sigma_i) = F_{0i}(\cdot | \sigma_i) , \ 1 \le i \le s,$$

and $G(\cdot | \tau) = G_0(\cdot | \tau);$

that is, the distributions of the random effects and errors, up to a set of unknown variance components $\sigma_1^2, \ldots, \sigma_l^2, \tau^2$, are as assumed. The model does not fit if we reject this hipothesis.



Mixed Models in Practice

Mixed Models in Practice



- Mixed model is very popular
- **Google search** (April 25th, 2016):
 - "mixed model" → 629,000 results (0.46 seconds)
 - "robust statistics" \rightarrow 197,000 results (0.40 seconds)
 - "AMMI" → 11,900 results (0.45 seconds) (Additive Main effect Multiplicative Interaction)
- The mixed models have been in many areas such as:
 - Biology, Education, Social and economics, Humanity, Physics and astronomy, Environment
- Next is several works in Bogor on mixed models which are promising but still in early stages.



- **Example 1.** (Notodiputro and Yahya, 2016)
- "The Gamma Mixed Model for Analyzing Length of Stay of Tuberculosis Patients at Gorontalo Hospital"
- Gorontalo is a province in eastern part of Indonesia and is considered having high incidence of tuberculosis.
- The length of stay of a patient is the response variable or the outcomes (Y measured in days).
- Usually distribution of duration (length of stay, Y) is skewed, and we may approximate using gamma distribution.
- The data was collected from a hospital in Gorontalo province and it consists of many predictors, namely room type, age category, sex, secondary diseases, rural or urban, doctor.
- The sample size is 612 patients, these are the patients during 2015.



- **Example 1.** (Notodiputro and Yahya, 2016)
- "The Gamma Mixed Model for Analyzing Length of Stay of Tuberculosis Patients at Gorontalo Hospital"
- Response var.: length of stay in the hospital (y) assumed to have (conditional on random effects) gamma distributions.
- Explanatory var.: room type (p), age category (u), sex (s), secondary diseases (k), rural or urban (a), doctor (d)
- The linear predictor:

$$\begin{split} \eta_{ijklmno} &= \mu + p_i + u_j + s_k + k_l + a_m + (pu)_{ij} \\ &+ (ps)_{ik} + (pk)_{il} + (pa)_{im} + (us)_{jk} + (uk)_{jl} \\ &+ (ua)_{jm} + (sk)_{kl} + (sa)_{km} + (ka)_{lm} + d_o \end{split}$$

where the distribution of random effects $d \sim N(0, \sigma_{\delta}^2)$





- Example 1. (Notodiputro and Yahya, 2016)
- The anova table

Source	DF	F	p-value
Room type (p)	1	5.51	0.019
Sex (<i>s</i>)	1	0.01	0.908
Age category (u)	2	0.03	0.972
Secondary disease (k)	1	3.58	0.059
Rural or urban (a)	1	8.22	0.004
<i>p*s</i>	1	0.30	0.582
<i>p*u</i>	2	0.95	0.388
p*k	1	13.91	0.000
p*a	1	8.78	0.003
s*u	2	0.54	0.581
s*k	1	0.01	0.925
s*a	1	0.06	0.809
u*k	2	2.06	0.129
u*a	2	0.11	0.900
k*a	1	1.39	0.239

n = 612 patients $s_{\delta}^2 = 0.00679$

 $s_{c}^{2} = 2.52600$

- Variation within doctor is smaller than variation between doctors
- Room type is significant: the patients in better room stay longer
- Interaction p^*k and p^*a are significant \rightarrow further exploration is required



• Example 1. (Notodiputro and Yahya, 2016)



- Example 2. (Notodiputro and Adabiyah, 2016)
- "Nested Linear Mixed Model and Parametric Stability Analysis for Multilocation Experiments of Shorgum Genotypes"



- 10 shorghum genotypes were evaluated based on their yields
- 3 sorghum were used as control
- Randomized block experiments were carried out in two seasons and 5 different locations within the seasons.
- Eventually, we want to know the performance of these genotypes and which of them that can produce stable yields.



- Example 2. (Notodiputro and Adabiyah, 2016)
- "Nested Linear Mixed Model and Parametric Stability Analysis for Multilocation Experiments of Shorgum Genotypes"
- Response var.: Shorgum yield (y) assumed to have (conditional on random effects) normal distributions
- Explanatory var.: genotypes (G), seasons (S), locations (L), blocks (B)
- The linear predictor:

 $\eta_{ijkl} = \mu + G_i + S_j + L_{k(j)} + (GS)_{ij} + (GL)_{ik(j)} + B_{l(jk)} + \varepsilon_{ijkl}$

where the distribution of random effects $B \sim N(0, \sigma_B^2)$; $L \sim N(0, \sigma_L^2)$; and $GL \sim N(0, \sigma_{GL}^2)$





Indonesia

- Locations:
- Yogyakarta
- Depok
- Wonosari
- Telukbetung
- Mataram
- Bogor
- Boyolali
- Tj.Karang
- Pekanbaru



- Example 2. (Notodiputro and Adabiyah, 2016)
- The anova table

13x3x5x2 = 390 experimental units



Random Effects	Var.est	Std eror	Z	p-value
Location (Season)	0.2248	0.1203	1.8700	0.0309
Genotype*Location (Season)	0.0882	0.0316	2.7900	0.0026
Block (Season*Location)	0.0000	•		
Error	0.3554			

	Fixed Effects	DF for numerator	DF for denominator	F	P-value
Č	Genotype	12	96	54.4700	<.0001
	Season	1	8	0.3100	0.5923
	Genotype*Season	12	96	1.3500	0.2027

Example 2. (Notodiputro and Adabiyah, 2016)

I

Ď

5.5

6.0

6.5

Canatinaa		Yield					
Genotypes	CV_i	W_i^2	d_i^2				
G1	19.5814	21.3965	10.5619	o _		•	
G2	20.4037	21.3650	8.8558	Ŋ	•		
G3	22.2235	25.1734	12.2653		IV)	
G4	17.9484	23.1228	7.1587	- 5			
G5	6.4129	14.4406	4.1602	CV_i			
G6	14.6978	24.7314	4.6860	-	M		
G7	6.9182	11.8964	3.1412	- 15	III		K
G8	8.0645	12.1867	4.5061				Ē
G9	19.2566	19.7342	7.4759				
G10	12.7184	29.7786	5.3025		35 40	4.5	5.0
Kawali (K)	9.5080	15.3093	4.6290		0.0 4.0	ч.5	7: -1-1
Mandau (M)	11.4640	27.5873	2.0947)	ield
Pahat (P)	9.5229	14.9245	3.4744				

 CV_i is coefficient of variation; W_i^2 is ecovalence; and d_i^2 is genotypic stability

Example 3. (Angraini and Notodiputro, 2016)





- **Example 3.** (Angraini and Notodiputro, 2016)
- "Generalized Linear Mixed Models for Analyzing Fish Stock at Na Thap River"
- Response var.: Fish stock (y) assumed to have (conditional on random effects) Gamma distributions
- Explanatory var.: Zone (Z), WDEPTH (W), SAL (L), DO (D)
 BOD (B), and Site (S).
- The linear predictor:

 $\eta_{ij} = \mu + \beta_1 Z_i + \beta_2 W_i + \beta_3 L_i + \beta_4 D_i + \beta_5 B_i + S_j$

where the distribution of random effects $S_i \sim N(0, \sigma_s^2)$;

- Or, η = Xβ + Zs and s ~ N(0,G)
- Link function: log







- Example 3. (Angraini and Notodiputro, 2016)
- The pairwise confidence intervals:

Site		Substracted from Site			
	2	3	4		
1	(-7909,4741)	(-774,11876)	(-8336,4315)		
2		(810,13460)	(-6751,5899)		
3			(-13886,-1236)		

- The power plant is in site 3
- Site 2 vs 3 is significant (site 2 > site 3)
- Site 3 vs 4 is significant (site 4 > site 3)



- Example 4. (Arisanti and Notodiputro, 2016)
- "Bias Reduction in Estimating Variance Components of Phytoplankton Abundance at Na Thap River based on Logistic Linear Mixed Models"
- Response var.: Phytoplankton abundance (y) assumed to have (conditional on random effects) binomial distributions
- Explanatory var.: SAL (L), DO (D) BOD (B), and Site (S).
- The linear predictor:

$$\eta_{ij} = \mu + \beta_3 L_i + \beta_4 D_i + \beta_5 B_i + S_j$$

where the distribution of random effects $S_i \sim N(0, \sigma_s^2)$;

- Or, η = Xβ + Zs and s ~ N(0,G)
- Link function: logistics



- **Example 4.** (Arisanti and Notodiputro, 2016)
- It is well known that variance estimates of MLE are biased
- We follow the idea of Firth (1993) to adjust variance components in GLMM with logistic link function



If $U(\theta)$ is the score function in ML estimation, and $\hat{\theta}$ is subject to a positive bias $b(\theta)$, the score function is shifted downward at each point θ by amount $i(\theta)b(\theta)$, where $-i(\theta)=U'(\theta)$ is the local gradient, then the adjusted score function:

 $U^{*}(\theta) = U(\theta) - i(\theta)b(\theta)$

<u>O</u>r

$$U^{*}(\theta) = U(\theta) + A(\theta)$$



- Example 4. (Arisanti and Notodiputro, 2016)
- The adjusted score function:

 $U^{*}(\theta) = U(\theta) - i(\theta)b(\theta) \text{ or } U^{*}(\theta) = U(\theta) + A(\theta)$

- The solution to this adjusted score function is θ^*
- Firth adjustment is based on F and H are Fisher's information and Hessian matrix

$$\boldsymbol{A}_{\theta j} = -\frac{1}{2} tr \big(\boldsymbol{F}^{-1} \boldsymbol{E} \big[\boldsymbol{U}_{\theta j} (\boldsymbol{H} - \boldsymbol{U} \boldsymbol{U}') \big] \big)$$



- Example 4. (Arisanti and Notodiputro, 2016)
- The ML estimate of variance component of site:

 $\hat{\sigma}^2$ site = 0.986

$$n = 466$$
 (Synedra)

The adjusted estimate is

 $\tilde{\sigma}^2$ site = 0.735

- The variance component has been reduced by 26%
- The analysis of variance table:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	2.0288	0.3682	5.511	0.0000
DO	-0.3153	0.0775	-4.067	0.0000
BOD	0.1325	0.0998	1.327	0.1840
Salinitas	-0.0108	0.0106	-1.021	0.3070



- **Example 5.** (Angraini and Notodiputro, 2016)
- "A Hierarchical Approach to Generalized Linear Mixed Model for Analyzing Fish Species Abundance at Na Thap River"
- Response var.: Fish species abundance (y) assumed to have (conditional on random effects) Poisson (A) distributions.
- Explanatory var.: SAL (L), DO (D) BOD (B), and Site (S).
- The linear predictor:

 $\eta_{ijk} = \mu + \beta_1 Z_i + \beta_2 W_j + \beta_3 L_j + \beta_4 D_j + \beta_5 B_j + S_k$ where the distribution of random effects $S_k \sim Gamma(\alpha)$

- Link function: log
- The hierarchical likelihood

 $L(y|\lambda)+L(\alpha;u) = \sum (y_{ijk} ln\lambda_{ij} - \lambda_{ij}) + \sum \{\alpha lnu_{ij} + \alpha ln \alpha - \alpha u_{ij} - \alpha ln \Gamma(\alpha)\}$ IPB - PSU Collaborative Research RPM ID16287; The Na Thap river Project



• Example 5. (Angraini and Notodiputro, 2016)

		Hierarchical Model				
	Estimate	Std.error	Т	Pr(> t)		
Intersep	3.021	0.180	16.826	0.000		
Zonefreshwater	-0.242	0.156	-1.558	0.120		
Zonesaline	0.344	0.167	2.056	0.040	n	
WTEMP	0.002	0.004	0.552	0.581		
WDEPTH	-0.010	0.010	-0.945	0.345		
SAL	0.006	0.001	8.004	0.000		
DO	0.026	0.007	3.507	0.000		
BOD	0.025	0.007	3.619	0.000		
σ^2 Site	0.709	Variation	mongeitee	ovidont		
σ^2 Sisaan	0.041	variation a	iniong sites was	evident		

n = 531



Concluding Remarks

Concluding Remarks



References



- 1. Firth, D. 1993. *Bias reduction of maximum likelihood estimation. Biometrika,* Vol. 80 No. 1. pp 27-38
- 2. Jiang, J. 2007. *Linear and Generalized Linear Mixed Models and Their Applications*. Springer.
- Lee, Y. and Nelder. J.A. 1996. Hierarchical generalized linear models. *Journal of the Royal Statistical Society*. Series B, (Methodological), Vol. 58, No. 4(1996), pp. 619-678
- Mehari M., Alamerew S., Lakew B., Yirga H., Tesfay M. 2014. Parametric stability analysis of malt barley genotypes for grain yield in Tigray, Ethiopia. *World Journal of Agricultural Sciences*. 10(5): 210-215
- 5. Saheem, N. (2015). *Statistical Modeling of Aquatic Animal Abundance in the Na Thap River*. Thesis. Prince of Songkla Univ.
- 6. Stroup, W. W. 2013. *Generalized Linear Mixed Models: Modern Concepts, Methods, and Applications.* Chapman & Hall







Syiah Kuala University conference in Aceh



AND THEIR APPLICATIONS

The 12th ICMSA 2016 **ABSTRACT SUBMISSION**

The Department of Mathematics of Syiah Kuala University invite you to join the 12th ICMSA conference that will be held in Banda Aceh, Indonesia on October 4-6, 2016. This conference will be held in

http://icmsa.unsyiah.ac.id