

# A simple, fast sampler for simulating spatial data and other Markovian data structures + intRo

Andee Kaplan

Iowa State University  
ajkaplan@iastate.edu

November 29, 2016

# A simple, fast sampler for simulating spatial data and other Markovian data structures

# Goal

- Markov random field models are possible for spatial or network data
- Rather than specifying a joint distribution directly, a model is specified through a set of full conditional distributions for each spatial location
- Assume the spatial data are on a regular lattice (wrapped on a torus)

**Goal:** A new, provably fast approach for simulating spatial/network data.

# Spatial Markov random field (MRF) models

## Notation

- Variables  $\{Y(\mathbf{s}_i) : i = 1, \dots, n\}$  at locations  $\{\mathbf{s}_i : i = 1, \dots, n\}$
- Neighborhoods:  $N_i$  specified according to some configuration
- Neighboring Values:  $\mathbf{y}(N_i) = \{y(\mathbf{s}_j) : \mathbf{s}_j \in N_i\}$
- Full Conditionals:  $\{f_i(y(\mathbf{s}_i) | \mathbf{y}(N_i), \boldsymbol{\theta}) : i = 1, \dots, n\}$ 
  - $f_i(y(\mathbf{s}_i) | \mathbf{y}(N_i), \boldsymbol{\theta})$  is conditional pmf/pdf of  $Y(\mathbf{s}_i)$  given values for its neighbors  $\mathbf{y}(N_i)$
  - Often assume a common conditional cdf  $F_i = F$  form ( $f_i = f$ ) for all  $i$

# Exponential family examples

- 1 Conditional Gaussian (3 parameters):

$$f_i(y(\mathbf{s}_i) | \mathbf{y}(N_i), \alpha, \eta, \tau) = \frac{1}{\sqrt{2\pi\tau}} \exp\left(-\frac{[y(\mathbf{s}_i) - \mu(\mathbf{s}_i)]^2}{2\tau^2}\right)$$

$Y(\mathbf{s}_i)$  given neighbors  $\mathbf{y}(N_i)$  is normal with variance  $\tau^2$  and mean

$$\mu(\mathbf{s}_i) = \alpha + \eta \sum_{\mathbf{s}_j \in N_i} [y(\mathbf{s}_j) - \alpha]$$

- 2 Conditional Binary (2 parameters):

$Y(\mathbf{s}_i)$  given neighbors  $\mathbf{y}(N_i)$  is Bernoulli  $p(\mathbf{s}_i, \kappa, \eta)$  where

$$\text{logit}[p(\mathbf{s}_i, \kappa, \eta)] = \text{logit}(\kappa) + \eta \sum_{\mathbf{s}_j \in N_i} [y(\mathbf{s}_j) - \kappa]$$

In both examples,  $\eta$  represents a dependence parameter.

# Concliques

## Cliques – Hammersley and Clifford [1971]

Singletons and sets of locations such that each location in the set is a neighbor of all other locations in the set

Example: Four nearest neighbors gives cliques of sizes 1 and 2

## The Converse of Cliques – Concliques

Sets of locations such that no location in the set is a neighbor of any other location in the set

4 Nearest  
Neighbors

.	*	.
*	<b>S</b>	*
.	*	.

Concliques  
4 Nearest  
Neighbors

1	2	1	2
2	1	2	1
1	2	1	2
2	1	2	1

8 Nearest  
Neighbors

*	*	*
*	<b>S</b>	*
*	*	*

Concliques  
8 Nearest  
Neighbors

1	2	1	2
3	4	3	4
1	2	1	2
3	4	3	4

# Generalized spatial residuals

## Definition

- $F(y|\mathbf{y}(N_i), \boldsymbol{\theta})$  is the conditional cdf of  $Y(\mathbf{s}_i)$  under the model
- Substitute random variables,  $Y(\mathbf{s}_i)$  and neighbors  $\{Y(\mathbf{s}_j) : \mathbf{s}_j \in N_i\}$ , into (continuous) conditional cdf to define residuals:

$$R(\mathbf{s}_i) = F(Y(\mathbf{s}_i) | \{Y(\mathbf{s}_j) : \mathbf{s}_j \in N_i\}, \boldsymbol{\theta}).$$

## Key Property

Let  $\{\mathcal{C}_j : j = 1, \dots, q\}$  be a collection of cliques that partition the integer grid. Under the conditional model, **spatial residuals within a clique are iid Uniform(0, 1)-distributed**:

$$\{R(\mathbf{s}_i) : \mathbf{s}_i \in \mathcal{C}_j\} \stackrel{iid}{\sim} \text{Uniform}(0, 1) \quad \text{for } j = 1, \dots, q$$

(Kaiser et al. [2012])

# Conclique-based Gibbs sampler

Using the conditional independence of random variables at locations within a conclique along with the probability integral transform we propose a conclique-based Gibbs sampling algorithm for sampling from a MRF.

- ① Split locations into  $Q$  disjoint concliques,  $\mathcal{D} = \cup_{i=1}^Q \mathcal{C}_i$ .
- ② Initialize the values of  $\{Y^{(0)}(\mathbf{s}) : \mathbf{s} \in \{\mathcal{C}_2, \dots, \mathcal{C}_Q\}\}$ .
- ③ Sample from the conditional distribution of  $Y(\mathbf{s})$  given  $\{Y(\mathbf{t}) : \mathbf{t} \in \mathcal{N}(\mathbf{s})\}$  for  $\mathbf{s} \in \mathcal{C}_1$ ,
  - ① Sample  $\{U(\mathbf{s}) : \mathbf{s} \in \mathcal{C}_1\} \stackrel{iid}{\sim} Unif(0, 1)$
  - ② For each  $\mathbf{s} \in \mathcal{C}_1$ ,  $Y^{(i)}(\mathbf{s}) = F^{-1}(U(\mathbf{s}) | Y^{(i-1)}(\mathbf{t}), \mathbf{t} \in \mathcal{N}(\mathbf{s}))$
- ④ Sample from the conditional distribution of  $Y(\mathbf{s})$  given  $\{Y(\mathbf{t}) : \mathbf{t} \in \mathcal{N}(\mathbf{s})\}$  for  $\mathbf{s} \in \mathcal{C}_j; j = 2, \dots, Q$ ,
  - ① Sample  $\{U(\mathbf{s}) : \mathbf{s} \in \mathcal{C}_j\} \stackrel{iid}{\sim} Unif(0, 1)$
  - ② For each  $\mathbf{s} \in \mathcal{C}_j$ ,  $Y^{(i)}(\mathbf{s}) = F^{-1}(U(\mathbf{s}) | \{Y^{(i)}(\mathbf{t}), \mathbf{t} \in \mathcal{N}(\mathbf{s}) \cap \mathcal{C}_k \text{ where } k < j\}, \{Y^{(i-1)}(\mathbf{t}), \mathbf{t} \in \mathcal{N}(\mathbf{s}) \cap \mathcal{C}_k \text{ where } k > j\})$



# It's (provably) fast!

- ① In many (commonly used) four-nearest neighbor models (including Gaussian and binary), the conclique-based Gibbs sampler is provably **geometrically ergodic**.
- ② Because we are using batch updating vs. sequential updating of each location, this approach is also **computationally fast**.
- ③ A flexible R package using Rcpp (called `conclique`, to appear on CRAN) that implements a conclique-based Gibbs sampler while allowing the user to specify an arbitrary model.

# Preliminary simulations

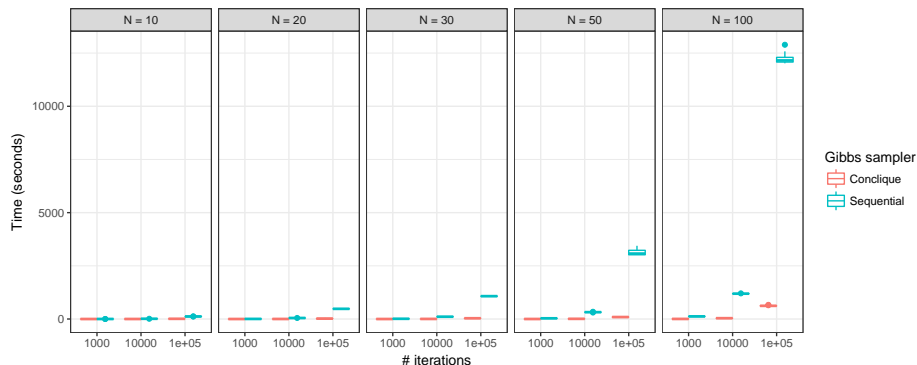


Figure 1: Comparisons of timing for simulation of 4NN Gaussian Markov Random Field data on a lattice of size  $N \times N$  for various size grids,  $N = 10, 20, 30, 50, 100$ , using sequential and conclave-based Gibbs samplers.

# Ideas and connections

- Goodness-of-fit test for network data
  - The model-based method of resampling re-frames network into a collection of (Markovian) neighborhoods by using covariate information
  - Creates cliques on a graph structure
  - Use a conditionally specified network distribution (Casleton et al. [2016]) to sample network data in a blockwise clique-based Gibbs sampler.
- Extension for multilayer networks
  - Layer-level dependence parameter
  - Utilization of single layer for neighborhood creation

## Other ideas

- Implementation of existing methods in Rcpp (e.g. ebLink, SMERED)
- Block bootstrap of dynamic networks over time
  - Potential nonparametric method of determining distributional properties of dynamic network statistics
  - Combine with GOF tests
- Dual record linkage problem - privacy

# intRo: statistical analysis software for teaching

# Do we really need another statistical software package?

Short answer: **yes**

- R is great, but requires students to have some knowledge/interest in programming
- JMP, Deducr, Rcmdr are powerful, but too big
  - Licenses and installation
- New tools recently released to spark an interest in R
  - Swirl and DataCamp teach R programming
  - Project MOSAIC facilitates learning, but assumes knowledge of R
- Want students to focus on data analysis rather than fight with software

# What is intRo?

- A simple **web-based** application for performing basic data analysis and statistical routines and accompanying utility package
- Built using R and Shiny
- Extensible modular structure
- Designed for a first statistics class student
- Assists in the learning of statistics rather than acting as a stand-alone deliverer of statistics education

- Focused on aspects of the user interface (UI) and output that make it easy to pick up without training
- Minimal necessary functionality for an introductory statistics course
- Organized around specific tasks a student may perform in the process of a data analysis



# Exciting

- Fun, easy to use (available on the web)
- Interactive plots using `ggvis`

**Ultior motive:** get students excited about programming

- By navigating about the user interface of `intRo`, students are creating a fully-executable R script that they can download and run locally
- Viewing their script change real-time within the application

# Extensible

- User interaction with `intRo` is split into bitesize chunks that we call *modules*
- Each module is a self contained set of R code that is dynamically added to the application at run time
- `intRo` can be easily extended by the addition of modules within the frame-work underlying the application
- Allows instructors/collaborators to tailor `intRo` to the needs of a particular course

Take a look <http://intro-stats.com>

# intRo package

## Installation:

```
devtools::install_github("rstudio/shinyapps")  
devtools::install_github("gammarama/intRo")`
```

## Functions:

- `download_intRo` - Downloads a current revision of `intRo` to your machine
- `run_intRo` - Runs an `intRo` session locally with the specified options
- `deploy_intRo` - Deploys an instance of `intRo` to ShinyApps.io with the specified options

## Future work - Module creation

*Modularity* is a key feature of `intRo`, but module creation is currently:

- Undocumented
- Entirely manual
- Unnecessarily lengthy

**Idea:** Include functionality in current R package to automate creation of `intRo` modules

# References I

- Emily Casleton, Daniel Nordman, and Mark Kaiser. A local structure model for network analysis. *Statistics and Its Interface, Forthcoming*, 2016.
- John M Hammersley and Peter Clifford. Markov fields on finite graphs and lattices. 1971.
- Mark S Kaiser, Soumendra N Lahiri, and Daniel J Nordman. Goodness of fit tests for a class of markov random field models. *The Annals of Statistics*, pages 104–130, 2012.