Contributeurs

MAS-II: Modern Actuarial Statistics II

aut., cre. Alec James van Rassel

Référence (manuels, YouTube, notes de cours) En ordre alphabétique :

Contributeurs

Contents	The Marginal Linear Model
	The Implied Marginal Model

		The Implied Marginal Model	12
A Prerequisites	3	Estimation REML estimation	12 12
Distributions Discrete Distributions	3 5	Algorithms Troubleshooting	13 13
Joint Distributions	5 5 5 5	Hypothesis Testing Likelihood Ratio Tests	14 14 14 14
Statistics	5	Model-Building Strategies	14
B Introduction to Credibility	6	Checking Model Assumptions	15
Basic Framework of Credibility Full Credibility	6 6	EBLUPS Intra Correlation Coefficient	15 15 15
Claim Severity	7 7	Information Criteria	15
Partial Credibility	7	Checking model assumptions	16
Bühlmann Credibility Basic framework	8 8 8	Graphical Tests	17
Credibility factors	8	D Bayesian Analysis and Markov Chain Monte Carlo	19
Bayesian Credibility Basic framework	9 9	E Statistical Learning	20
Conjugate distributions	9	K-Nearest Neighbors	20
Tomparametric empirical Dayes method	3	Decision Trees	21
C Linear Mixed Models	10	Principal Components Analysis (PCA)	22
General Theory Model Assumptions	10 10 10 11	Clustering	23

A

Prerequisites

Distributions

Context

We typically use 3 types of random variable to describe losses:

Frequency or number of losses	always discrete	
Severity or amount of losses (payment)	usually continuous, can be discrete or mixed too	
Aggregate or total loss from summing a number (Frequency) of Severity vari- ables	same as the severity	

$\Pr(N = n)$ = $\frac{e^{-\lambda}\lambda^n}{n!}$ E[N] = λ Var(N) = λ

- > If N_1 and N_2 are independent Poisson r.v., then $N_1 + N_2 \sim \text{Poisson}(\lambda_1 + \lambda_2)$.
- > The $e^{-\lambda}$ term makes the probabilities sum to 1 as the Taylor series for e^{λ} is

$$e^{\lambda} = 1 + \lambda + \frac{\lambda^2}{2!} + \dots + \frac{\lambda^n}{n!} + \dots$$

Discrete Distributions

Context

Discrete random variables are usually counting (frequency) variables, meaning their possible values are $\{0,1,2,\dots\}$

Probability Mass Function (PMF)

N is a discrete random variable if it has a **probability mass function** p_k such that $p_k = \Pr(N = k)$

Definition	Domain	Condition
$p_k = \Pr(N = k)$	$p_k \in [0,1]$	$\sum_{k} p_{k} = 1$

▼ Poisson Distribution

Notation	Parameters	Domain
$N \sim \text{Poisson}(\lambda)$	$\lambda > 0$	$n = 0, 1, 2, \dots$

▼ Binomial Distribution

Context

A binomial r.v. N has m independent trials each having a probability q of a loss where n is the total number of losses.

Notation	Parameters	Domain
$N \sim \operatorname{Bin}(m,q)$	$q \in (0,1); m \in \mathbb{N}$	$n = 0, 1, 2, \dots$

Pr(N = n)	$= \binom{m}{n} q^n (1-q)^{m-n}$
E[N]	= mq
Var(N)	= mq

- > If N_1 and N_2 are independent binomial r.v. with the same q then $N_1 + N_2 \sim \text{Bin}(m_1 + m_2, q)$.
- \succ The case where m=1 corresponds to a $\boldsymbol{Bernoulli}$ r.v.

∨ Geometric Distribution

Context

A geometric r.v. N with mean β can be obtained by setting n as the number of years **before** the <u>first</u> loss. Given the geometric distribution is memoryless, each year independently has a loss with probability

$$\underbrace{\Pr(N=0)}_{\text{probability of a}} = \frac{1}{1+\beta}.$$
loss the first year

Notation	Parameters	Domain
$N \sim \text{Geo}(\beta)$	$\beta > 0$	$n = 0, 1, 2, \dots$

Pr(N=n)	$= \left(\frac{\beta}{1+\beta}\right)^n \frac{1}{1+\beta}$
$\Pr(N \ge n)$	$= \left(\frac{\beta}{1+\beta}\right)^n$
E[N]	$=\beta$
Var(N)	$=\beta(1+\beta)$

> Like the exponential distribution, the geometric distribution is memoryless:

$$\Pr(N = d + n | N \ge d) = \Pr(N = n)$$

$$E[N - d|N \ge d] = E[N]$$

lacksquare

Negative Binomial Distribution

Context

A negative binomial r.v. N represents the number of years n with no loss before the $r^{\rm th}$ year with a loss. We obtain a negative binomial r.v. $N \sim {\rm NBin}(r,\beta)$ by summing r iid geometric r.v., N_1,N_2,\ldots,N_r , all with the same mean β .

Notation	Parameters	Domain
$N \sim NBin(\beta)$	$r, \beta > 0$	$n = 0, 1, 2, \dots$

$\Pr(N=n)$	$= {r+n-1 \choose r-1} \left(\frac{\beta}{1+\beta}\right)^n \left(\frac{1}{1+\beta}\right)^r$
$\Pr(N \ge n)$	$=\left(rac{eta}{1+eta} ight)^n$
E[N]	$=r\beta$
Var(N)	$=r\beta(1+\beta)$

 \rightarrow A geometric r.v. is a negative binomial r.v. with r=1.

Distribution	Mean		Variance
Binomial	mq	>	mq(1-q)
Poisson	λ	=	λ
Geometric	β	<	$\beta(1+\beta)$
Negative Binomial	rβ	<	$r\beta(1+\beta)$

Severity Distributions

Joint Distributions

Conditional Distributions

Aggregate Distributions

Normal, Uniform, Pareto, Exponential, and Gamma

$$\gamma(1/2) = \sqrt{\pi}$$

Statistics

Mode

Context

The mode is the value that occurs the most often. A non-mathematical example of the concept is looking at the most used letter in the English alphabet. The letter E is the most used letter in the dictionary and as such is the mode of the English language.

In mathematical terms, the mode is the point which maximises the PMF/PDF.

Finding the mode of a continuous r.v. can be done by calculating the derivative of the PDF and finding the point where it equals 0. If the distribution is

- > unimodal, i.e. it has a hump, then mode = x s.t. f'(x) = 0.
- > strictly increasing or decreasing, the mode will be one of the 2 extremes.
 - For example, the exponential distribution is strictly decreasing and its mode is always 0.

For discrete variables, there are some ways to simplify it's calculation:

- > Using the table function on the calculator and seeing where the probabilities peak.
- \gt Using the algebraic approach of looking at $p_k/p_{k-1}.$
 - $-p_k > p_{k-1}$ iff $p_k/p_{k-1} > 1$.
 - The mode is the largest k s.t. $p_k > p_{k-1}$.

Note In the exam, it's best to use the calculator approach.

\mathbf{B}

Introduction to Credibility

Basic Framework of Credibility

Context

The *limitation fluctuation credibility* approach, or *classical credibility* approach, calculates an updated prediction (U) of the **loss measure** as a weighted (Z) average of recent claim experience (D) and a rate (M) specified in the manual. Thus, we calculate the *premium* paid by the *risk group* as U = ZD + (1 - Z)M.

Notation

M Predicted loss based on the "manual".

D Observed losses based on the recent experience of the risk group.

Z Weight assigned to the recent experience D called the *credibility factor* with $Z \in [0,1]$.

U Updated prediction of the premium.

Terminology

Risk group block of insurance policies, covered for a period of time upon payment of a *premium*.

Claim frequency The number of claims denoted N.

Claim severity The amount of the i^{th} claim denoted X_i .

Aggregate loss The total loss denoted S where $S = X_1 + X_2 + ... + X_N$.

Pure premium The pure premium denoted P where P = S/E with E denoting the number of exposure units.

Exam tips

Typical questions about this involve being given 3 of M, D, Z, and U then finding the missing one.

Context

With $\min\{D, M\} \le U \le \max\{D, M\}$, we can see that the credibility factor determines the relative importance of the claim experience of the risk group D relative to the manual rate M.

If Z=1, we obtain $\overline{Full\ Credibility}$ where the predicted premium depends only on the data $\overline{(U=D)}$. It follows that with Z<1, we obtain $Partial\ Credibility$ as the weighted average of both D and M.

Full Credibility

Contexte

The classical credibility approach determines the $minimum\ data\ size$ required for the experience data (D) to be given $full\ credibility$. The minimum data size, or $standard\ for\ full\ credibility$, depends on the loss measure.

Claim Frequency

The claim frequency random variable N has mean μ_N and variance σ_N^2 . If we assume $N \approx \mathcal{N}(\mu_N, \sigma_N^2)$, then the probability of observing claim frequency

within
$$k$$
 of the mean is $\Pr(\mu_N - k\mu_N \le N \le \mu_N + k\mu_N) = 2\Phi\left(\frac{k\mu_N}{\sigma_N}\right) - 1$.

We often assume that the claim frequency $N \sim \text{Pois}(\lambda_N)$ and then apply the normal approximation to find the standard for full credibility for claim frequency λ_F . First, we impose that the probability of the claim being with k of the mean must

be at least $1 - \alpha$. Then, we rewrite $\frac{k\mu_N}{\sigma_N} = k\sqrt{\lambda_N}$ and set $\lambda_N \ge \left(\frac{z_{1-\alpha/2}}{k}\right)^2$ where

$$\lambda_F = \left(\frac{z_{1-\alpha/2}}{k}\right)^2$$

Claim Severity

We assume that the loss amounts $X_1, X_2, ..., X_N$ are independent and identically distributed random variables with mean μ_X and variance σ_X^2 . Full credibility is

attributed to $D = \bar{X}$ if $2\Phi\left(\frac{k\mu_X}{\sigma_N/\sqrt{N}}\right) - 1 \ge 1 - \alpha$.

Similarly to claim frequency, we apply the normal approximation with

$$\bar{X} \approx \mathcal{N}\left(\mu_X, \sigma_X^2/N\right)$$
. Then, we find $N \geq \left(\frac{z_{1-\alpha/2}}{k}\right)^2 \cdot \left(\frac{\sigma_X}{\mu_X}\right)^2 = \lambda_F C V_X^2$ where the

standard for full credibility for claim severity is $\lambda_F CV_X^2$.

Aggregate Loss

For the aggregate loss $S=X_1+X_2+\ldots+X_N$, we have $\mu_S=\mu_N\mu_X$ and $\sigma_S^2=\mu_N\sigma_X^2+\mu_X^2\sigma_N^2$.

With the same normality assumptions for the Poisson distributed N, we find

$$\lambda_N \geq \left(\frac{z_{1-\alpha/2}}{k}\right)^2 \cdot \left(\frac{\mu_X^2 + \sigma_X^2}{\mu_X^2}\right) = \lambda_F (1 + CV_X^2)$$
 where the **standard for full cred**-

ibility for claim severity is $\lambda_F(1+CV_X^2)$.

Note The conditions are the same for the $\it Pure \ Premium$ as for the aggregate loss.

Partial Credibility

The $\boldsymbol{credibility\ factor}$ for :

Claim Frequency is $Z = \sqrt{\frac{\lambda_N}{\lambda_F}}$

Claim Severity is $Z = \sqrt{\frac{N}{\lambda_F C V_X^2}}$

Aggregate Loss and Pure Premium is $Z = \sqrt{\frac{\lambda_N}{\lambda_F(1+CV_X^2)}}$

Bühlmann Credibility

Context

Buhlmann's approach, a.k.a. the greatest accuracy approach or the least squares approach, estimates the future loss measure X_n

Basic framework
Variance components
Credibility factors

Bayesian Credibility

Basic framework

Premium

Conjugate distributions

Nonparametric empirical Bayes method

\mathbf{C}

Linear Mixed Models

Context

What distinguishes a linear mixed model is that it may include both **fixed-effect parameters** and **random effects**. The mix of these gives the linear *mixed* model its name. Fixed-effect parameters describes the relationships of the covariates to the dependant variable for an *entire population*. Random effects are specific to clusters or subjects *within a population*. Random effects are thus directly used in modelling the random variation in the dependant variable at different levels of the data.

Fixed factors are categorical or classification variables for which all levels (conditions) that are of interest are included. Random factors can be thought of as being randomly sampled from a population of levels being studied. The text gives as an example the Dental Veneer case study where if we specified the tooth being sampled, selected teeth would become a fixed factor. This would however limit inferences by teeth rather than generalizing to "teeth within a patient".

The case studies use 3 types of data:

clustered The dependant variable is measured once per subject (unit of analysis), and the units are grouped into/nested within clusters of units.

- > We can have data sets that are two-level (e.g. rat pup data set), three-level (e.g. classroom data), etc.
- > For MAS-II, we shouldn't have beyond three levels.

repeated-measures The dependant variable is measured more than once (on the same unit of analysis) across levels of a repeated-measures factor(s). (e.g. time, measurement conditions, etc.)

longitudinal The dependant variable is measured at several points in time for each unit of analysis.

- > Clustered longitudinal data combines features of both. (e.g. Dental Veneer data set).
- > Each unit is measured more than once, but those units of analysis are nested within clusters.

Context

These 3 are *hierarchical* data sets as the observations can be placed into levels of a hierarchy in the data.

Generally:

Level 1 most detailed level; subjects, repeated measures on the same unit of analysis.

Level 2 clusters of units, units of analysis.

Level 3 clusters of clusters, clusters of units.

Levels are emphasized in the text because they help to conceptualize LMM as simple models defined at each level of the data hierarchy.

General Theory

Model Assumptions



Fixed Effects



Random Effects

Model specification

$$Y_i = \underbrace{X_i oldsymbol{eta}}_{ ext{fixed}} + \underbrace{Z_i u_i}_{ ext{random}} + arepsilon_i$$

where:

- \rightarrow Y_i is the $n_i \times 1$ vector of continuous responses for the *i*-th subject.
- $\rightarrow X_i$ is the $n_i \times p$ design matrix.
- $\gt Z_i$ is the $n_i \times q$ design matrix.
 - In a LMM in which only the intercepts are assumed to vary randomly from subject to subject, the Z_i matrix would simply be a column of 1's.
- $> u_i$ is the $q_i \times 1$ matrix (vector) of random effects.

- \rightarrow D $q \times q$ variance-covariance matrix of u_i .
- $\succ \varepsilon_i$ is the $n_i \times 1$ matrix (vector) of random effects.
- $\rightarrow R_i n_i \times n_i$ variance-covariance matrix of ε_i .

Note When LMMs are specified in terms of an explicitly defined hierarchy of simpler models, they are often referred to as hierarchical linear models (HLMs) or multilevel models (MLMs).

Note The elements of both R_i and D are defined as functions of another set of \rightarrow Implies observations closer to each other have a higher correlation than those that covariance parameters θ .

Random Covariance Structures

unstructured D matrix with no additional constraints on the values of its elements.

> Often used for random coefficient models.

diagonal Each random effect u_i has its own variance, and all covariances in D are assumed to be zero.

Residual Covariance Structures

Diagonal

$$R_i = \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix}$$

- > Assumes residuals from the same subject are *uncorrelated* with equal variance.
- > Often the default structure.

Compound Symmetry

$$R_{i} = \begin{bmatrix} \sigma^{2} + \sigma_{1} & \sigma_{1} & \dots & \sigma_{1} \\ \sigma_{1} & \sigma^{2} + \sigma_{1} & \dots & \sigma_{1} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1} & \sigma_{1} & \dots & \sigma^{2} + \sigma_{1} \end{bmatrix}$$

- > Assumes equal correlation between residuals from the same individual. (e.g. repeated trials under the same condition in an experiment).
- > Good for clustered or repeated measures data.

First Order Auto-Regressive (AR(1))

$$R_{i} = \begin{bmatrix} \sigma^{2} & \sigma^{2}\rho & \dots & \sigma^{2}\rho^{n_{i}-1} \\ \sigma^{2}\rho & \sigma^{2} & \dots & \sigma^{2}\rho^{n_{i}-2} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma^{2}\rho^{n_{i}-1} & \sigma^{2}\rho^{n_{i}-2} & \dots & \sigma^{2} \end{bmatrix}$$

- \rightarrow Means adjacent residuals have a covariance of $\sigma^2 \rho$.
- > Good for longitudinal observations with equal time between observations.
- are further apart in time.

Note We can allow heterogeneous variances for different groups of subjects (e.g. males and females). We could assume the same structure but with different values.

The Marginal Linear Model

Context

Random effects are explicitly used in LMMs to explain the between-subject (between-cluster) variation, but they're not used in specifying marginal models. We thus refer to LMMs as subject-specific models, and marginal models as population-averaged models.

We specify the model as $Y_i = \underbrace{X_i oldsymbol{eta}}_{ ext{fixed}} + oldsymbol{arepsilon}_i$ where the $marginal\ residual\ errors$ $\varepsilon_i \sim \mathcal{N}(0, V_i^*).$

The big difference here is that the entire random part of the marginal model is described in terms of the marginal residuals ε_i^* only.

Note All structures used for R_i can be used for V_i^* . There are others that can be used, such as the one defined for the implied marginal model.

The Implied Marginal Model

Context

The concept of the implied marginal model is important for at least 2 reasons specified in the text:

- 1. The framework of the implied marginal model is used to estimate the fixed-effect and covariance parameters in the LMM.
- 2. When we obtain an invalid estimate of the **D** matrix from a software, we can try to fit the implied marginal model (with fewer restrictions) to potentially diagnose problems with nonpositive-definiteness.

The marginal model implied by the LMM has the variance-covariance matrix $V_i = Z_i D Z_i' + R_i.$

- > Both the LMM and the implied marginal model involve the same set of covariance parameters θ , however there are more restrictions imposed on it in the LMM.
- > Whereas both R_i and D have to be positive-definite, only V_i has to be for the implied marginal model.

Estimation

When θ is assumed to be known, we apply the method of generalized least squares (GLS) to estimate β with MLE. The estimate $\hat{\beta}$ is the **best linear unbi**ased estimator (BLUE) of β . The Empirical Best Linear Unbiased Estimator (EBLUE) of β replaces V_i by its estimate \hat{V}_i .

We obtain that
$$\operatorname{Var}(\hat{\pmb{\beta}}) = \left(\sum_i \pmb{X}_i' \hat{\pmb{V}}_i^{-1} \pmb{X}_i\right)^{-1}$$
.

REML estimation

REML is often preferred to ML estimation because it produces unbiased estimates of covariance parameters. It does this by taking into account the loss of df that results from estimating the fixed effects in β .

The resulting estimated $\hat{\boldsymbol{\beta}}$ and $\operatorname{Var}(\hat{\boldsymbol{\beta}})$ differ because $\hat{\boldsymbol{V}}_i$ is different.

Context

The variance of the estimated fixed effects (diagonal elements of $Var(\hat{\beta})$) are biased downward in both ML and REML estimation. This is because neither methods takes into account the uncertainty introduced replacing V_i by \hat{V}_i .

It follows that $se(\hat{\beta})$ is also biased downward.

Algorithms

Expectation Maximization

The underlying assumption behind the EM algorithm is that optimization of the complete data log-likelihood function is simpler than optimization of the likelihood based on the observed data.

Its main drawback is its slow rate of convergence. Also, the precision of the estimators is overly optimistic. This because the estimators are based on the likelihood from the last maximization step which uses complete data instead of observed data.

Usually used to provide starting values for other algorithms.

Newton-Raphson

Most commonly used in ML and REML estimation of LMMs.

While iterations are more time consuming (given Hessian matrix calculations), but convergence is quicker than the EM algorithm. Another advantage is that the last iteration's Hessian matrix can be used to obtain an asymptotic variance-covariance matrix for the estimated covariance parameters $\boldsymbol{\theta}$ to calculate the se($\hat{\boldsymbol{\theta}}$).

Fisher Scoring Algorithm

Modification of the N-R that uses the $\boldsymbol{expected}$ Hassian matrix rather than the observed one.

Its advantages are that it's more stable numerically, more likely to converge, and has simpler calculations at each iteration. However, it is not recommended to obtain the final estimates. It's primary disadvantage is that it may be difficult to determine the expected value of the Hessian matrix owing to difficulties identifying the appropriate sampling distribution.

Troubleshooting

If there are problems fitting the model, these are the steps you can take:

1. Choose alternative starting values for covariance parameter estimates

- 2. Rescale the covariates
 - > Improves numerical stability of the optimization algorithm and may circumvent convergence problems.
- 3. Simply and remove some random effects
 - > Generally, start with higher-order terms.
 - > Though, it can be valid to remove lower-order terms, but requires thorough justification.
- 4. Fit the implied marginal model
- 5. Fit the marginal model with an unstructured covariance matrix

Methods 4 and 5 shift a more restrictive requirement for the D and R_i matrices to be positive-definite to a less restrictive requirement that V_i be positive-definite.

Hypothesis Testing

Likelihood Ratio Tests

Mixture of Chi Squares REML

Non-Likelihood Ratio Tests



We compensate for the fact that the standard error is underestimated with REML/ML estimation by using an estimated number of df. We approximate the df hence is is not simply n - p.

F-test

To test $\mathcal{H}_0: L\beta = 0$ vs $\mathcal{H}_1: L\beta \neq 0$.

Context

Degrees of freedom of the numerator correspond

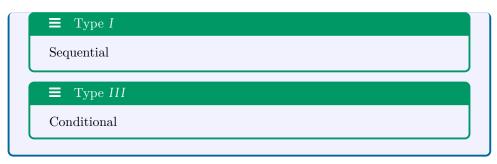
We get that the test statistic $t \approx F_{\text{num. df,den. df}}$ where the numerator df corresponds to the number of parameters being tested and the denominator degrees of freedom is obtained from R.

The particularity of the F-test is that we must make adjustments for it due to:

- 1. Random Effects
- 2. Potential correlation between residuals
- 3. Estimate covariance matrix

We have a few ways of approximating them:

- \succ Method used by R.
- > Method used by SAS.



Using tests:

- 1. Compute test statistic
 - > F-statistic would be too hard to compute, would have to be provided.
- \rightarrow For t-test, may just give components of the calculation and have us compute t to the compare it to the CV.
- > For both tests, the number of df would have to be provided.
- 2. Look up critical value table
- 3. Reject null / keep effects if test statistic > CV

Other tests

Omnibus Wald Test (good) similar to F-test. Numerator of the F-test statistic. test statistic asymptotically χ^2

Wald z-test (not good) only good asymptotically and breaks in some situations text recommends LRT instead

Model-Building Strategies

1

Top-Down Strategy

Start with a model that includes the maximum number of fixed effects that we wish to consider. The steps to build the *loaded* mean structure are:

- 1. Start with a well-specified mean structure for the model
 - > Want to ensure that the systematic variation in the responses is well explained before investigating various covariance structures to describe random variation in the data.
- 2. Select a structure for the random effects in the model

- > Selecting set of random effects to include in the model with REML-based LRT.
- 3. Select a covariance structure for the residuals in the model
 - > Variation remaining after both fixed and random effects have been added to the model due to residual error.
 - > Investigate an appropriate covariance structure for the those residuals.
- 4. Reduce the model
 - > Use appropriate statistical tests to determine whether certain fixed-effect parameters are needed in the model.

2 Step-Up Strategy

- 1. Start with an "unconditional" (or means-only) Level 1 model for the data
- 2. Build the model by adding Level 1 covariates to the Level 1 model. In the Level 2 model, consider adding random effects to the equations for the coefficients of the Level 1 covariates.
- 3. Build the model by adding Level 2 covariates to the Level 2 model. For 3-level models, consider adding random effects to the Level 3 equations for the coefficients of the Level 2 covariates.

Checking Model Assumptions

In general, raw conditional residuals in their basic form are not well suited for verifying model assumptions and detecting outliers. They tend to be correlated and their variances may be different for different subgroups of individuals.

Standardized residuals are conditional residuals divided by their true standard deviations. Unfortunately, the true SDs are rarely known in practice and we use estimated SDs instead to obtain studentized residuals.

Scaling the residuals by dividing them by the estimated SD of the dependent variable produces Pearson residuals. This is appropriate when we assume we can ignore the variability of $\hat{\beta}$.

internal vs external studentization.

Influence diagnostics are formal techniques to identify observations that heavily influence estimates of the parameters in either β or θ .

EBLUPS

Intra Correlation Coefficient

$$ICC_{\mathrm{whatever}} = \frac{\mathrm{variance\ in\ common}}{\mathrm{total\ variance}}$$

2 level model $ICC_{\text{group}} = \frac{\sigma_{\text{lvl 2}}^2}{\sigma_{\text{lvl 2}}^2 + \sigma^2}$ 3 level model $ICC_{\text{lvl 3 group}} = \frac{\sigma_{\text{lvl 3}}^2}{\sigma_{\text{lvl 3}}^2 + \sigma_{\text{lvl 2}}^2 + \sigma^2}$ $ICC_{\text{lvl 2 group}} = \frac{\sigma_{\text{lvl 3}}^2 + \sigma_{\text{lvl 2}}^2 + \sigma^2}{\sigma_{\text{lvl 3}}^2 + \sigma_{\text{lvl 2}}^2 + \sigma^2}$

EBLUPS

EBLUP

 ${f E}$

 ${f B}$ Best i.e. lowest variance among all such unbiased estimators

L Linear as functions of y_i

U Unbiased with $E[\hat{u}_i] = u_i$

Р

> Typically tedious to calculate so we use computers unless we calculate only for 1 random effect.

Use Buhlmann's formula where:

M Average predicted value from the implied marginal model

 \bar{Y} Average observed value from group

$$\sigma_{HM}^2 \operatorname{Var}(u_j) = \sigma_{int}^2$$

$$\mu_{PV} \operatorname{Var}(\varepsilon_{ij}) = \sigma^2$$

Prediction is for $M + u_i = M + Z_i(\bar{Y} - M)$.

Information Criteria

Context

When comparing 2 nested models, the more complex will be better than the simpler model. While the <u>Likelihood Ratio Tests</u> checks if the simpler model is sufficient, it does not enable us to directly compare the 2 models. In addition, with the LRT we are limited to nested models. The AIC and BIC measures permit us to compare several models which don't have to be nested. They do so by adding a penalty to the likelihood for a model's complexity via the amount of parameters it has.

We wish to maximize the likelihood of our observations. As observed for the LRT, maximizing the likelihood is equivalent to minimizing the loglikelihood or a function thereof. Namely, $-2 \times \ell(\theta)$ (a.k.a. the deviance, see <u>Graphical Tests</u>). In both cases, we add a penalty to the measure we wish to minimize.

■ Akaike Information Criteria (AIC)

The AIC penalizes models which have more parameters by adding twice the number of estimated parameters p in the model to twice the negative log-likelihood: $AIC = -2\ell(\theta) + 2p$.

We choose the model with the smallest AIC.

Context

The disadvantage of the AIC lies in that for 2 nested models the probability of choosing the simpler model knowing it's the true model does not tend towards 1 when the number of observations increases towards infinity. We thus consider it an *inconsistent* measure.

In comparison, the BIC **is** a *consistent* measure given its parameters penalty is a function of the number of observations.

That being said, in both cases, the probability of rejecting the simpler model while the true model is somewhere in between tends towards 1.

■ Bayesian Information Criteria (BIC)

The BIC penalizes more severely models which have more parameters given its penalty is a function of the number of observations n: $BIC = -2\ell(\theta) + \ln(n)p$.

To better understand the difference between the AIC and BIC penalty, we can use log rules to rewrite the measures:

$$AIC = -2 \ln |\mathcal{L}(\theta)| + 2p$$

$$= -2 \ln |\mathcal{L}(\theta)| + \ln \left(e^{2p}\right)$$

$$= -\left[\ln |\mathcal{L}(\theta)^{2}| - \ln \left| (e^{p})^{2} \right| \right]$$

$$= -\ln \left| \frac{\mathcal{L}(\theta)^{2}}{\left(e^{p}\right)^{2}} \right|$$

$$BIC = -2 \ln |\mathcal{L}(\theta)| + \ln |n|p$$

$$= -\left[\ln \left| \mathcal{L}(\theta)^{2} \right| - \ln |n^{p}| \right]$$

$$= -\ln \left| \frac{\mathcal{L}(\theta)^{2}}{n^{p}} \right|$$

Context

There's no agreement on which is better for LMM and the text tends to build models piecewise, testing between steps with LRT. So, we probably won't use them much.

Fundamentally, the AIC tries to to find the model that best describes the data under the belief that there is no "correct" model. In contrast, the BIC tries to find the "correct" model under the belief that such a model exists.

Intuitively, we may think we'd prefer the AIC given that it's typically unrealistic to believe there exists a "correct" model. However, some feel the BIC often gives better results. *However*, part C on *Graphical Tests* has other information criterion that are more complicated but arguably better.

Notes:

- > REML criterion at convergence is the deviance $(-2\ell(\theta))$.
 - Likelihoods are typically <1 given they're probability densities.
 - Thus, loglikelihoods are typically negative.
 - Thus a positive output suggests they already multiplied by -2.

Checking model assumptions

≡ Likelihood distance / displacement

The change in Maximum log-Likelihood for all the data with the parameter of interest ψ estimated with all vs reduced data is

$$LD_{(u)} = 2\left(\ell(\hat{\boldsymbol{\psi}}) - \ell(\hat{\boldsymbol{\psi}}_{(u)})\right)$$

≡ Restricted likelihood distance / displacement

The change in Restricted Maximum log-Likelihood for all the data with the parameter of interest ψ estimated with all vs reduced data is

$$LD_{(u)} = 2\left(\ell_R(\hat{\boldsymbol{\psi}}) - \ell_R(\hat{\boldsymbol{\psi}}_{(u)})\right)$$

\equiv Cook's D

For parameter of interest $\pmb{\beta}$ ou $\pmb{\theta}$, the scaled change in the entire estimated $\pmb{\beta}$ ou $\pmb{\theta}$ vector is

$$D(\pmb{\beta}) = \left| \left(\hat{\pmb{\beta}} - \hat{\pmb{\beta}}_{(u)} \right)' \widehat{\operatorname{Var}} \left[\hat{\pmb{\beta}} \right]^{-1} \left(\hat{\pmb{\beta}} - \hat{\pmb{\beta}}_{(u)} \right) / rank(\pmb{X}) \right|.$$

$$D(\boldsymbol{\theta}) = \left[\left(\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{(u)} \right)' \widehat{\operatorname{Var}} \left[\hat{\boldsymbol{\theta}} \right]^{-1} \left(\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{(u)} \right) \right].$$

■ Multivariate DFFITS Statistic

For parameter of interest $\boldsymbol{\beta}$ ou $\boldsymbol{\theta}$ the scaled change in the entire estimated $\boldsymbol{\beta}$ ou $\boldsymbol{\theta}$ vector, using externalized estimates of $\operatorname{Var}(\hat{\boldsymbol{\beta}})$, is

$$MDFFITS(\boldsymbol{\beta}) = \left[\left(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(u)}\right)'\widehat{\operatorname{Var}}\left[\hat{\boldsymbol{\beta}}_{(u)}\right]^{-1}\left(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(u)}\right)/rank(\boldsymbol{X})\right].$$

$$MDFFITS(\boldsymbol{\theta}) = \left[\left(\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{(u)} \right)' \widehat{\operatorname{Var}} \left[\hat{\boldsymbol{\theta}}_{(u)} \right]^{-1} \left(\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{(u)} \right) \right].$$

Trace of covariance matrix

For parameter of interest $\boldsymbol{\beta}$ ou $\boldsymbol{\theta}$, the change in precision of estimated $\boldsymbol{\beta}$ ou $\boldsymbol{\theta}$ vector, based on trace of $\operatorname{Var}(\hat{\boldsymbol{\beta}})$, is

$$COVTRACE(\boldsymbol{\beta}) = \left| \left| trace(\widehat{\operatorname{Var}} \left[\hat{\boldsymbol{\beta}} - \widehat{\operatorname{Var}} \left[\hat{\boldsymbol{\beta}}_{(u)} \right]^{-1} \right]^{-1}) - rank(\boldsymbol{X}) \right| .$$

$$COVTRACE(\boldsymbol{\theta}) = \left| trace(\widehat{\operatorname{Var}} \left[\hat{\boldsymbol{\theta}} - \widehat{\operatorname{Var}} \left[\hat{\boldsymbol{\theta}}_{(u)} \right]^{-1} \right]^{-1}) - q \right|$$

≡ Covariance ratio

For parameter of interest β ou θ , the change in precision of estimated β ou θ vector, based on determinant of $Var(\hat{\beta})$, is

$$COVRATIO(oldsymbol{eta}) = egin{array}{c} \widehat{\operatorname{Var}}[\hat{oldsymbol{eta}}_{(u)}] \\ \widehat{\operatorname{Var}}[\hat{oldsymbol{eta}}] \end{array}$$

$$COVRATIO(\boldsymbol{\theta}) = \begin{bmatrix} \widehat{\operatorname{Var}}[\hat{\boldsymbol{\theta}}_{(u)}] \\ \widehat{\operatorname{Var}}[\hat{\boldsymbol{\theta}}] \end{bmatrix}$$

≡ Sum of squared PRESS residuals

The sum of PRESS residuals calculated by deleting observations in U is $PRESS_{(u)} = \sum_{i \in u} (y_i - x_i' \hat{\beta}_{(u)})$.

Graphical Tests

Not heavily tested. Case study will have some graphs and there will be some questions about case study which may need graphs interpretation

marginal residual residual leftover plugging in estimated fixed effects rarely used conditional (textbook) / response (R) / raw (typical) residuals residual from estimate of everything In LMM, variance of residual ε_{ij} can vary based on other factors So, still not residual we want

standardized / normalized residuals conditional residual / estimated SD for that residual almost always prefer standardized residual important:

- > Use residual plots for normality testing
- > raw data plots are useless; ignore them.
- > Standardized residuals adjust the data so we can tell if a residual is an outlier because it's from a high variance group or because it's really an outlier.

implied marginal model is LMM w/o random effects but with same variance structure (var Yij same for both) marginal model is with just same variance for everything

 \mathbf{D}

Bayesian Analysis and Markov Chain Monte Carlo

 \mathbf{E}

Statistical Learning

K-Nearest Neighbors

Decision Trees

Principal Components Analysis (PCA)

Clustering