



Bayesian Model Checking in GLSM for Count Data

Liang Jing and Victor De Oliveira

liang.jing@utsa.edu, victor.deoliveira@utsa.edu



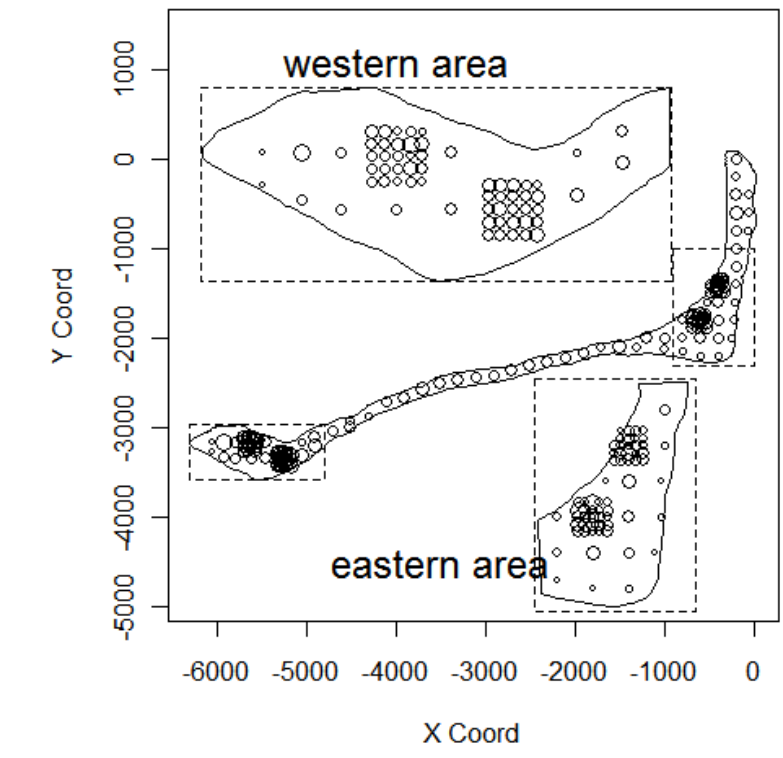
ABSTRACT

Model checking and selection in hierarchical models remain difficult problems due to the unobservable latent process. Here Bayesian model checking methods are introduced and compared with a new model checking method based on transformed residuals. Our simulation study reveals that Bayesian model checking methods fails on GLSM due to difficulty of choosing appropriate diagnostic statistics and conservatism of posterior predictive p-value, while the new model checking method provides more promising results.

GEOSTATISTICAL DATA

Each observation consists of **two attributes**:

- sampling location s_i , where the response is observed;
- observed value Y_i , count of certain event.



Ronglap data set were collected from Rongelap Island, 2500 miles south-west of Hawaii. U.S. nuclear weapons testing programme generated heavy fallout over the island in the 1950s and Rongelap island has been uninhabited since 1985. There are 157 observations:

- y_i is the photon emission count
- s_i identifies spatial location
- l_i is the time (in seconds) over which y_i was accumulated

HIERARCHICAL MODEL

Generalized linear spatial models were first proposed by Diggle et al. (1998) and widely used since then,

$$Y_i|\theta_i \stackrel{\text{id}}{\sim} \text{Pois}(\cdot|e^{\theta_i}), \quad i = 1, \dots, n;$$

$$\theta|\eta \sim \text{MVN}(D\beta, \Sigma), \quad \Sigma_{ij} = \sigma^2\rho(u_{ij}; \phi, \kappa)$$

$$\eta = (\beta, \sigma^2, \phi, \kappa)$$

- $\theta(s)$ is stationary Gaussian process;
- $D\beta$ is the mean structure where D is covariate matrix (usually related to locations) and β is coefficient vector;
- $\rho(u_{ij}; \phi, \kappa)$ is a correlation function, determined by Euclidean distance $u_{ij} = \|s_i - s_j\|$ between the i^{th} and j^{th} locations s_i, s_j . For example, matern family:

$$\rho(u) = [2^{\kappa-1}\Gamma(\kappa)]^{-1}(u/\phi)^\kappa K_\kappa(u/\phi)$$

where $K_\kappa(\cdot)$ denotes a modified Bessel function of order κ .

R PACKAGE DEVELOPMENT

► **Features** (under development)

- performs posterior sampling for parameter estimation, prediction, and model checking in hierarchical models with correlated latent variables;
- C++ programs are seamlessly embedded to handle heavy computational tasks of Markov chain generation and large matrices computation;
- parallel computing techniques are implemented to further speed up estimation and prediction;
- results are displayed by a combination of numerical and graphical summaries.

ROBUST MCMC ALGORITHM

Common Hastings-within-Gibbs algorithms fail to generate well-behaved posterior samples because of the large number of latent variables that are highly correlated and influenced by response variables. We implemented robust MCMC algorithm proposed by Christensen et al. (2006)

► **Data-based parameterization**

after parameterization, the components of $\tilde{\theta}, \tilde{\beta}, \tilde{\eta}$ are approximately uncorrelated, and have zero means and unit variances

$$\theta \rightarrow \tilde{\theta}(\theta; \beta, \eta, Y)$$

$$\beta \rightarrow \tilde{\beta}(\beta; \eta, Y)$$

$$\eta \rightarrow \tilde{\eta}(\eta; Y)$$

► **Langevin-Hastings algorithm**

the gradient information of the target density is used into the proposal density to improve convergence.

► **Group Updating with Gibbs Sampler**

- for N iterations do
- 1 Update $\tilde{\theta} \sim p(\tilde{\theta}|\tilde{\beta}, \tilde{\eta}, Y)$
 - 2 Update $\tilde{\beta} \sim p(\tilde{\beta}|\tilde{\theta}, \tilde{\eta}, Y)$
 - 3 Update $\tilde{\eta} \sim p(\tilde{\eta}|\tilde{\theta}, \tilde{\beta}, Y)$

► **Prior**

Christensen (2002) suggested flat prior and provided conditions under which the posterior is proper.

BAYESIAN MODEL CHECKING

Bayarri and Castellanos (2007) stated that **existing Bayesian model checking methods** can be seen to correspond to particular choices of the following three components:

- a **diagnostic statistic** $T(\mathbf{y})$ that summarizes a feature of the data
- a **specified reference distribution** $h(t)$ that represents an “predictive” distribution for T
- a **way to measure conflict** between $t^{\text{obs}} = T(\mathbf{y}^{\text{obs}})$ and the reference distribution

► **Two ways to measure conflict**

$$p\text{-value} = P^{h(\cdot)}(T(\mathbf{Y}) \geq T(\mathbf{y}^{\text{obs}})), \quad RPS = \frac{h(t^{\text{obs}})}{\sup_t \{h(t)\}}$$

small values indicate incompatibility.

► **The choices of reference distribution**

$$h(t) \leftarrow h(\mathbf{y}) = \int p(\mathbf{y}|\theta)\pi^*(\theta) d\theta$$

thus, choosing $h(t)$ amounts to choosing $\pi^*(\theta)$.

- Empirical Bayes distribution: $\pi_{\text{prior}}^{EB}(\theta) = \pi(\theta|\eta = \hat{\eta})$
 - Posterior predictive distribution: $\pi(\theta, \eta|\mathbf{y}^{\text{obs}}) \propto p(\mathbf{y}^{\text{obs}}|\theta)\pi(\theta, \eta)$
 - Partial posterior predictive distribution
- $$\pi_{\text{ppp}}(\theta) = \pi(\theta|\mathbf{y}^{\text{obs}} \setminus t^{\text{obs}}) \propto \frac{p(\mathbf{y}^{\text{obs}}|\theta)\pi(\theta)}{p(t^{\text{obs}}|\theta)}$$

► **Results**

Simulated Data vs. Assumed Model				
	T_1	T_2	T_3	T_4
“explt” vs. “explt”	0.49/0.97	0.53/0.99	0.47/0.99	0.12/0.58
“explt” vs. “expnt”	0.48/0.99	0.52/1.00	0.45/0.97	0.12/0.63
“expQt” vs. “expnt”	0.52/1.00	0.54/1.00	0.54/0.99	0.24/0.70
“expSin” vs. “expnt”	0.51/1.00	0.59/0.95	0.49/0.99	0.73/0.77

TRANSFORMED RESIDUAL CHECKING

► **Transformed residuals** are defined as

$$e_i = \Phi^{-1}(F_i(Y_i))$$

where the CDF $F_i(\cdot)$ for GLSM comes from

$$f(y_i|\eta) \leftarrow f(\mathbf{y}|\eta) = \int \prod_i p(y_i|\theta_i)\pi(\theta|\eta) d\theta.$$

► **Property**

When

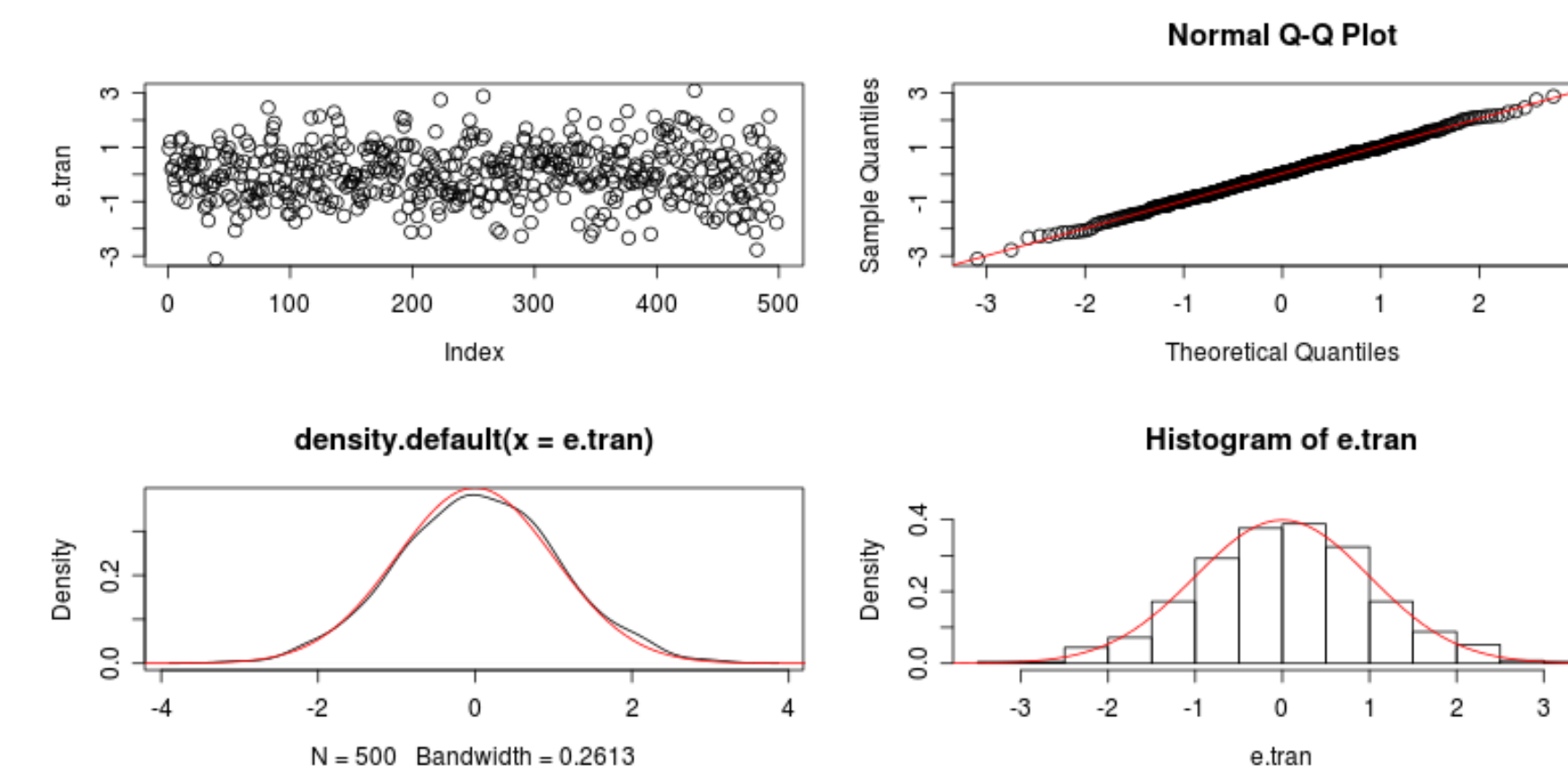
- (1) the model is correctly specified and
- (2) η is the true parameter,

transformed residuals will follow standard normal distribution.

Thus, by examining the normality of transformed residuals the information of goodness of fit can be obtained.

► **Comparing Distributions**

- Graphical methods: QQ-plot, histogram, relative density plot



- Numerical methods: we propose to use Hellinger distance. **Hellinger distance** is defined as

$$d_H(P, Q) = \frac{1}{2} \int |\sqrt{dP} - \sqrt{dQ}|^2$$

where dP and dQ are two PDFs to be compared.

► **The Procedure of Calculation**

- calculate the Hellinger distance between empirical distribution of transformed residuals and standard normal.

$$(e^{\text{obs}}, N(0, 1)) \rightarrow d^{\text{obs}}$$

- obtain the “baseline distribution”, which is the distribution of Hellinger distances computed between theoretical standard normal distribution and empirical distributions of standard normal samples.

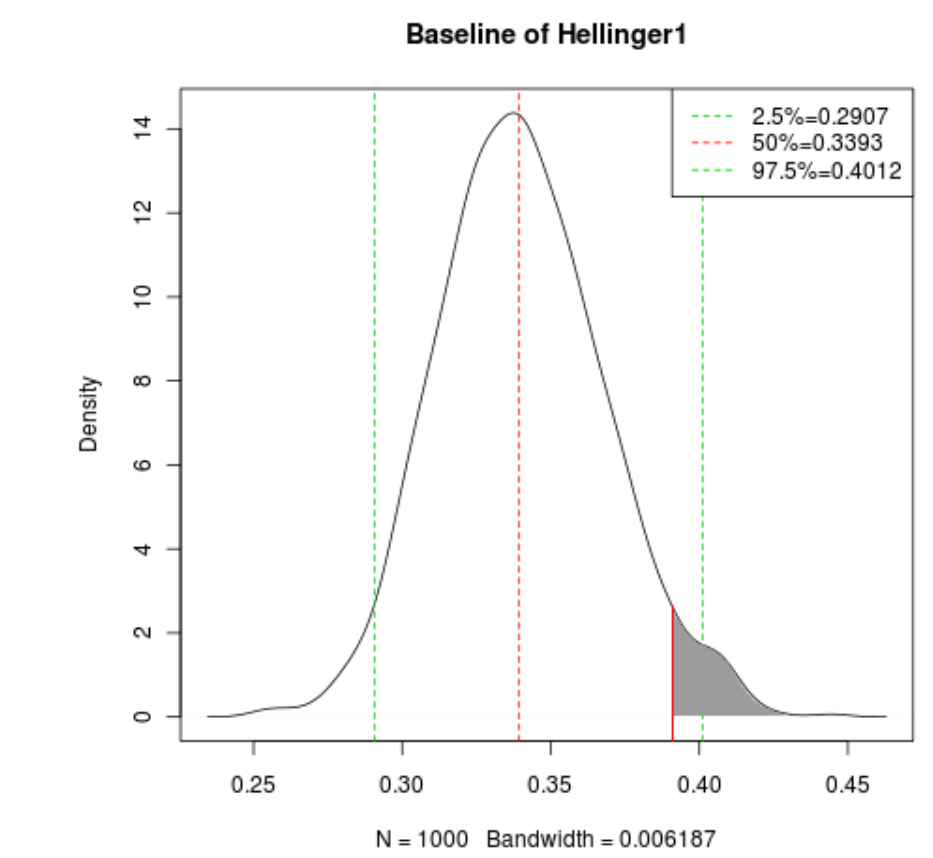
$$x_i^{(k)} \sim N(0, 1) \quad i = 1, \dots, n$$

$$(x^{(k)}, N(0, 1)) \rightarrow d^{(k)} \quad k = 1, \dots, N$$

- calculate “one-side p-value”, $p = P(d \geq d^{\text{obs}})$

REFERENCES

- [1] Diggle, P. J., Tawn, J. A. and Moyeed, R. A. (1998). Model-based geostatistics. *J. Roy. Statist. Soc. Ser. C*, 47, 299-350.
- [2] Christensen, O. F., Roberts, G. O. and Sköld, M. (2006). Robust Markov chain Monte Carlo methods for spatial generalized linear mixed models. *J. of Computational and Graphical Statistics*, Vol. 15, No. 1, 1-17.
- [3] Bayarri, M. J. and Castellanos, M. E. (2007). Bayesian checking of the second levels of hierarchical models. *Statistical Science*, Vol.22 No. 3. 322-343.



Small p-value rejects the normality of transformed residuals which means the observed data is NOT compatible with the assumed model.

► **Results**

Mis-fit			Correct-fit		
	p-value	r.rate		p-value	r.rate
“explt” vs. “expnt”	0.02	0.92	“explt” vs. “explt”	0.26	0.01
“expQt” vs. “expnt”	0.00	1.00	“expQt” vs. “expQt”	0.28	0.01
“expSin” vs. “expnt”	0.00	1.00	“expSin” vs. “expSin”	0.16	0.08

Ronglap data			Weed data		
	“expnt”	“explt”		“expnt”	“explt”
$\hat{\beta}$	1.83	(1.87, -0.49, 0.31)	$\hat{\beta}$	4.07	(4.35, 0.10, -1.51)
$\hat{\sigma}$	0.55	0.60	$\hat{\sigma}$	1.09	1.22
$\hat{\phi}$	0.02	0.03	$\hat{\phi}$	0.17	0.19
p-value	0.027	0.000	p-value	0.315	0.004
r.rate	0.90	1.00	r.rate	0.005	1.00

► **Conclusions**

- Bayesian model checking can detect model failure only when the diagnostic statistic is well chosen which is usually very difficult without deep understanding of data.

- Transformed residual checking has much better performance. In the case the mean structure is misspecified, model failures are successfully detected.
- Transformed residual checking suggests GLSM without any trend in mean structure for “Weed” data and rejects GLSM with either linear trend or no trend for “Rongelap” data.

► **Improvement & Future Work**

Problem: however, when the normality of \mathbf{e} is rejected, we cannot determine the rejection is due to condition (1) or (2), or both.

Goal: find a way to determine

- if condition (1) is satisfied?
- if yes, if condition (2) is satisfied?