

# Take Me Out to (Analyze) the Ballgame

## Visualization and Analysis Techniques for Big Spatial Data

Chris Comiskey

Oregon State University

August 30, 2017

# Baseball, Baseball, Baseball

- Chris, rookie year

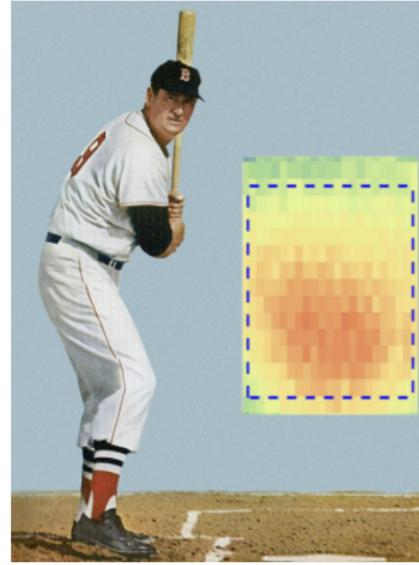
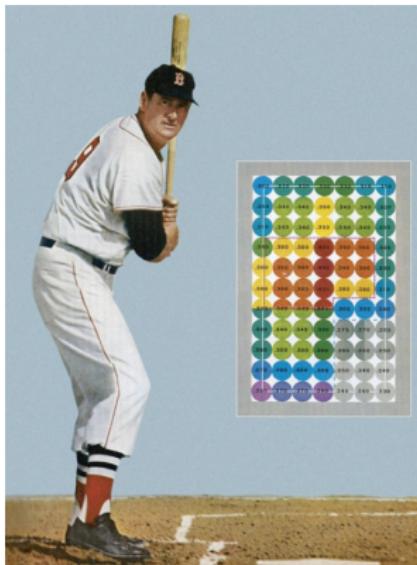


- Chris, Boston Red Sox



# Hitting Analytics

- “The Science of Hitting”<sub>1970</sub>
- The Statistics of Hitting



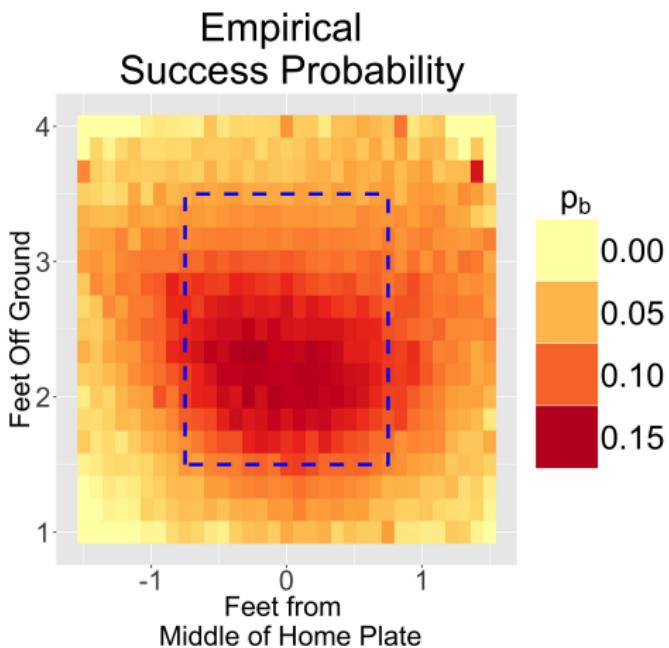
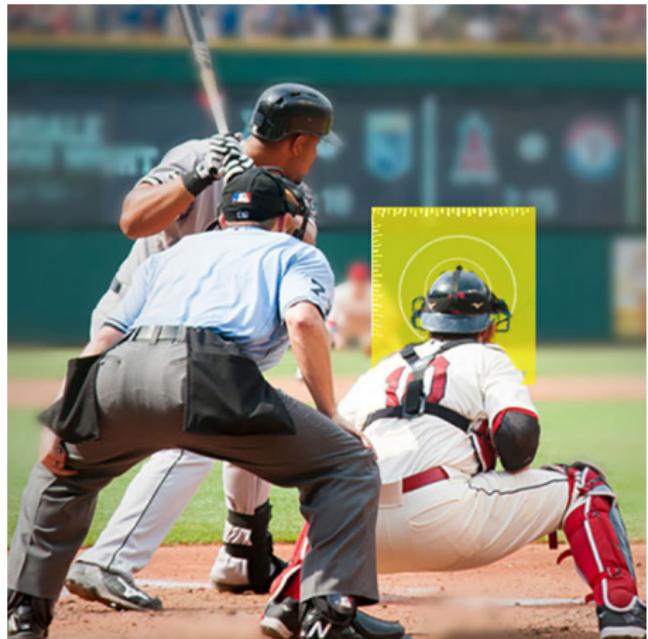
- Iconic breakthrough, hitter
- No data

- PITCHf/x data, R, heat maps
- SGLMMs, Stan, PPMs, INLA

# Outline

- 1 Variable-Resolution Heat Maps
- 2 Interactive Heat Map Confidence Intervals
- 3 Approaches to Big Data Spatial Mixed Models for Baseball Data
  - Computational Optimization in Stan
  - Predictive Process Models
  - Integrated Nested Laplace Approximation

# Empirical Success Probability Heat Map



# Empirical Success Probability Heat Map

- Swings:  $i = 1, 2, \dots, N$

- Bernoulli( $\pi_i$ ) trials:

$$S_i = \begin{cases} 1; & \text{swing success} \\ 0; & \text{swing failure} \end{cases}$$

- Location  $(x_i, y_i)$

- Grid boxes  $G_1, G_2, \dots, G_B$

- Box  $b$  totals:

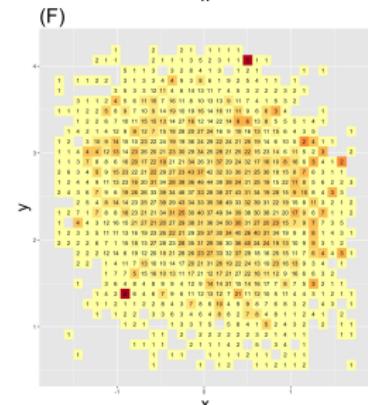
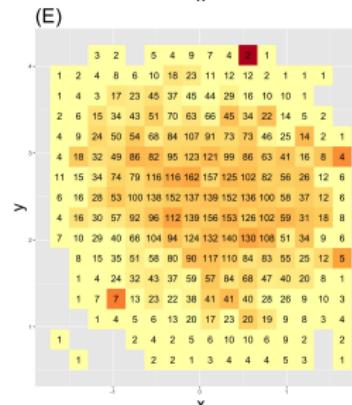
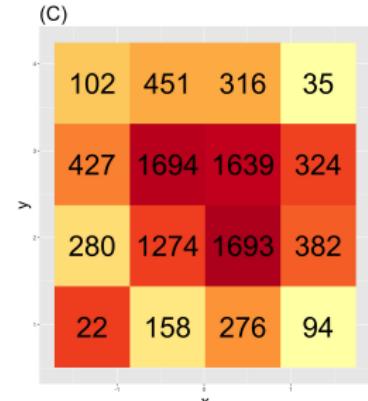
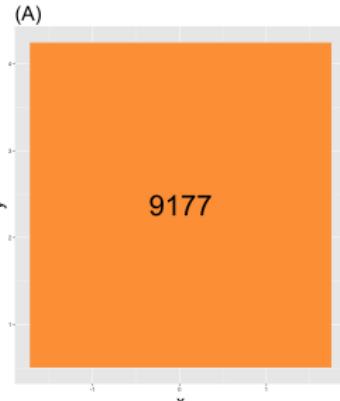
$$N_b = \sum_{i=1}^N I_{(x_i, y_i) \in G_b}$$

- Box  $b$  empirical success:

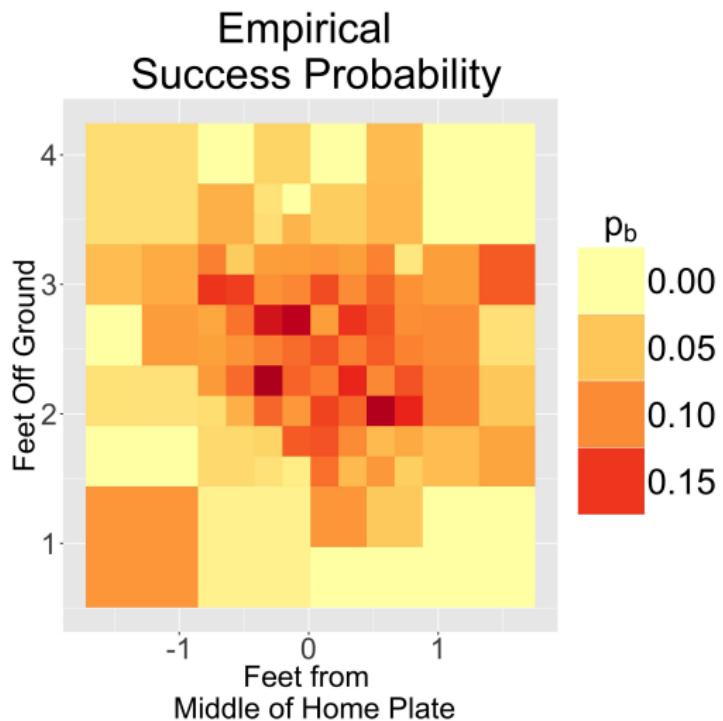
$$p_b = \frac{1}{N_b} \sum_{i=1}^N S_i I_{(x_i, y_i) \in G_b}$$

# Heat Map Resolution Selection

Jhonny Peralta



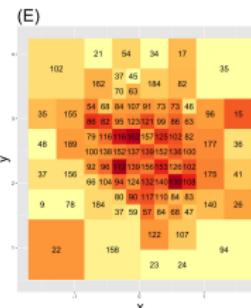
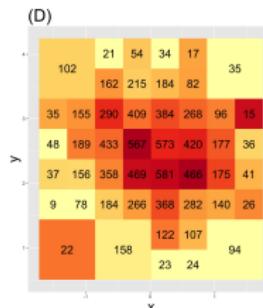
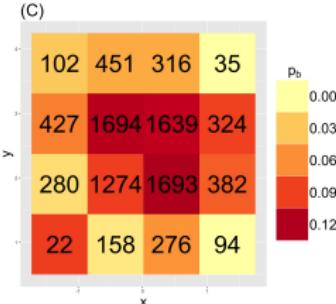
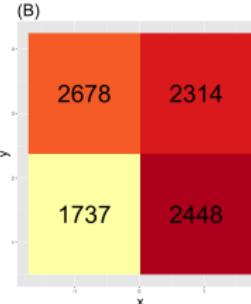
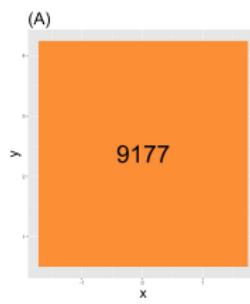
# A Thing of Beauty



# Variable-Resolution Heat Maps

Jhonny Peralta

- Stopping rule
- Subdivision method
- VR algorithm



# Variable-Resolution Heat Maps

## Combine Resolutions

- Stopping rule: sample size threshold
- Combine resolutions
- Resolution conveys data density
- Improvements:
  - ▶ Subdivision methods
  - ▶ Subdivision criteria

# varyres(...) in **varyres**

`varyres {varyres}`

R Documentation

## A variable-resolution heat map generator

### Description

This function creates variable resolution heat maps according to a stopping rule

### Usage

```
varyres(dataset, cutoff, fun = mean, max = 6)
```

### Arguments

**dataset** data frame with spatial data: x-coordinates (x), y coordinates (y), and Bernoulli responses at those locations (res)

**cutoff** Box subdivisions cease when a box sample size drops below the cutoff

**fun** Function to apply to responses in each box

**max** The maximum number of subdivision iterations the algorithm will perform

### Value

A list containing a data frame for each iteration of the subdivision algorithm; and a vector of the number of boxes eligible for subdivision at each iteration.

### Examples

```
data(hitter)
data <- varyres(hitter, mean, cutoff = 200, max = 6)
mapit(data[[4]])
```

# Outline

- 1 Variable-Resolution Heat Maps
- 2 Interactive Heat Map Confidence Intervals
- 3 Approaches to Big Data Spatial Mixed Models for Baseball Data
  - Computational Optimization in Stan
  - Predictive Process Models
  - Integrated Nested Laplace Approximation

# Generalized Linear Model

$$Y_i | \mathbf{X}_i(\mathbf{s}_i) \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\pi_i)$$

$$\text{logit}(\pi_i | \mathbf{s}_i) = \mathbf{X}_i(\mathbf{s}_i)\beta$$

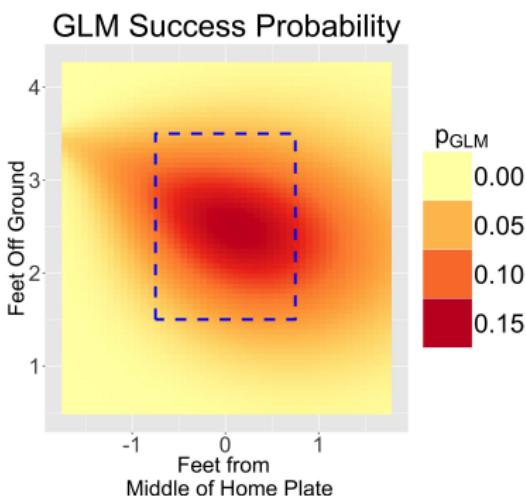
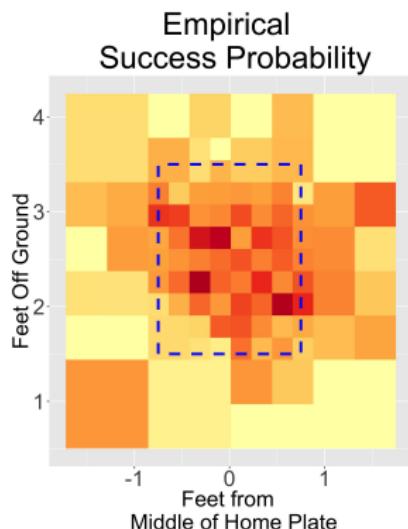
- Swings:  $i = 1, 2, \dots, N$
- Pitch location:  $\mathbf{s}_i = (x_i, y_i)$
- Biomechanical covariates:  $\mathbf{X}_i(\mathbf{s}_i)$
- Fit to Jhonny Peralta data

# Logistic Regression Model

Jhonny Peralta

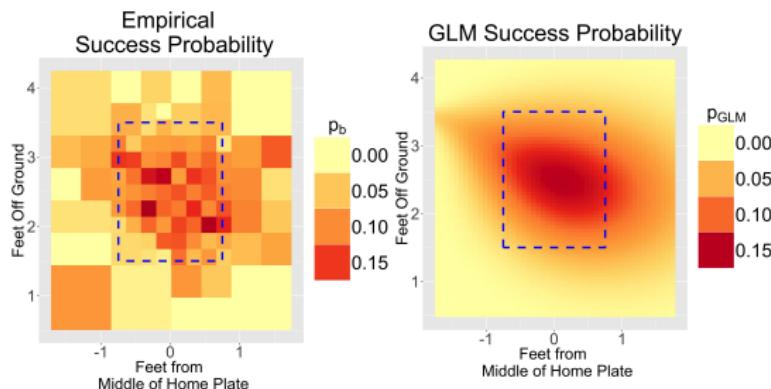
$$Y_i | \mathbf{X}_i(\mathbf{s}_i) \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\pi_i)$$

$$\text{logit}(\pi_i | \mathbf{s}_i) = \mathbf{X}_i(\mathbf{s}_i)\beta,$$



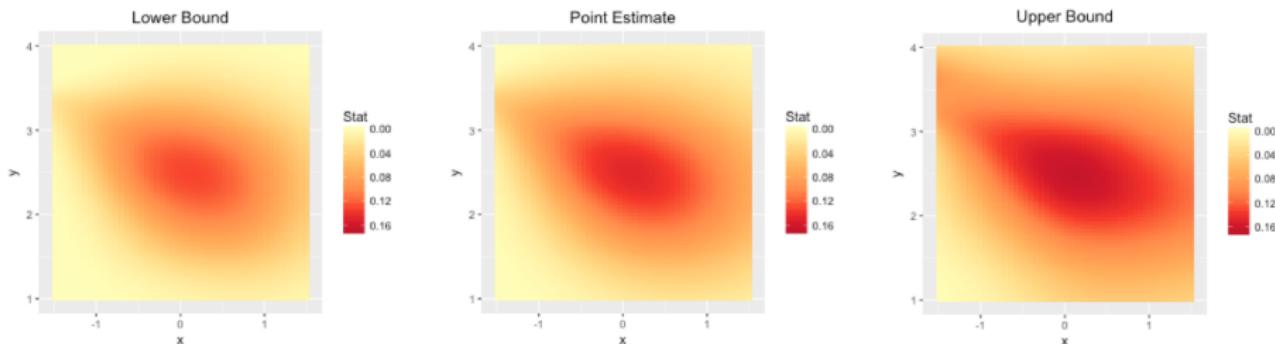
# Logistic Regression Model

Jhonny Peralta



- Hosmer-Lemeshow (logistic regression) GOF test
  - ▶  $H_0$ : Well fit
  - ▶  $H_A$ : Lack of fit
  - ▶ p-value = 0.8217
- Confidence intervals

# Interactive Confidence Intervals



Confidence Interval



This is the 99 % confidence interval layer.

**mapapp:** An R Package

# mapapp: An R Package

- Challenge: Where does user stop and package begin?
- Future improvement

```
all_in_one <- get_CI(model, x, y, levels)
shinyHMCI(all_in_one)
```

- ➊ `get_CI(...)` — creates proper data structure
- ➋ `shinyHMCI(all_in_one)` — creates Shiny app

# Outline

- 1 Variable-Resolution Heat Maps
- 2 Interactive Heat Map Confidence Intervals
- 3 Approaches to Big Data Spatial Mixed Models for Baseball Data
  - Computational Optimization in Stan
  - Predictive Process Models
  - Integrated Nested Laplace Approximation

# Spatial Generalized Linear Mixed Model (SGLMM)

$$Y_i | \mathbf{X}_i(\mathbf{s}_i) \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\pi_i)$$

$$\text{logit}(\pi_i | \mathbf{s}_i) = \mathbf{X}_i(\mathbf{s}_i)\boldsymbol{\beta} + w(\mathbf{s}_i)$$

- $w(\mathbf{s}_i)$  — spatial random effect, location  $\mathbf{s}_i$ .
- $w(\mathbf{s}) = (w(\mathbf{s}_1), w(\mathbf{s}_2), \dots, w(\mathbf{s}_N))$  — vector
- $w(\mathbf{s})$  - Gaussian Random Field (GRF)

# Gaussian Random Field: $\mathbf{w}(\mathbf{s})$

$$\mathbf{w}(\mathbf{s})|\boldsymbol{\theta} \sim \text{MVN}(\mathbf{0}, \Sigma(\boldsymbol{\theta}))$$

$$\Sigma(\phi, \sigma^2)_{i,k} = \sigma^2 \exp(-||\mathbf{s}_i - \mathbf{s}_k||/\phi)$$

- Spatial exponential covariance
  - ▶  $||\mathbf{s}_i - \mathbf{s}_k||$  - Euclidean distance
  - ▶  $\sigma^2$  - scale parameter
  - ▶  $\phi$  - range parameter.
- Notice:  $\Sigma(\boldsymbol{\theta}) — n \times n$

## Computational Cost: “Big N” Problem

- Peralta:  $n = 9177$
- $\mathbf{w}(\mathbf{s})$ :  $9177 \times 9177$  cov. matrix
- MCMC iterations require:  $\Sigma^{-1}$ , determinant of  $\Sigma$
- $\mathcal{O}(n^3)$  rate of increase:

$$t(n) \leq M \cdot n^3$$

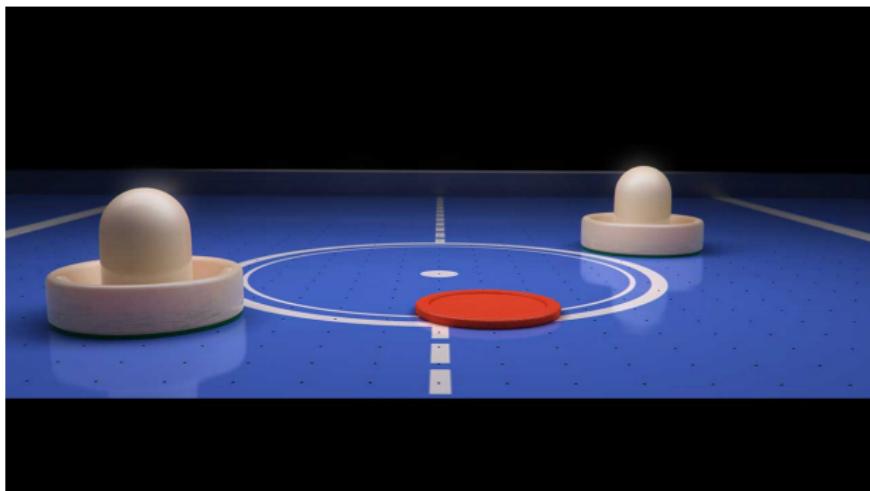
as  $n \rightarrow \infty$

# Outline

- 1 Variable-Resolution Heat Maps
- 2 Interactive Heat Map Confidence Intervals
- 3 Approaches to Big Data Spatial Mixed Models for Baseball Data
  - Computational Optimization in Stan
  - Predictive Process Models
  - Integrated Nested Laplace Approximation

# Hamiltonian Monte Carlo (HMC)

- Stan uses Hamiltonian proposal mechanism
- **Short version:** disk on surface, randomly sample momentum (auxiliary), calculate new position (parameters) — that's your Metropolis proposal.



# Stan Computational Optimization

- $\phi, \beta$ : informative/proper priors  $\rightarrow$  identifiability/cost/convergence
- QR factorization of X, Cholesky decomposition of  $\Sigma(\theta)$
- Matrices & vectors faster than loops & scalars
- n = 25: 40 seconds  $\rightarrow$  3 seconds
- n = 2000 — overnight, 350 draws



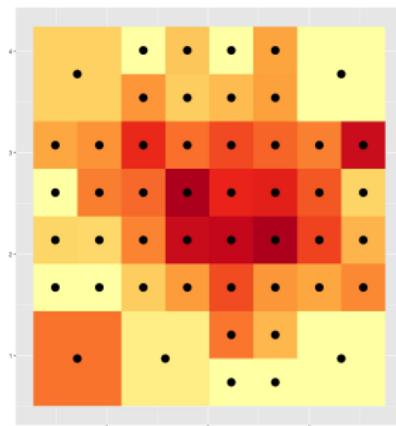
# Outline

- 1 Variable-Resolution Heat Maps
- 2 Interactive Heat Map Confidence Intervals
- 3 Approaches to Big Data Spatial Mixed Models for Baseball Data
  - Computational Optimization in Stan
  - **Predictive Process Models**
  - Integrated Nested Laplace Approximation

# Predictive Process Models (PPMs)

[Banerjee et al., 2008]

- Knots:  $\mathbf{S}^* = \{\mathbf{s}_1^*, \dots, \mathbf{s}_m^*\}$ 
  - ▶  $m \ll n$



$$\text{logit}(\pi_i | \mathbf{s}_i) = \mathbf{X}_i(\mathbf{s}_i)\boldsymbol{\beta} + \tilde{\mathbf{w}}(\mathbf{s}_i)$$

$$\tilde{\mathbf{w}}(\mathbf{s}) \sim \text{MVN}\{0, \tilde{\Sigma}(\boldsymbol{\theta})\}$$

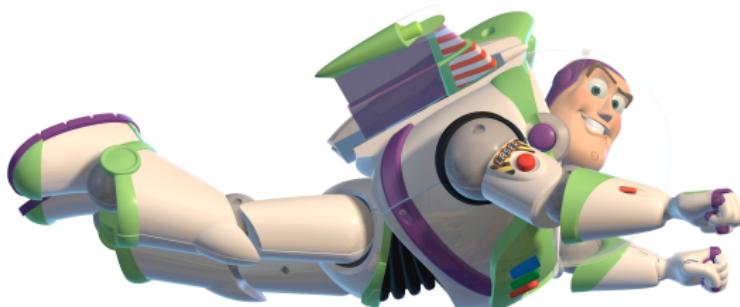
$$\tilde{\Sigma}(\boldsymbol{\theta})_{i,j} = \sigma^{*T}(\mathbf{s}_i; \boldsymbol{\theta}) \cdot \Sigma^{*-1}(\boldsymbol{\theta}) \cdot \sigma^*(\mathbf{s}_j; \boldsymbol{\theta})$$

- $\sigma^*(\mathbf{s}_i; \boldsymbol{\theta}) = \text{Cov}(\mathbf{s}_i, \mathbf{S}^*)$

- $\Sigma^*(\boldsymbol{\theta}) = \text{Var}(\mathbf{S}^*)$

# PPM Results

- Implement in **spBayes** [Finley et al., 2013].
- MCMC chains **did not converge**.
  - ▶  $n = 1000$ , knots = 97, 10K samples,  $\approx 6.7$  mins
  - ▶  $n = 1000$ , knots = 49, 30K samples,  $\approx 7$  mins
  - ▶  $n = 3000$ , knots = 49, 80K samples,  $\approx 54$  mins
- Speed issue for extending MCMC chains
- Buzz Lightyear is fast.



# Outline

- 1 Variable-Resolution Heat Maps
- 2 Interactive Heat Map Confidence Intervals
- 3 Approaches to Big Data Spatial Mixed Models for Baseball Data
  - Computational Optimization in Stan
  - Predictive Process Models
  - Integrated Nested Laplace Approximation

# INLA Overview

$$\text{logit}(\pi_i) = \mathbf{X}_i(\mathbf{s}_i)\boldsymbol{\beta} + \mathbf{w}(\mathbf{s}_i)$$
$$\mathbf{w}(\mathbf{s})|\boldsymbol{\theta} \sim \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}(\boldsymbol{\theta}))$$

- Bayesian hierarchical models w/ latent GRF
- Assume Matérn covariance
- Two parts (continuous domain)
  - ▶ Part 1: Represent GRF as GMRF (SPDE)
  - ▶ Part 2: Integrated Nested Laplace Approximation (INLA)



# Integrated Nested Laplace Approximation

[Rue et al., 2009]

- Parameter vector:  $\rho = (\beta^T, \tilde{w}^T)^T$ 
  - ▶  $Q$ : precision matrix of  $\rho$
- Hyperparameter vector:  $\theta = (\kappa, \sigma)$
- ① Gaussian approximation:

$$p(\rho|\theta, y) \propto p(y|\rho, \theta)p(\rho|\theta)p(\theta)$$

$$p(\rho|\theta, y) \propto \exp\left(-\frac{1}{2}\rho^T Q \rho + \sum_i \log p(y_i|\rho, \theta)\right)$$

$$p_G(\rho|\theta, y) \propto \exp\left(-\frac{1}{2}(\rho - \mu)^T (Q + \text{diag}(c))(\rho - \mu)\right)$$

- ▶  $c$  and  $\mu$  depend on second order Taylor expansions of  
 $f(\rho) = \sum_i \log p(y_i|\rho, \theta)$

# Integrated Nested Laplace Approximation

[Rue et al., 2009]

- Fact:  $P(A, B) = P(A|B)P(B) = P(B|A)P(A)$

$$p(\mathbf{y}|\boldsymbol{\theta}) = \frac{p(\mathbf{y}|\boldsymbol{\rho}, \boldsymbol{\theta})p(\boldsymbol{\rho}|\boldsymbol{\theta})}{p(\boldsymbol{\rho}|\mathbf{y}, \boldsymbol{\theta})}$$

- Bayes proportionality:

$$\begin{aligned} p(\boldsymbol{\theta}|\mathbf{y}) &\propto p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \\ &\propto \frac{p(\mathbf{y}|\boldsymbol{\rho}, \boldsymbol{\theta})p(\boldsymbol{\rho}|\boldsymbol{\theta})}{p(\boldsymbol{\rho}|\mathbf{y}, \boldsymbol{\theta})} \cdot p(\boldsymbol{\theta}) \end{aligned}$$

- For a given  $\boldsymbol{\theta}$ , let  $\boldsymbol{\rho}_0 = \operatorname{argmax}_{\boldsymbol{\rho}} p(\boldsymbol{\rho}|\mathbf{y}, \boldsymbol{\theta})$ . Then,

$$\tilde{p}(\boldsymbol{\theta}|\mathbf{y}) \propto \frac{p(\mathbf{y}|\boldsymbol{\rho}_0, \boldsymbol{\theta})p(\boldsymbol{\rho}_0|\boldsymbol{\theta})}{p_G(\boldsymbol{\rho}_0|\mathbf{y}, \boldsymbol{\theta})} \cdot p(\boldsymbol{\theta})$$

# Integrated Nested Laplace Approximation

[Rue et al., 2009]

- ③ Numerical integration:

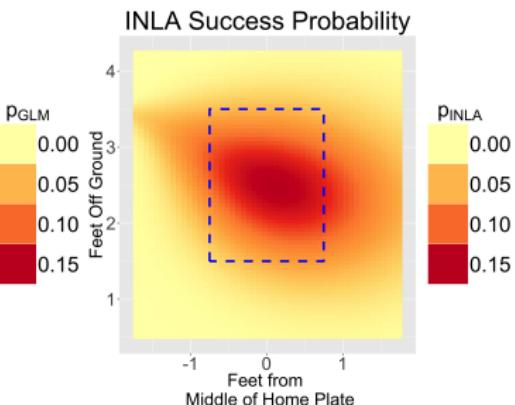
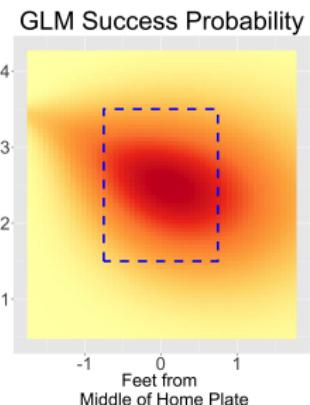
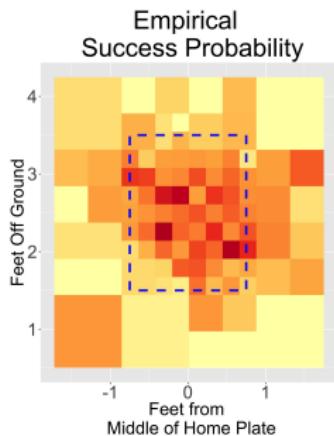
$$p(\rho_j | \mathbf{y}) \approx \int p_G(\rho_j | \boldsymbol{\theta}, \mathbf{y}) \tilde{p}(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}$$

- ④ Numerical integration:

$$p(\theta_k | \mathbf{y}) \approx \int \tilde{p}(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}_{-k}$$

# INLA Model Fit: 34 seconds

34 seconds!



# Estimates

GLM			
Covariate	$\beta_i$	MLE	SE
1	$\beta_0$	-4.08	0.70
r	$\beta_1$	1.19	0.51
$\theta$	$\beta_2$	-1.93	1.90
$r * \theta$	$\beta_3$	-1.64	0.70
$r^2$	$\beta_4$	-0.32	0.09
$\theta^2$	$\beta_5$	-3.92	1.10
$r^2 * \theta^2$	$\beta_6$	-0.46	0.21

INLA			
Covariate	$\rho_i$	$\hat{\rho}_i$	SE
N/A	$\kappa$	3.23	$\pm 1 \text{ SE}: (1.35, 7.54)$
N/A	$\sigma$	0.11	$\pm 1 \text{ SE}: (0.05, 0.26)$
1	$\beta_0$	-4.14	0.76
r	$\beta_1$	1.25	0.55
$\theta$	$\beta_2$	-1.90	1.96
$r * \theta$	$\beta_3$	-1.70	0.93
$r^2$	$\beta_4$	-0.33	0.10
$\theta^2$	$\beta_5$	-3.93	1.14
$r^2 * \theta^2$	$\beta_6$	-0.48	0.22

- $\hat{\kappa}$  — long range correlation
- $\text{SE}(w(s)) = 0.11$ .
- $\rightarrow 0.15 \pm 1 \cdot \text{SE} = (0.137, 0.165)$
- $\rightarrow 0.15 \pm 2 \cdot \text{SE} = (0.125, 0.181)$

# Summary

- Resolution selection? Variable-resolution heat maps and **varyres**
- Heat map confidence intervals? Interactive HMCIs and **mapapp**
- Fitting big data SGLMMs to baseball data?
  - ▶ Stan - insufficient
  - ▶ PPM - Slow and did not converge; to be continued...
  - ▶ INLA - Fast, successful; to be continued...

Thanks for listening.

Sudipto Banerjee, Alan E Gelfand, Andrew O Finley, and Huiyan Sang.  
Gaussian predictive process models for large spatial data sets.  
*Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(4):825–848, 2008.

Andrew O Finley, Sudipto Banerjee, and Alan E Gelfand. spbayes for  
large univariate and multivariate point-referenced spatio-temporal  
data models. *arXiv preprint arXiv:1310.8192*, 2013.

Håvard Rue, Sara Martino, and Nicolas Chopin. Approximate bayesian  
inference for latent gaussian models by using integrated nested  
laplace approximations. *Journal of the royal statistical society: Series b (statistical methodology)*, 71(2):319–392, 2009.

# Moving Forward

- Improve packages
- Submit packages to CRAN
- Biomechanists at OSU
- Compare GLM and SGLMM with scoring rules
- Exit velocity as response variable
- Independence assumption, pitch sequences
- Randomized pitch selection