# When Growth Mixture Models Break:

Identifiability Failures and Misleading Model Evaluation

**Presenter:** Doria Xiao

**Affiliations:** 

Ph.D., University of California, Berkeley;

Postdoctoral Fellow, Stanford University

**Conference Details:** 

**BEAR Seminar** 

September 9, 2025



Berkeley Evaluation and Assessment Research Center

#### **Outline**

- 1. Motivation
  - Why model selection is critical in GMMs
  - The stakes when criteria fail
- 2. Frequentist Failures
  - Local maxima and misleading AIC
- 3. Bayesian Failures
  - The likelihood choice issue
  - Negative DIC penalties reveal hidden nonidentifiability
- 4. Consequences of Pathologies
  - Minuscule-class behavior
  - Twinlike-class behavior
- 5. Role of Priors
  - How vague priors exacerbate problems
  - Informative priors as remedies
- 6. Practical Recommendations
  - Using criteria as diagnostics
  - Rethinking evaluation workflow
- 7. Conclusion
  - o Failures are signals, not just noise

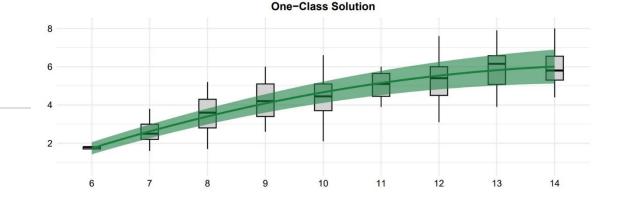
#### What Are GMMs

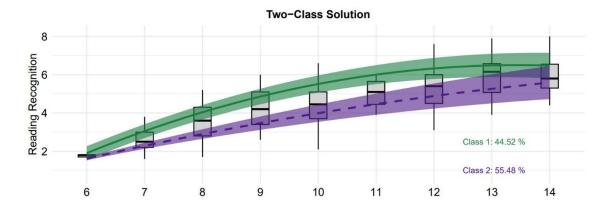
1-class: qualitatively homogeneous growth

• 2-class: early vs. late developers

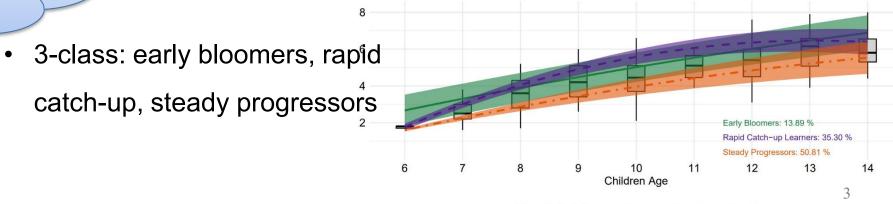
Same data, different interpretations Which model should we trust?

catch-up, steady progressors





Three-Class Solution





#### **How Each Information Criterion Works**

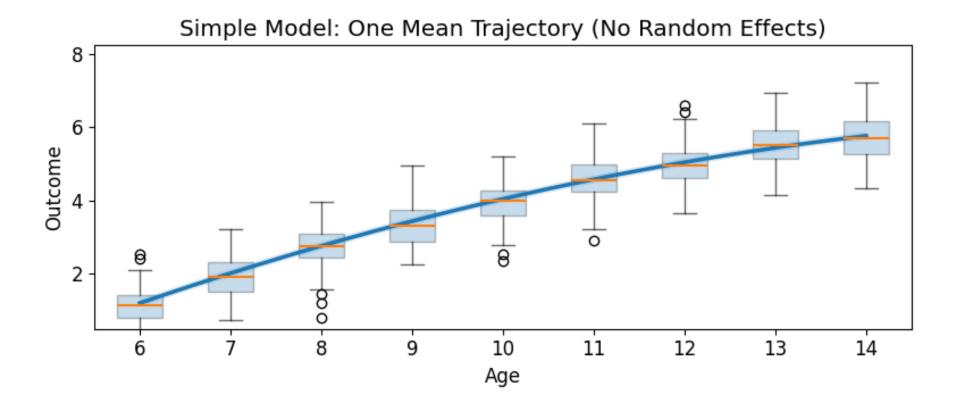
Criterion		Estimator Form	Penalty/complexity p		
	Fit term + Penalty term		Optimism due to using data twice		
<b>AIC</b> (Akaike, 1973)	<sup>1</sup> Stata, Mplus, R/Ime4	$-2\log f(y\mid \hat{\theta}(y)) + 2p$	Number of parameters $p$		
DIC (Spiegelhalter et al., 2002)	<sup>2</sup> Open BUGS, JAGS	$-2\log f(y\mid \tilde{\theta}(y)) + 2p_D$	Mean deviance minus plug-in deviance $p_D = \overline{D} - D(\overline{\theta})$		
<sup>3</sup> Stan <b>WAIC</b> (Watanabe, 2010)	-2	Plug-in deviance $D(\bar{\theta})$ $\sum_{i} \log p_{\text{post}}(y_i \mid y) + 2p_{\text{WAIC}}$	Posterior variance of log-likelihood contrib. $p_{\text{WAIC}} = \sum_{i=1}^{N} \text{Var}_{\text{post}} \left[ \log f(y_i \mid \theta) \right]$ Posterior		
LOO-CV (Vehtari et al., 2016)	-2	$\sum_{i} \log p_{\mathrm{post}}\left(\left.y_{i} \mid y_{-i}\right.\right)$ , via PSIS	NA*, leave-one-out reweighting		

Fit term ≈ log-likelihood or log predictive density

Penalty term ≈ effective number of parameters (adjusts for reuse of data)

PSIS: Pareto-Smoothed Importance Sampling (Vehtari et al., 2017)

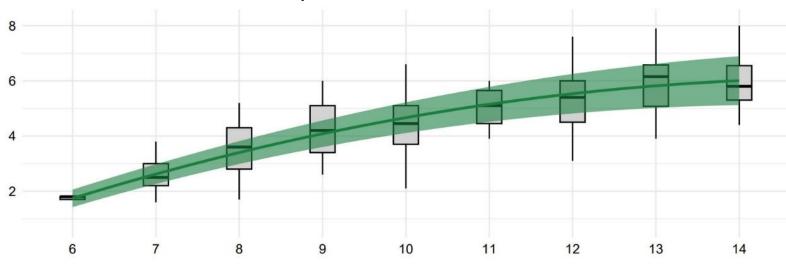
## From Simple to Hierarchical to GMM



$$f(\mathbf{y} \mid \boldsymbol{\theta}) = \prod_{i=1}^{N} f(y_i \mid \boldsymbol{\theta})$$

## From Simple to Hierarchical to GMM

Hierarchical Model: One Population Curve with Individual Variation



Conditional Likelihood: 
$$f_c(\mathbf{y} \mid \boldsymbol{\beta}, \mathbf{r}) = \prod_{j=1}^J f(\mathbf{y}_j \mid \boldsymbol{\beta}, \mathbf{r}_j) = \prod_{j=1}^J \prod_{i=1}^{n_j} f(\mathbf{y}_{ij} \mid \boldsymbol{\beta}, \mathbf{r}_j)$$

 $\succ$ Conditions on random effects  $r_i$  -> Predict at individual level\*

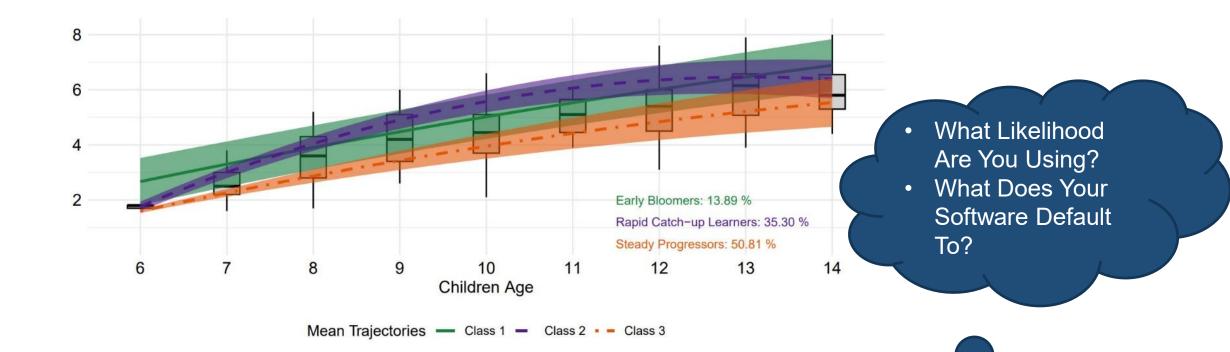
Marginal Likelihood: 
$$f_m(y \mid \beta, \Sigma) = \prod_{j=1}^J f(y_j \mid \beta, \Sigma) = \prod_{j=1}^J \int f(y_j \mid \beta, r_j) p(r_j \mid \Sigma) dr_j$$

➤Integrates over random effects -> Predict at population level\*

- Vaida, F., & Blanchard, S. (2005, 06). Conditional Akaike information for mixed-effects models. Biometrika
- Spiegelhalter, D. J., Best, N. G., Carlin, B. P., &van der Linde, A. (2002). Bayesian measures of model complexity and fit. Journal of the Royal Statistical Society, Series B,
- Merkle, E., Furr, D., & Rabe-Hesketh, S. (2019). Bayesian comparison of latent variable models: Conditional versus marginal likelihoods.

Psychometrika

## From Simple to Hierarchical to GMM



$$y_{ij}|_{w_j} = k, r_{1j}, r_{2j} \sim N \left(\mu_{ij}^{(k)}, \sigma_e\right);$$

Categorical latent variable (class label)  $w_i$ 

Multinomial distribution with probability parameters  $\{\lambda^{(1)}, ..., \lambda^{(K)}\}$ 

#### **Continuous** latent variables

•  $(r_{1j}, r_{2j})' | w_j = k \sim N(\mathbf{0}, \mathbf{\Sigma}^{(K)})$  with class-specific covariance matrix  $\mathbf{\Sigma}^{(K)}$ 

## Bayesian: The Right Likelihood for GMM for the Right Purpose

Likelihood Type	What Is Integrated or Conditioned?	Prediction Target	Valid for	Common Software
Marginal	Integrates over both latent classes and random effects	Predict outcomes in new clusters	Class enumeration	Stan (which marginalizes over discrete parameters)
Conditional	Conditions on class memberships and random effects	Predict outcomes for in-sample clusters	Model comparison for in-sample clusters	OpenBUGS, JAGS
Hybrid	Integrates over classes but conditions on random effects	Ambiguous: in- sample clusters, but use prior for class prob.	Theoretically incoherent	Often occurs by default in Stan when conditioning on random effects

#### **Bottom line:**

Only the **marginal likelihood** aligns with the population-level goal of **class enumeration**. Conditional usable for some model comparisons, but **hybrid not valid for model comparison**.

## Bayesian: Why DIC Breaks in GMMs, How to Fix It





#### **Problem: Traditional DIC**

DIC = 
$$\overline{D}$$
 +  $p_D$ , where  $p_D = \overline{D} - D(\overline{\theta}) \rightarrow$   
DIC =  $\overline{D}$  +  $\overline{D}$  -  $D(\overline{\theta})$ =  $\overline{D}$  +  $p_D$ 

#### In GMMs:

- Skewed / multimodal posteriors, label switching or degenerate nonidentifiability (Xiao, Rabe-Hesketh, & Skrondal, 2015), leads to poor estimate  $\bar{\theta}$
- Plug-in deviance too large
- $p_D$  can be **negative or unstable** (Spiegelhalter et al., 2002; Gelman et al., 2013)

#### **More Stable Alternatives**

• DIC $_{pV2}$  (Gelman , Hwang, and Vehtari, 2014):

$$\mathrm{DIC}_{pV2} = D(\bar{\theta}) + p_V$$

- Variance-based penalty,
- BUT retains plug-in deviance
- DIC<sub>pV</sub> (we proposed):

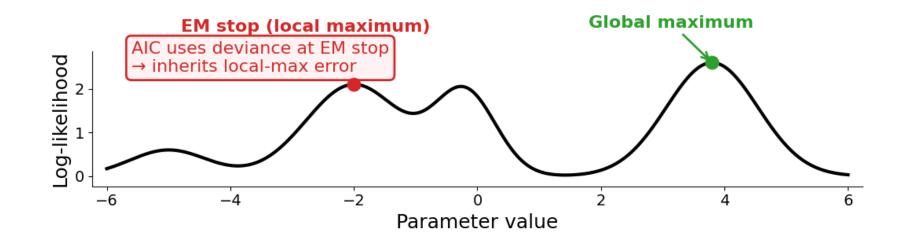
$$\mathrm{DIC}_{pV} = \overline{D} + p_V$$

- No plug-in deviance
- Fully posterior-based

## Why AIC Breaks in GMMs, How to Fix It

- Plug-in likelihood depends on EM
- EM can stop at local maxima
- Researchers assume "convergence = global"
- Fix: more iterations, more random starts
- Still no guarantee → trial-and-error

#### Local Maxima → AIC Computed at Wrong Solution



#### **Simulation Conditions & Structure**

### **Purpose**

When do model evaluation tools succeed or fail for enumeration?

### **Design Factors**

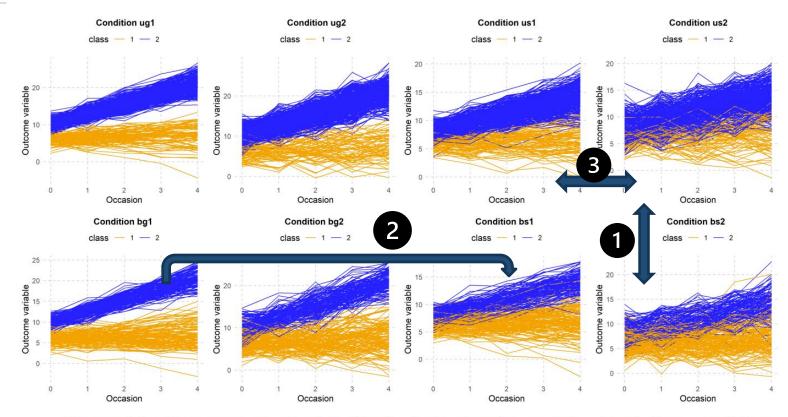


Figure 3.1: Comparison of subject trajectories by class across simulation conditions

Factor	Levels	
Class Probabilities and Level-2 Sample Sizes J	Balanced ( $\lambda$ = 0.5) with J = 250 vs. Unbalanced ( $\lambda$ = 0.2 / 0.8) with J = 400	
Class Separation	Strong vs. Weak slope/intercept differences	
Residual Variability	Low ( $\sigma_e$ = 1) vs. High ( $\sigma_e$ = 2)	

## **Estimation Setup & Fit Criteria**

#### **Design Summary**

- Labels: bg1, ug2, us2, etc. (8 simulation conditions)
- **Replications**: 50 datasets per condition
- **Time Points**: 5 per subject
- Models Fitted: 1- to 4-class GMMs
- Total Fits:  $8 \times 50 \times 4 = 1,600$  per method

#### **Bayesian Estimation (CmdStan 2.30)**

- MCMC Specs: 4 chains × 1,000 post-warmup iterations
- Target: Marginal likelihood
- Information Criteria: DIC, DIC\_pV, DIC\_pV2, WAIC, LOO-CV

#### Frequentist Estimation (MLE via flexmix)

- Engine: R flexmix (Grün & Leisch, 2023)
- Information Criterion: AIC

## **Strengthening EM Estimation**

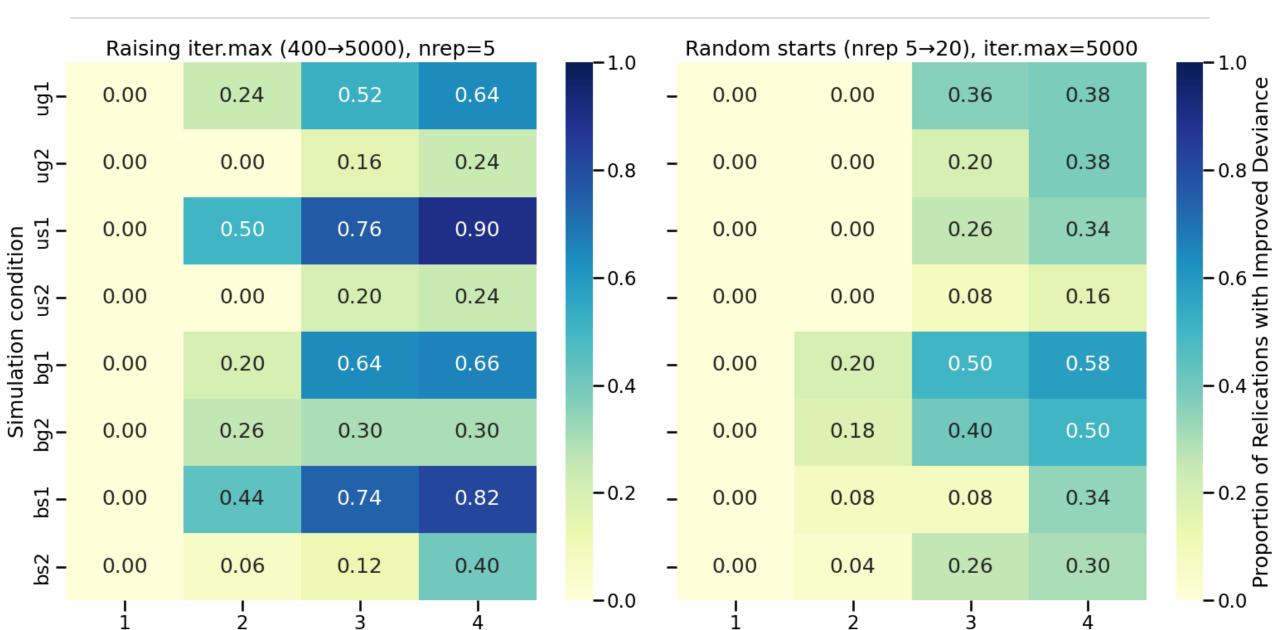
#### **Convergence settings (flexmix, Frequentist)**

- iter.max: 400 → 5,000
   Prevents premature stopping
- nrep: 5 → 20 (and beyond)
   More chances to escape local maxima

#### **Bottom note:**

These fixes improve robustness, but global maximum is not guaranteed.

## **EM** Convergence: Iteration Cap and Random Starts



## **EM Local Maxima, Fit Criteria Fail**

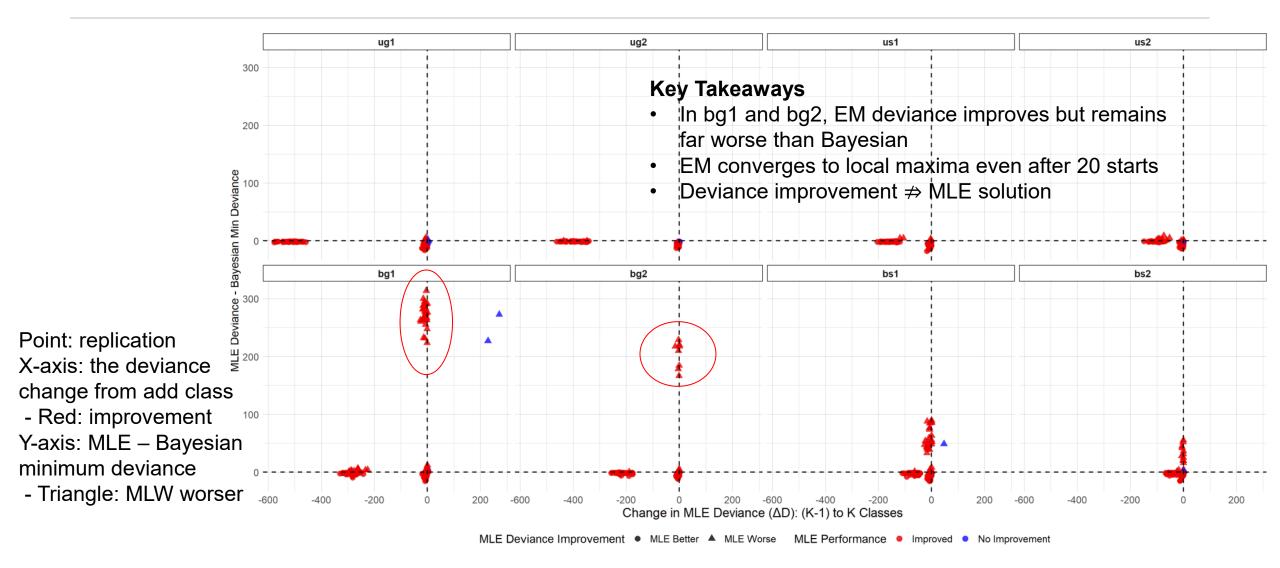
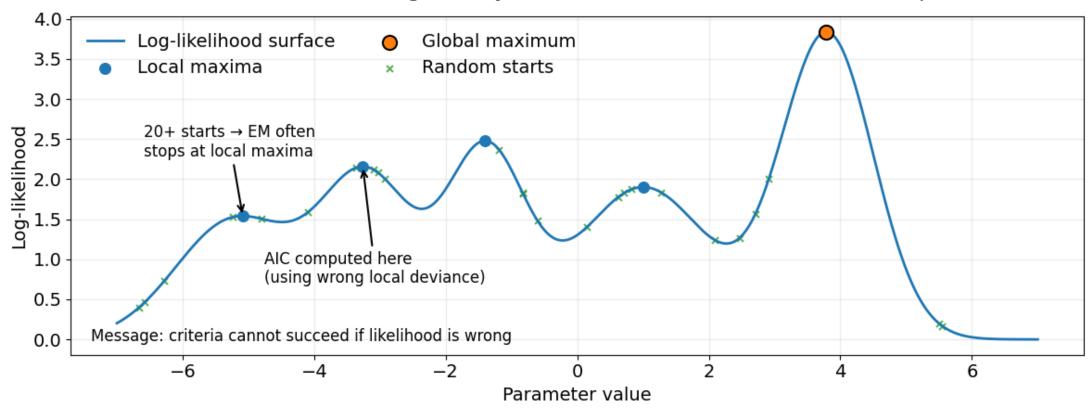


Figure 3.7: Comparison of MLE and Bayesian deviance across simulation conditions (based on nrep = 20).

## **Brute Force is Not Enough**

#### Brute Force Is Not Enough: Many Random Starts Still Miss the Global Optimum



## False confidence, wrong model.

## **Bayesian Estimation: Why DIC Variants Work When DIC Fails**

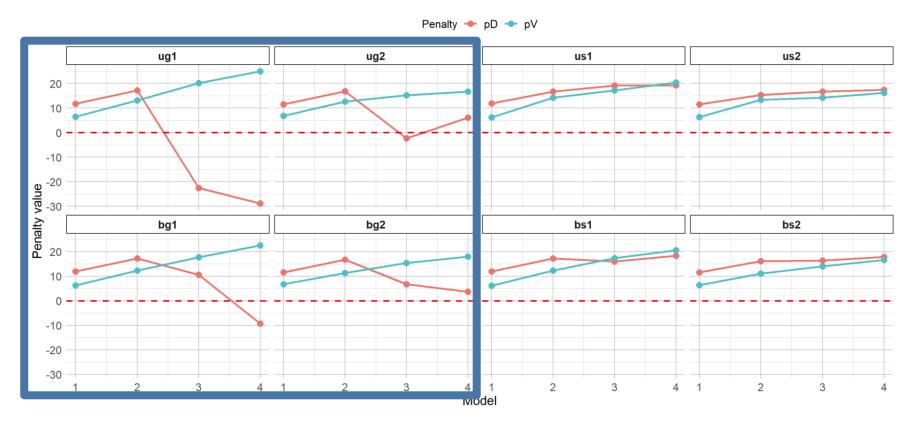


Figure 3.5: Comparison of average DIC penalty terms: Original  $(p_D)$  vs. variance-based  $(p_V)$  for each simulation Condition.

Traditional DIC penalty (pD) can be negative or unstable Variance-based penalty (pV) always positive and stable across all conditions

## Among DIC Variants, DIC\_pV Best Aligns with WAIC

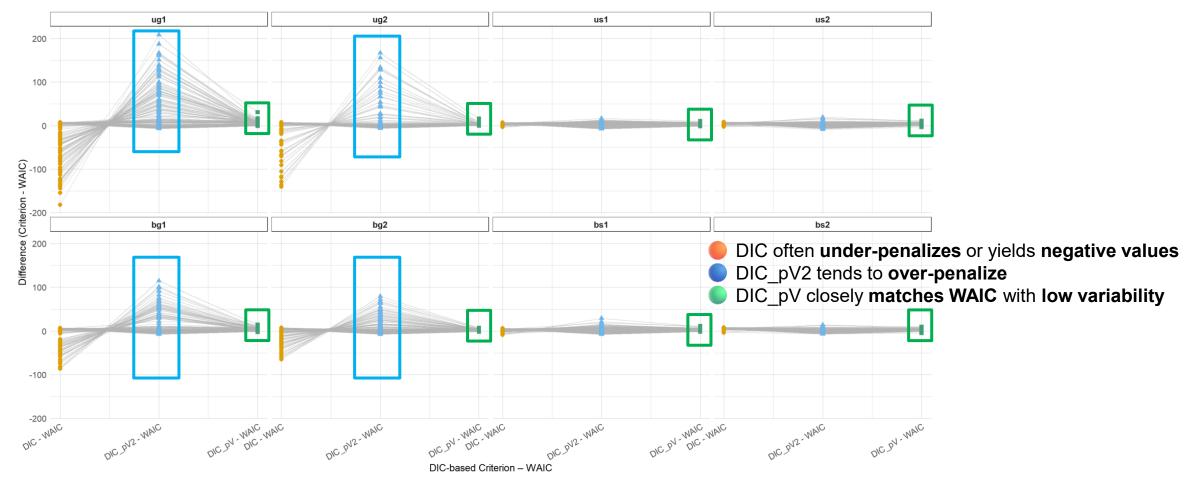
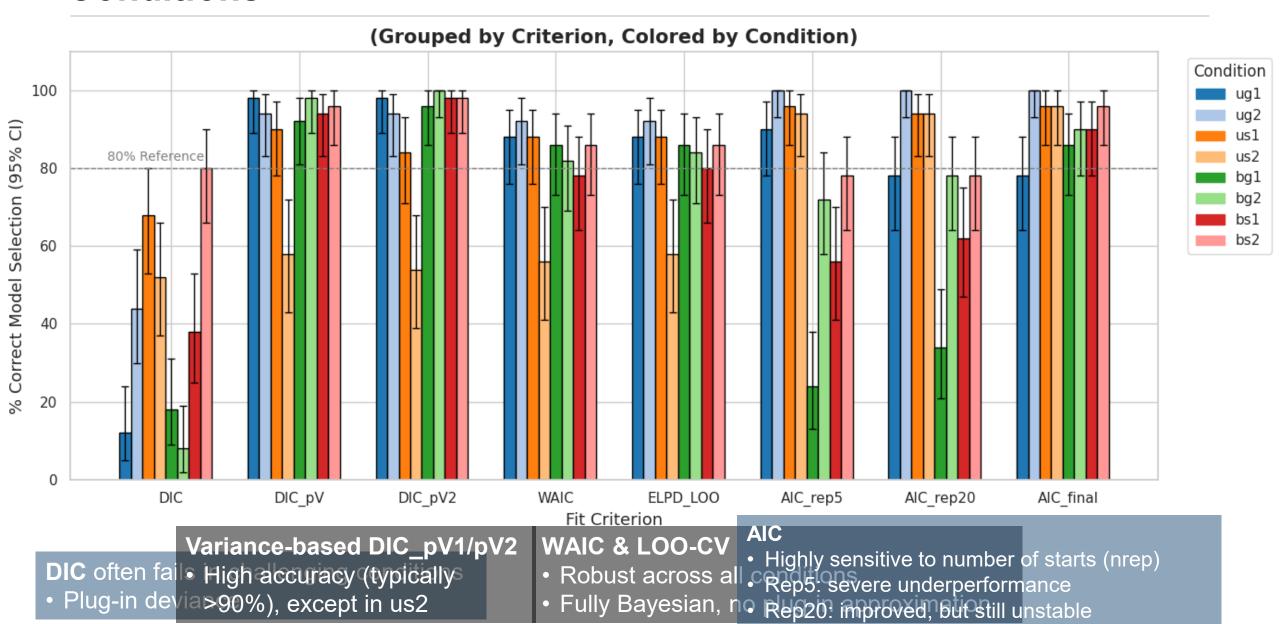
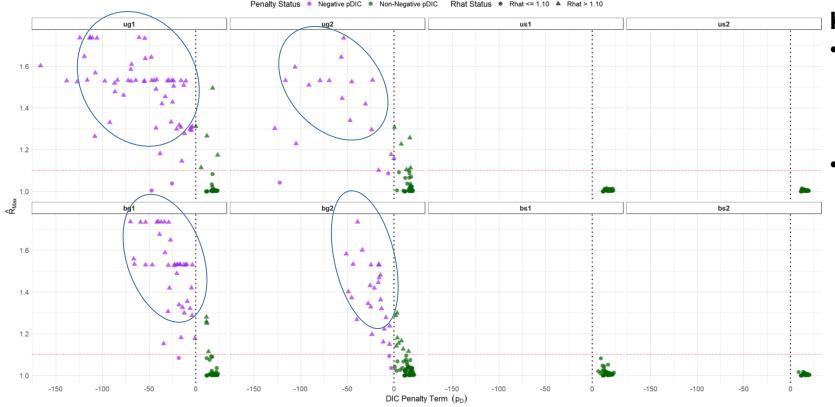


Figure 3.6: Replicate-level comparison of DIC criteria with WAIC for each simulation condition. Each replicate is represented as a path across DIC – WAIC (orange),  $DIC_{pV2}$  – WAIC (blue), and  $DIC_{pV}$  – WAIC (green).

# Information Criteria Performance Across Simulation Conditions



## **Beyond Selection: When DIC Penalty Reveals Failures**



#### **Key Takeaways**

- Negative  $p_D$  values strongly associated with poor convergence (purple triangles)
- Such failures often stems from **non-identifiability** (Xiao, X., Rabe-Hesketh, S., & Skrondal, A. (2025). Bayesian Identification and Estimation of Growth Mixture

Models. Psychometrika).

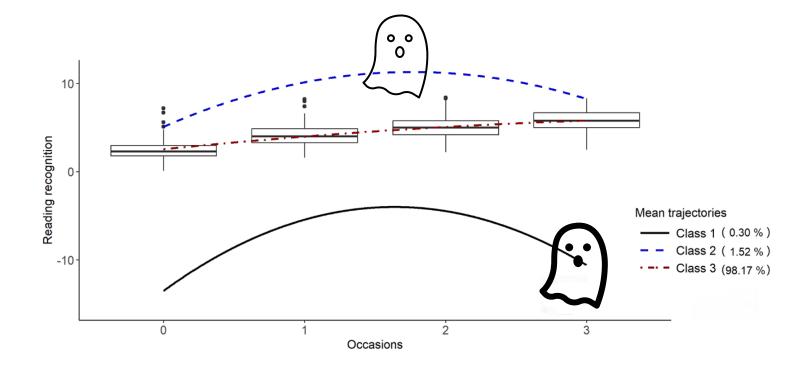
Figure 3.2: Relationship between DIC penalty term  $(p_D)$  and convergence diagnostic  $(\widehat{R}_{\text{Max}})$  across simulation conditions.

- Panels: simulation conditions
- X-axis: $p_D$
- Y-axis: the maximum convergence diagnostic  $\hat{R}_{\text{Max}}$
- Red dashed line at 1.10: the threshold for acceptable convergence
  - Black dotted line at 0: the boundary of  $p_D$
  - Point shape: convergence status (circle/triangle)
  - Point color: whether  $p_D$  is negative/non-negative (purple/green)



## Consequences of Minuscule-Class Behavior

- Class weight collapses toward zero ("ghost" class)
- Parameters drift from prior, not learned from data
- Subgroup trajectory becomes incompatible with observed data

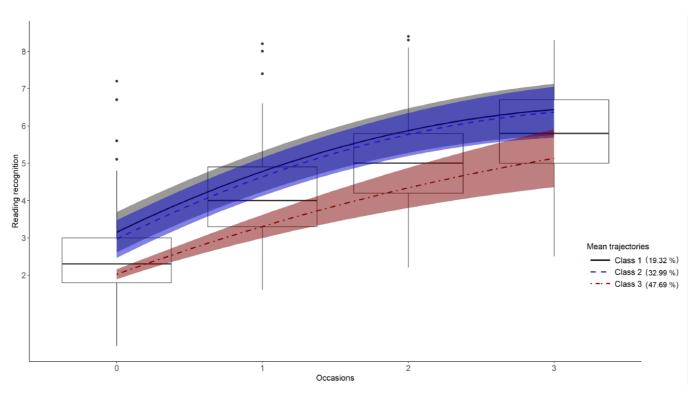


Two classes fit nonsensical curves (tiny class weights). Highlight "phantom" class trajectories with dotted lines.

## **Consequences of Twinlike-Class Behavior**

- One subgroup split into two nearly identical classes
- Creates illusion of distinct learner types
- Leads to overestimating heterogeneity





Same subgroup split in two  $\rightarrow$  inflates # of learner types.

## What Nonidentifiability Looks Like in GMMs



- Standard convergence (R) may not detect these
- Distinguishability Index (DI): near zero flags collapsed or indistinguishable classes

## Negative $p_D$ as a Diagnostic

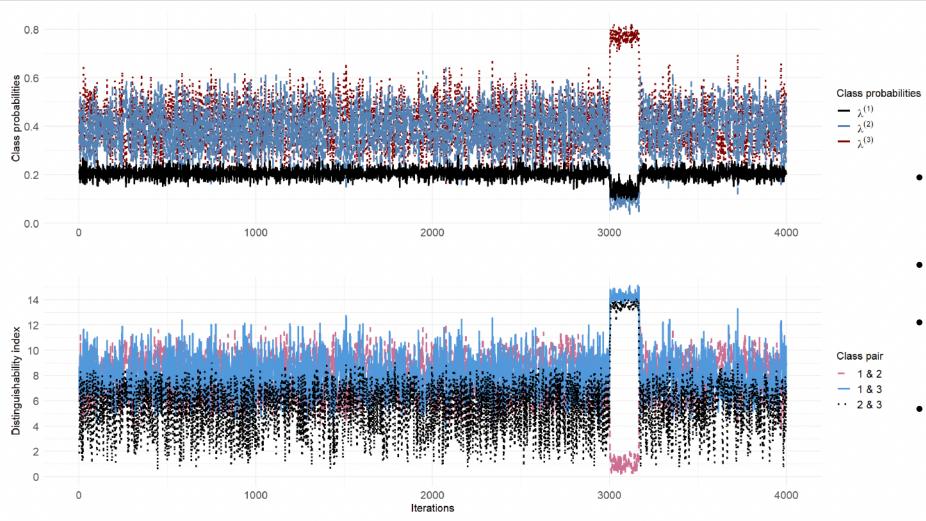


Figure 3.3: Traceplots of class probabilities and distinguishability index (DI) for Condition ug2, replicate 21, 3-class model.

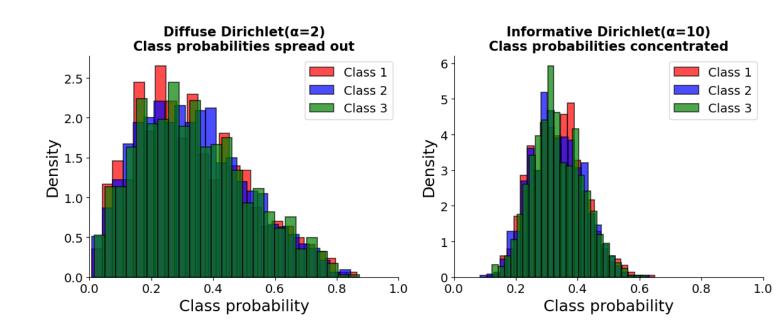
- Example: 3-class model with acceptable  $\hat{R}_{\text{Max}} = 1.04$
- But p<sub>D</sub> strongly negative -122
- Chain switches between degenerate solutions, undermining stability
- Distinguishability Index shows classes collapsing

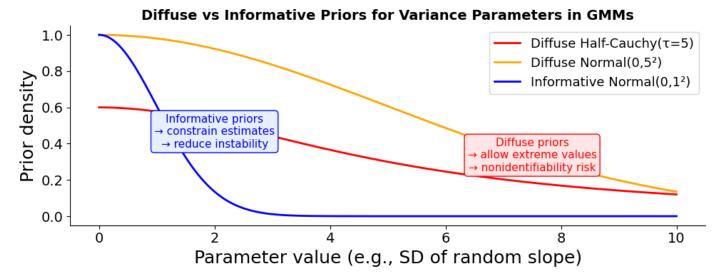
## From Weak Priors to Nonidentifiability

 Priors too weak → model can't separate classes.

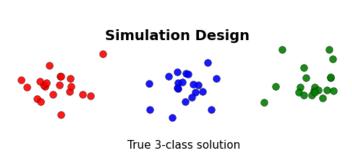
#### Example defaults

- Class probs: Dirichlet(α = 2)
- SDs: half-Cauchy(τ) or half-Normal(τ)



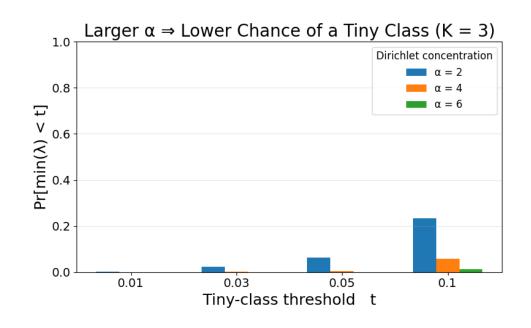


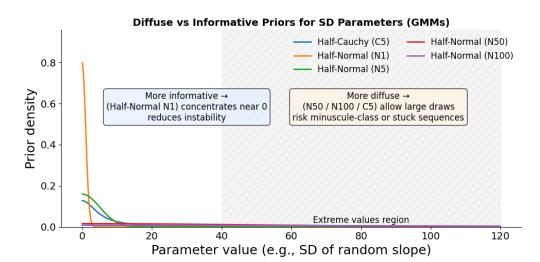
## Simulation Highlight: Testing Priors for Stability



200 chains × 1,000 iterations

Simulated responses for a real data using a well-behaved 3-class solution.





#### **Findings:**

- Vague priors → more stuck chains and minuscule-class behavior
- Informative priors → fewer failures, more stable estimation

#### **Practical Prior Choices**



Prefer half-Normal for SDs over half-Cauchy



Choose α to reflect plausible class balance

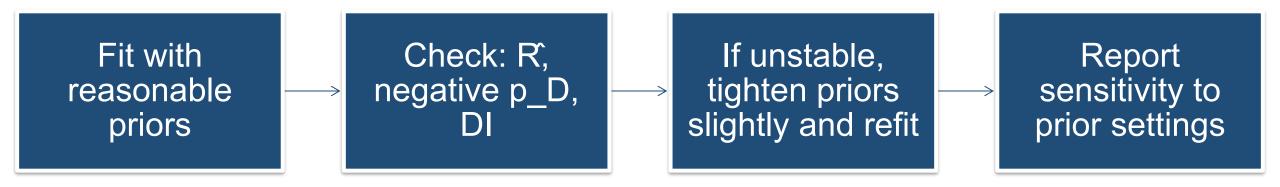


Use weakly informative fixed-effect priors with realistic scales

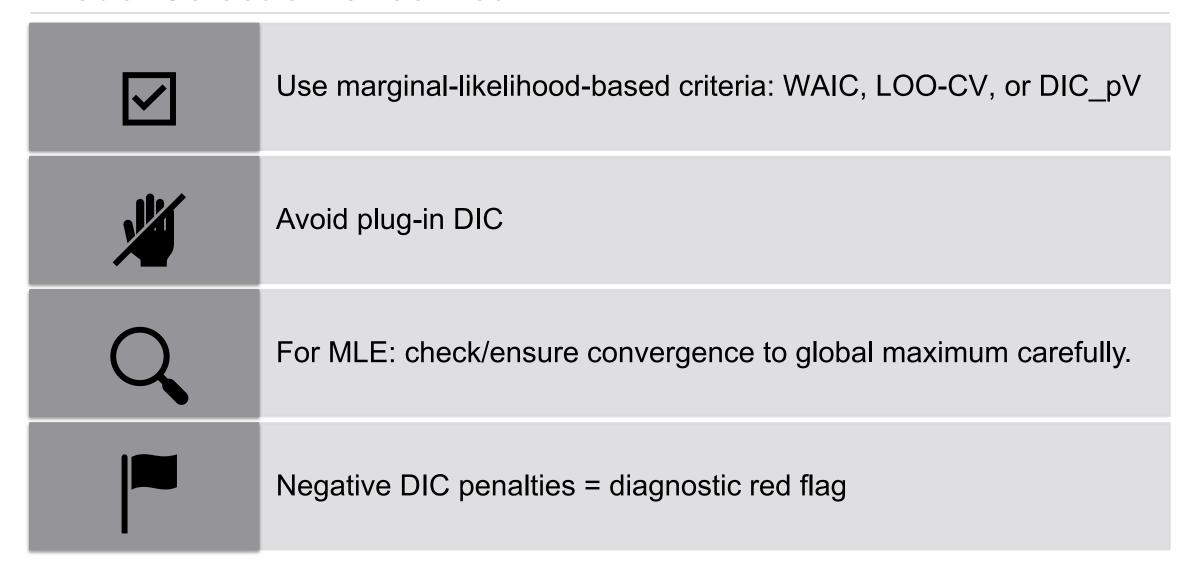


Tune priors in response to diagnostics

#### Minimal Workflow You Can Use Tomorrow



#### **Model-Selection Checklist**



# Thank You









Questions? Feedback welcome.

Seatbelts on, chairs filled, thermostat set, stethoscope ready.



https://doriaxiao.github.io/



#### References

- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using Ime4. *Journal of Statistical Software*, 67(1), 1–48. https://doi.org/10.18637/jss.v067.i01
- Gelman, A., Hwang, J., & Vehtari, A. (2014). Understanding predictive information criteria for Bayesian models. *Statistics and Computing*, 24, 997–1016.
- Grün, B., & Leisch, F. (2023). *flexmis: Flexible mixture modeling*. <a href="https://CRAN.R-project.org/package=flexmix">https://CRAN.R-project.org/package=flexmix</a> (R package version 2.3-19)
- Jung, T., & Wickrama, K. A. S. (2008). An introduction to latent class growth analysis and growth mixture modeling. *Social and Personality Psychology Compass*, 2(1), 302–317. <a href="https://doi.org/10.1111/j.1751-9004.2007.00054.x">https://doi.org/10.1111/j.1751-9004.2007.00054.x</a>
- Merkle, E. C., Furr, D., & Rabe-Hesketh, S. (2019). Bayesian comparison of latent variable models: Conditional versus marginal likelihoods. *Psychometrika*, 84(3), 802–829. <a href="https://doi.org/10.1007/s11336-019-09679-0">https://doi.org/10.1007/s11336-019-09679-0</a>
- Muthén, L. K., & Muthén, B. O. (1998-2017). Mplus user's guide. In www.statmodel.com (8th ed.). Muthén & Muthén.
- Plummer, M. (2017). JAGS version 4.3.0 User Manual. Retrieved from <a href="https://sourceforge.net/projects/mcmc-jags/files/Manuals">https://sourceforge.net/projects/mcmc-jags/files/Manuals</a>
- Spiegelhalter, D. J., Best, N. G., Carlin, B. P., & van der Linde, A. (2002). Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(4), 583–639. https://doi.org/10.1111/1467-9868.00353
- Stan Development Team. (2021). CmdStan User's Guide: Version 2.30. https://mc-stan.org/docs/cmdstan-guide/index.html
- StataCorp. (2023). Stata statistical software: Release 18. College Station: StataCorp LLC.
- Surhone, L. M., Tennoe, M. T., & Henssonow, S. F. (2010). OpenBUGS. Betascript Publishing.
- Vaida, F., & Blanchard, S. (2005). Conditional Akaike information for mixed-effects models. *Biometrika*, 92(2), 351–370. https://doi.org/10.1093/biomet/92.2.351
- Vehtari, A., Gelman, A., & Gabry, J. (2017). Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC. Statistics and Computing, 27, 1413–1432.
- Wardenaar, K. J. (2020). Latent class growth analysis and growth mixture modeling using R: A tutorial for two R-packages and a comparison with mplus. PsyArXiv. <a href="https://doi.org/10.31234/osf.io/m58wx">https://doi.org/10.31234/osf.io/m58wx</a>
- Xiao, X., Rabe-Hesketh, S., & Skrondal, A. (2025). Bayesian Identification and Estimation of Growth Mixture Models. *Psychometrika*, 1–34. doi:10.1017/psy.2025.11