

Revisiting the Gelman-Rubin Diagnostic¹: Improved Stability and a Principled Threshold

Christina Knudson, Ph.D.
University of St. Thomas
St. Paul, Minnesota

Symposium for Data Science and Statistics

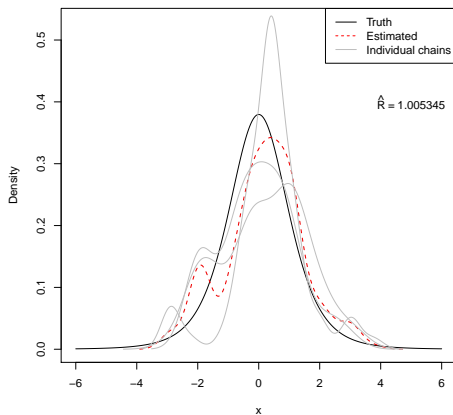
¹Joint work with Dootika Vats, Ph.D., University of Warwick 

Overview

Goal: Use MCMC to approximate a target distribution.
(e.g. an intractible posterior distribution)

Issue: After the chain has started sampling from the target distribution, how long should the sampler run to produce a decent approximation?

What if we terminate our MCMC sampler too early?



Density of T_5 estimated using 3 chains (each of length 150) produced with a Metropolis-Hastings sampler and proposal $N(\cdot, 2.6^2)$.

Gelman-Rubin Overview

Goal: run MCMC *long enough* to approximate a target distribution

Tool: Gelman-Rubin diagnostic (1992)

$$\hat{R} = \sqrt{\frac{\text{chain length} - 1}{\text{chain length}} + \frac{\text{between-chain variance}}{\text{within-chain variance}}}$$
$$\approx \sqrt{1 + \frac{\text{between-chain variance}}{\text{within-chain variance}}}$$

\hat{R} decreases to 1 as the chain length increases

Gelman-Rubin: $\hat{R} < 1.1$?

Gelman et al (2004):

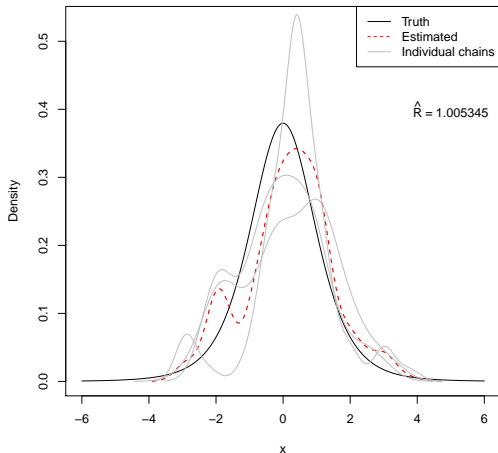
For most examples, values below 1.1 are acceptable, but for a final analysis in a critical problem, a higher level of precision may be required.

\hat{R} thresholds used in 100 papers from 2017:

\hat{R}	1.003	1.01	1.02	1.03	1.04	1.05	1.06	1.07	1.1	1.2	1.3
Freq.	1	12	9	9	2	11	2	1	43	9	1

Gelman-Rubin: $\hat{R} < 1.1$?

Reality: stopping at $\hat{R} = 1.1$ can be too early!



Vats and Knudson's Contributions

How can we improve the Gelman-Rubin diagnostic?

- 1 Stabilize the Gelman-Rubin statistic
- 2 Construct principled threshold for terminating simulation

Stabilizing the Gelman-Rubin Statistic

$$\hat{R} = \sqrt{\frac{\text{chain length} - 1}{\text{chain length}} + \frac{\text{between-chain variance}}{\text{within-chain variance}}}$$

Incorporate more sophisticated method of variance estimation

- Lugsail batch means (Vats and Flegal, 2019)
- More efficient (less variability in variance estimates)

Stabilizing the Gelman-Rubin Statistic

$$\hat{R} = \sqrt{\frac{\text{chain length} - 1}{\text{chain length}} + \frac{\text{between-chain variance}}{\text{within-chain variance}}}$$

Incorporate more sophisticated method of variance estimation

- Lugsail batch means (Vats and Flegal, 2019)
- More efficient (less variability in variance estimates)

Upshot: Less variability in variance estimates $\rightarrow \hat{R}$ stabilizes
 \rightarrow Time-to-termination stabilizes

Relevant R command: `stable.GR` in R package `stableGR`

(Can't cover details in 15 minutes. Sorry! Paper is on arXiv.)

Stabilizing the Gelman-Rubin Statistic

An AR(1) process

$$Y_t = .95 Y_{t-1} + \epsilon_t, \quad t = 1, 2, \dots$$
$$\epsilon_t \sim N(0, 1^2)$$

is the same as a Markov chain with distribution $N(0, 10.25641)$.

Stabilizing the Gelman-Rubin Statistic

An AR(1) process

$$Y_t = .95 Y_{t-1} + \epsilon_t, \quad t = 1, 2, \dots$$

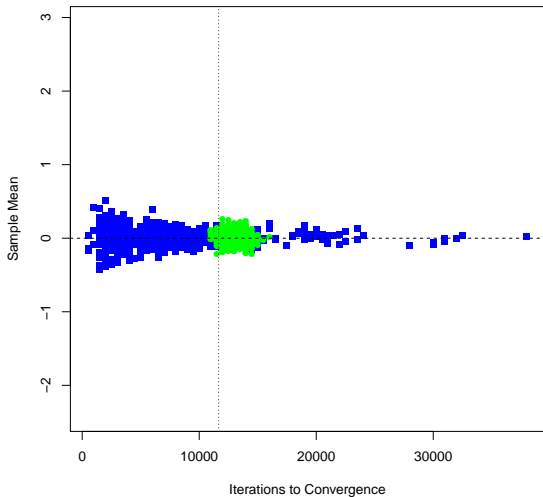
$$\epsilon_t \sim N(0, 1^2)$$

is the same as a Markov chain with distribution $N(0, 10.25641)$.

For each of 500 replications, we run five Markov chains until $\hat{R} < 1.001625$ using

- original GR \hat{R} calculation
- VK \hat{R} calculation

Stabilizing the Gelman-Rubin Statistic



A Principled Threshold for Terminating Simulation

Effective sample size: number of uncorrelated samples that produce the same precision as the correlated (MCMC) sample.

$$\hat{R} = \sqrt{\frac{\text{chain length} - 1}{\text{chain length}} + \frac{\text{number of chains}}{\text{effective sample size}}}$$

A Principled Threshold for Terminating Simulation

Effective sample size: number of uncorrelated samples that produce the same precision as the correlated (MCMC) sample.

$$\hat{R} = \sqrt{\frac{\text{chain length} - 1}{\text{chain length}} + \frac{\text{number of chains}}{\text{effective sample size}}}$$

Upshot:

- Threshold can be calculated *a priori*
 - Similar to introductory statistics sample size calculations for a desired width of a confidence interval
 - Gong and Flegal (2016) and Vats et al. (2019)
- \hat{R} threshold is easily-interpretable

Relevant commands in R package `stableGR`: `target.psrfr`, `n.eff`

A Principled Threshold for Terminating Simulation

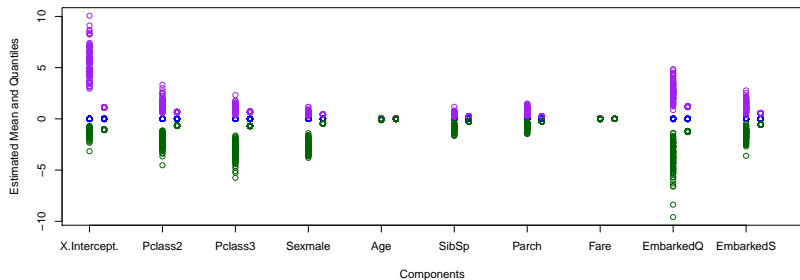
Model the log odds of surviving the Titanic's sinking with 10 predictors.

For each of 100 reps, we run 5 chains until convergence is diagnosed according to

- $\hat{R} < 1.1$
- VK's \hat{R} termination threshold

using VK's new \hat{R} calculation in both cases.

A Principled Threshold for Terminating Simulation



Centered posterior means (blue) and 95% credible interval estimates (green for lower bound, purple for upper bound).

Left points: $\hat{R} < 1.1$.

Right points: VK's new \hat{R} threshold.

Concluding Remarks

To review, we have:

- Stabilized the Gelman-Rubin statistic \hat{R} .
- Identified a one-to-one relationship between ESS and \hat{R} .
- Created an interpretable stopping rule to replace $\hat{R} < 1.1$.

Concluding Remarks

To review, we have:

- Stabilized the Gelman-Rubin statistic \hat{R} .
- Identified a one-to-one relationship between ESS and \hat{R} .
- Created an interpretable stopping rule to replace $\hat{R} < 1.1$.

Additional information:

- Diagnostic is usable for multiple chains or a single chain.
- We have also stabilized the multivariate version of the Gelman-Rubin statistic and produced an interpretable stopping rule for multivariate chains.
- R Package `stableGR` will be available on CRAN shortly. These slides end with instructions for installing it from Github.

cknudson.com

has these slides,
a link to the Github repo for R package `stableGR`,
instructions for installing `stableGR`,
and “Revisiting the Gelman-Rubin Diagnostic” (Vats and Knudson)

knud8583@stthomas.edu

References

- Brooks, S. P. and Gelman, A. (1998). General methods for monitoring convergence of iterative simulations. *Journal of Computational and Graphical Statistics*, 7:434-455.
- Gelman, A., Carlin, J. B., Stern, H. S., and Rubin, D. B. (2004). *Bayesian Data Analysis*. Chapman & Hall/CRC, Boca Raton, FL.
- Gelman, A. and Rubin, D. B. (1992). Inference from iterative simulation using multiple sequences (with discussion). *Statistical Science*, 7:457-472.
- Gong, L. and Flegal, J. M. (2016). A practical sequential stopping rule for high-dimensional Markov chain Monte Carlo. *Journal of Computational and Graphical Statistics*, 25:684-700.
- Vats, D. and Flegal, J. M. (2018). Lugsail lag windows and their application to MCMC. *arXiv e-prints*.
- Vats, D., Flegal, J. M, and Jones, G. L. (2019). Multivariate output analysis for Markov chain Monte Carlo. *Biometrika*, 106(2):321-337.

Commands for Installing R Package stableGR from Github

```
#get some required packages
install.packages("Rcpp")
install.packages("RcppArmadillo")
install.packages("devtools")

#install mcmcse from github (rather than CRAN)
library(devtools)
install_github("dvats/mcmcse")

#install stableGR package
install_github("knudson1/stableGR/stableGR")
library(stableGR)
```

Will be available on CRAN in a couple months.

Bayesian Logistic Regression

Model the log odds of surviving the Titanic's sinking.

Bayesian logistic regression with the following predictors:

- Fare class (3 categories)
- Sex (2 categories)
- Age (quantitative)
- Number of siblings/spouses aboard (quantitative)
- Number of parents/children aboard (quantitative)
- Port of embarkation (3 categories)

For each of 100 reps, run 5 chains til convergence is diagnosed (according to both $\hat{R} = 1.1$ and VK's new threshold).

Bayesian Logistic Regression

