

A Bayesian Joint Model for Spatial Point Process with Application to Basketball Shot Chart

Jieying Jiao

Department of Statistics, University of Connecticut

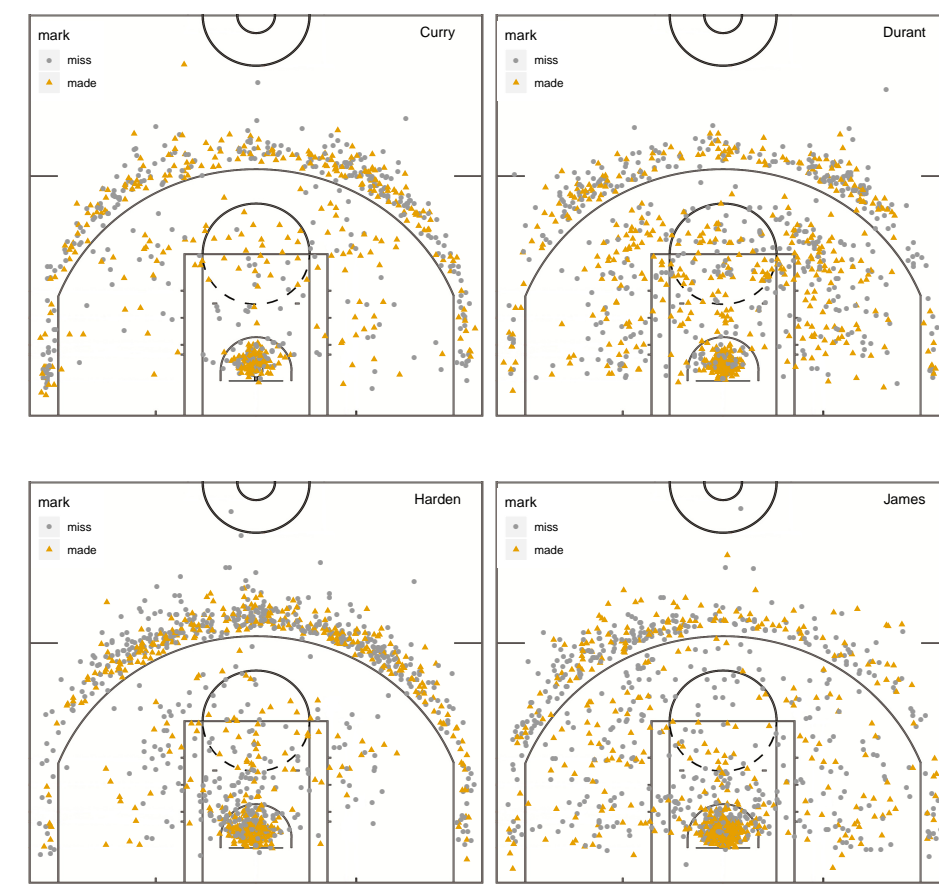
Joint work with Guanyu Hu and Jun Yan

Introduction

- Shot chart are important summaries for basketball coaches. But no work has been done to jointly model shot location and shot success indicator, which is called mark.
- Shot success rate maybe higher at locations where more shots are made.
- We propose a Bayesian joint model of marked spatial point process to model shot location and mark simultaneously. And use spike-slab prior to do variable selection.

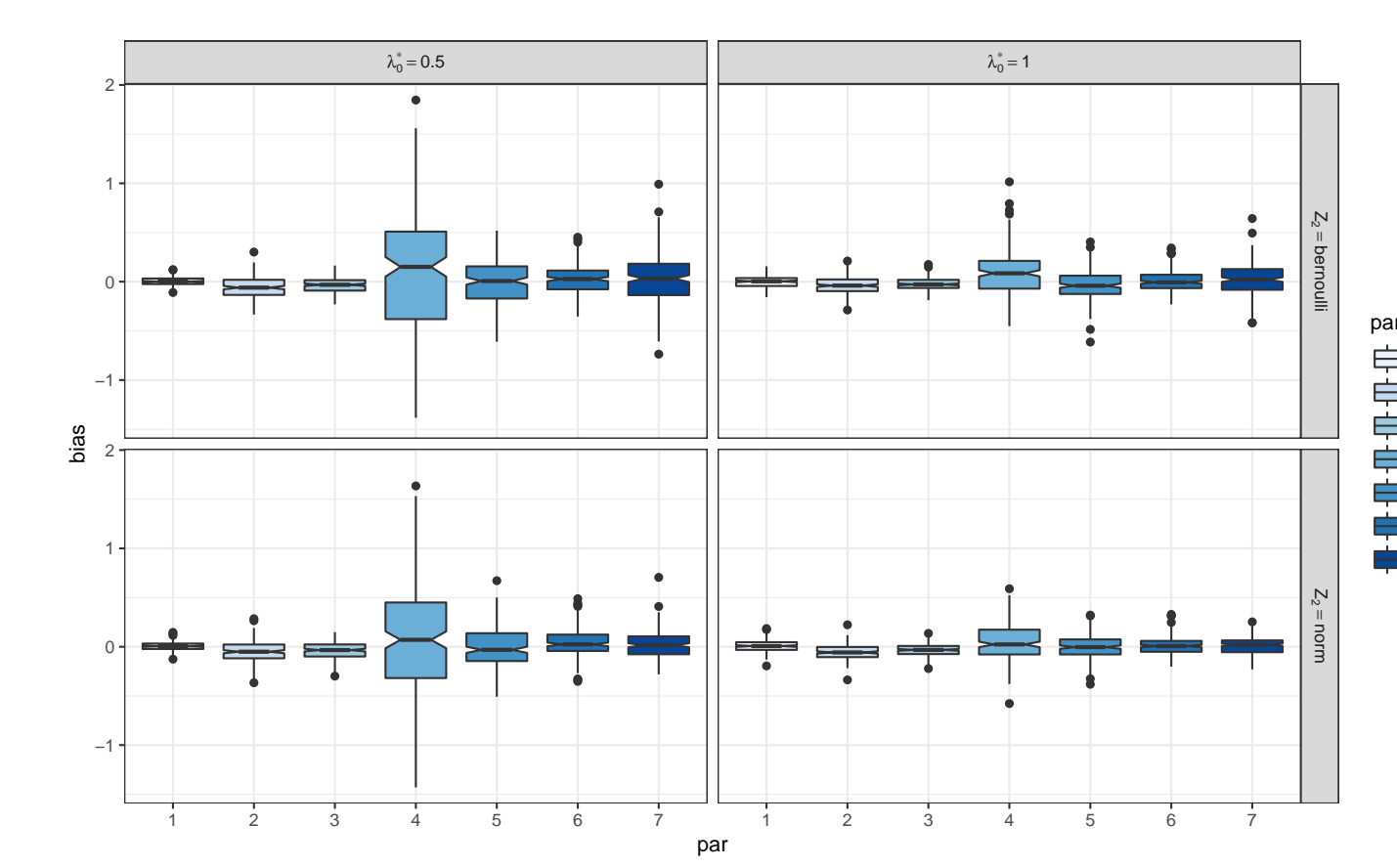
Data Overview

NBA 2018	player regular	shot	chat	data	of	2017-
Player Name	Shot Number	Success rate (%)	2-point shot(%)	Period(%)		
Curry	740	50.0	43.4	(35.0, 20.7, 34.1,	9.9, 0.4)	
Durant	1032	52.5	66.7	(30.9, 23.6, 30.4, 14.7,	0.3)	
Harden	1286	45.6	50.6	(28.7, 22.3, 27.7, 21.0,	0.3)	
James	1409	54.3	74.6	(29.1, 21.5, 25.3, 23.6,	0.4)	



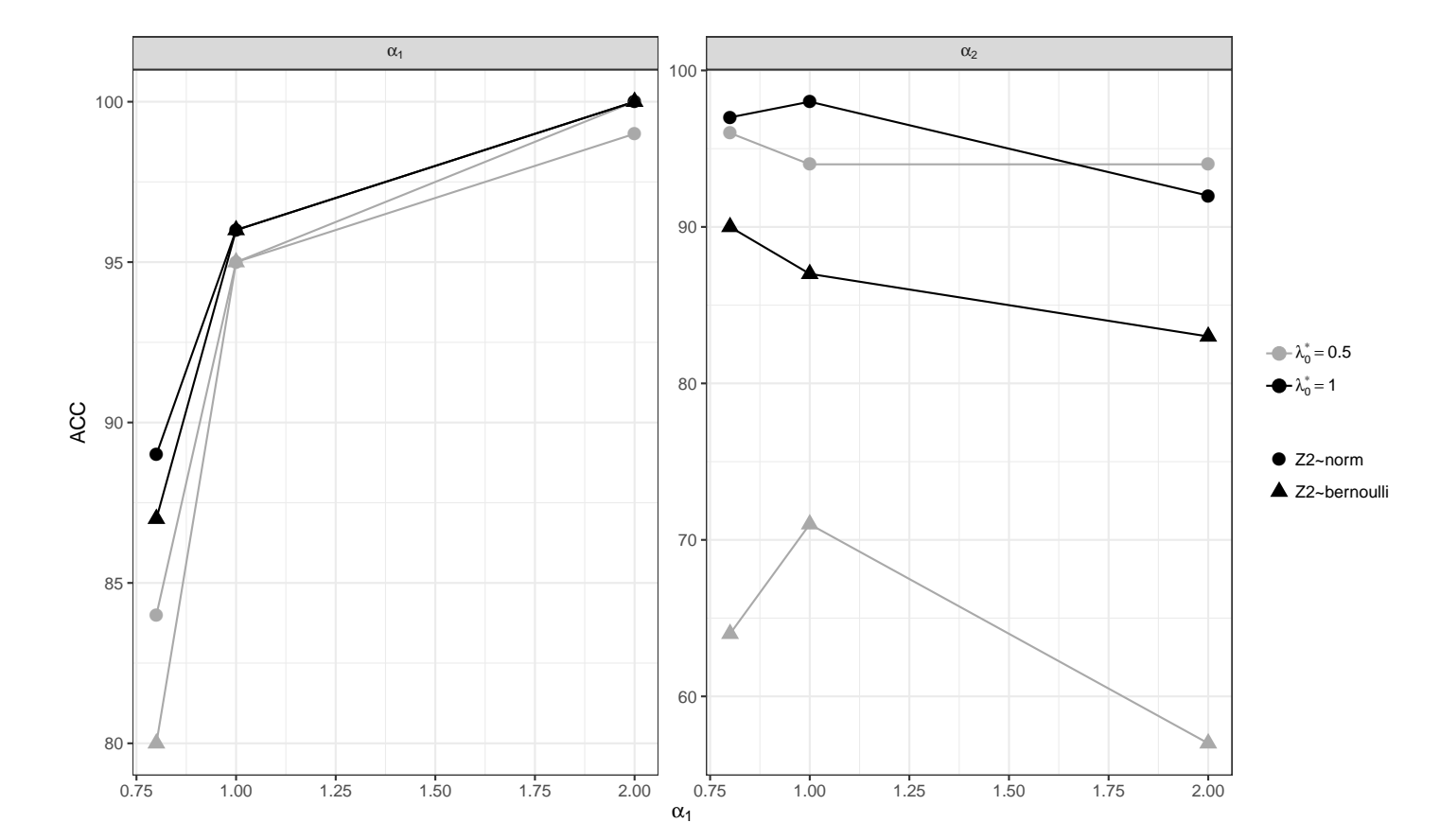
Simulation–Estimation

- On area $\mathcal{B} = [-1, 1] \times [-1, 1]$, $\mathbf{X}^\top(\mathbf{s}_i) = (x_i, y_i)$, $\mathbf{Z}^\top(\mathbf{s}_i) = (1, Z_{1i}, Z_{2i})$ from $N(0, 1)$ or $Bern(0.5)$, $\beta^\top = (2, 1)$, $\xi = \alpha_0 = 0.5$, $\alpha_2 = 1$, $\lambda_0 = 100\lambda_0^*$, $\lambda_0^* \in \{0.5, 1\}$, $\alpha_1 \in \{0.8, 1, 2\}$. No variable selection and vague prior was used. 20,000 MCMC chain with 10,000 burn-in.
- No obvious trend in estimation performance when α_1 increases. Plot below shows the results when $\alpha_1 = 1$.



Simulation–Variable Selection

- New covariates $Z_3, \dots, Z_4 \sim N(0, 1)$ in marks model, $\alpha_3, \dots, \alpha_4$ equals 1 or 0, number of significant Z'_i s are $\{2, 4, 6\}$ (Spike-Slab used). Others unchanged.
- Results shows that the selection accuracy rates for those zero valued coefficients are always 100%.
- Plot below shows the results when only α_1, α_2 are non-zero.



Model Setup

- With location data $\mathbf{S} = (s_1, \dots, s_N)$ and associated mark data $\mathbf{M} = (m(s_1), \dots, m(s_N))$, we propose to use non-homogeneous Poisson point process (Diggle, 2013) and logistic regression, respectively. Intensity dependent model (Ho and Stoyan, 2008) is incorporated for the joint model.
- Joint Likelihood: $L(\Theta | \mathbf{S}, \mathbf{M}) \propto L(\lambda_0, \beta | \mathbf{S}) \times L(\xi, \alpha | \mathbf{M}, \lambda)$, $\Theta = (\lambda_0, \beta, \xi, \alpha_0)$.

Intensity Model:

Mark Model:

$$L(\lambda_0, \beta | \mathbf{S}) \propto \left(\prod_{i=1}^N \lambda(s_i) \right) \exp \left(- \int_{\mathcal{B}} \lambda(s) ds \right),$$

$$\lambda(s_i) = \lambda_0 \exp(\mathbf{X}^\top(s_i)\beta),$$

$$L(\xi, \alpha | \mathbf{M}, \lambda) = \prod_{i=1}^N \theta(s_i)^{m(s_i)} (1 - \theta(s_i))^{1-m(s_i)},$$

$$\text{logit}(\theta(s_i)) = \xi \lambda(s_i) + \mathbf{Z}^\top(s_i)\alpha,$$

where λ_0 is the baseline intensity, $\mathbf{X}(s_i)$ is a $p \times 1$ spatially varying covariate vector, and β is the corresponding coefficient vector.

where $\lambda(s_i)$ is the intensity defined above with a scalar coefficient ξ , $\mathbf{Z}(s_i)$ is a $q \times 1$ spatially varying covariate vector, and α is a $q \times 1$ vector of coefficient.

Bayesian Inference

- Prior Specification: Conjugate prior for λ_0 and vague prior for other parameters.
- Variable Selection: Spike-Slab prior on covariates coefficients.

$$\lambda_0 \sim G(a, b), \quad \beta \sim \text{MVN}(\mathbf{0}, \sigma^2 \mathbf{I}_p), \quad \alpha_i \sim N(0, \delta_i^2), \quad \delta_i^2 = 0.01(1 - \gamma_i) + 100\gamma_i,$$

$$\xi \sim N(0, \delta^2), \quad \alpha \sim \text{MVN}(\mathbf{0}, \delta^2 \mathbf{I}_q), \quad \gamma_i \sim \text{Bernoulli}(\phi_i), \quad \phi_i \sim \text{Beta}(0.5, 0.5).$$

- MCMC Sampling Schemes: Metropolis-Hasting (MH) within Gibbs algorithm (Using R package nimble).

Sampler: `RW()` for ϕ_i 's, `binary()` for γ_i 's, and `RW_llFunction()` for others.

- Bayesian Model Selection Criteria: mDIC and mLPML on mark model only (Ma, Chen and Hu, 2018) to test intensity independent or dependent model ($\xi = 0$ or not).

Real Data Analysis

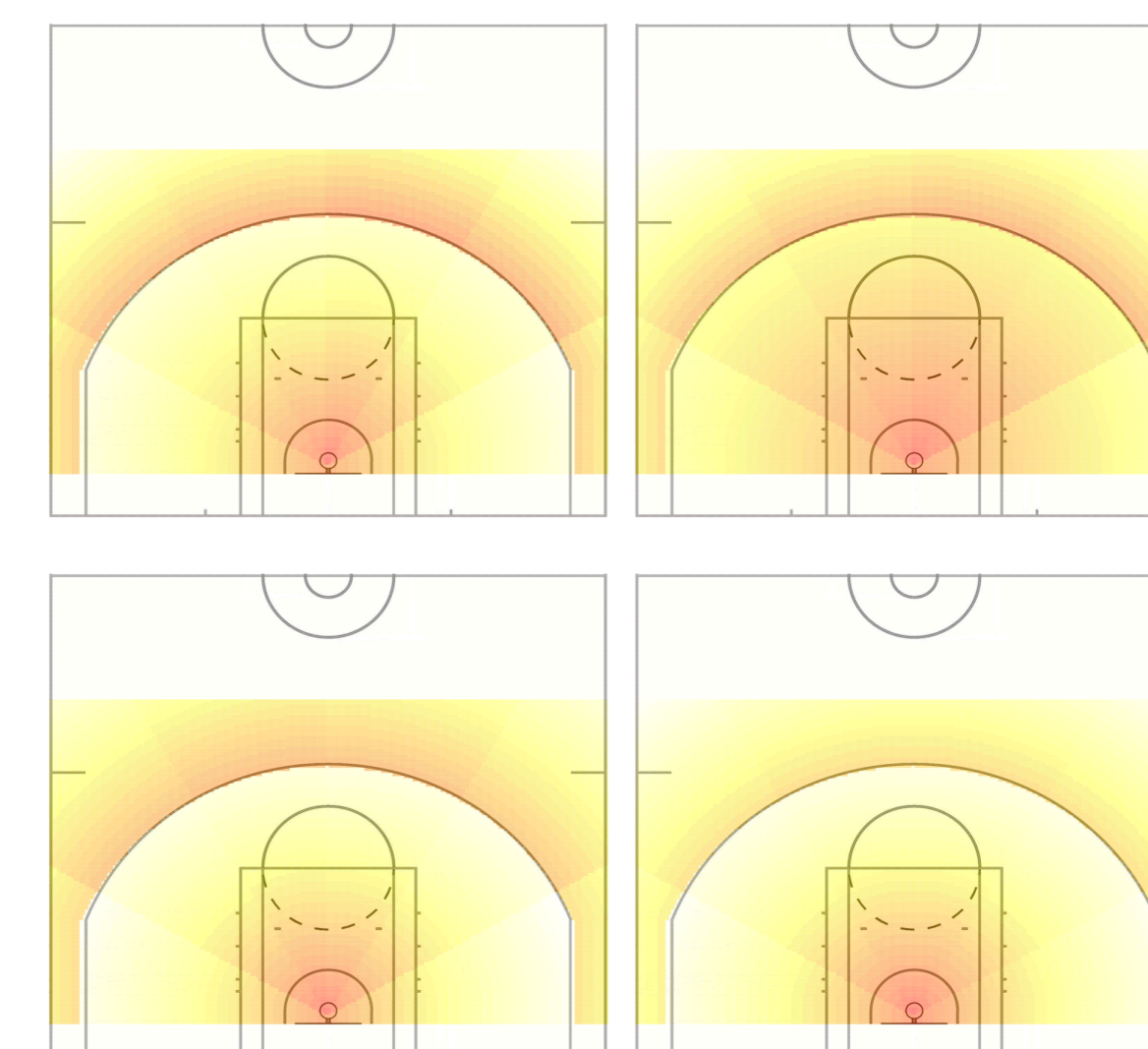
Intensity Model Covariates:

Covariates	Comment
3-point shot	binary
distance1	standardized, distance to origin for 2-point shot, 0 for 3-point shot
distance2	standardized, 0 for 2-point shot, distance to 3-point line for 3-point shot
angle	reference level is $[-\frac{\pi}{2}, \frac{\pi}{6}]$

Mark Model Covariates:

Covariates	Comment
intensity value	left towards the end of the period, scaled with 100
3-point shot	binary
period	reference level is first period
opponent	1 for playoff teams and 0 for others
distance to origin	standardized
angle	reference level is $[-\frac{\pi}{2}, \frac{\pi}{6}]$

- Fitted intensity heat plot for four players and parameters estimation results for Curry's data.



Model	Covariate	Posterior Mean	Posterior SD	95% Credible Interval
Intensity	baseline (λ_0)	0.74	0.06	(0.62, 0.87)
	3-point shot	-3.15	0.22	(-3.58, -2.74)
	distance1	-2.10	0.05	(-2.20, -2.01)
	distance2	-2.32	0.17	(-2.67, -1.99)
	angle $[\pi/6, \pi/3]$	0.42	0.09	(0.23, 0.60)
	angle $[\pi/3, \pi/2]$	0.58	0.09	(0.40, 0.75)
	angle $[\pi/2, 2\pi/3]$	0.58	0.09	(0.40, 0.76)
Mark	angle $[2\pi/3, 5\pi/6]$	0.41	0.09	(0.23, 0.60)
	angle $[5\pi/6, 3\pi/2]$	-0.22	0.11	(-0.43, 0.00)
	Intercept	-0.75	0.15	(-1.04, -0.46)
	λ	0.09	0.01	(0.06, 0.11)

Future Work

- Interaction among different players.
- Multivariate model for more players jointly.
- Non-parametric method for intensity model.
- DIC and LPML for joint model.