



---

# Bayesian spatial scan statistic adjusted for over-dispersion

Deepak Agarwal  
Yahoo! Research  
February 8<sup>th</sup>, 2007  
QMDNS, Virginia.



# Outline

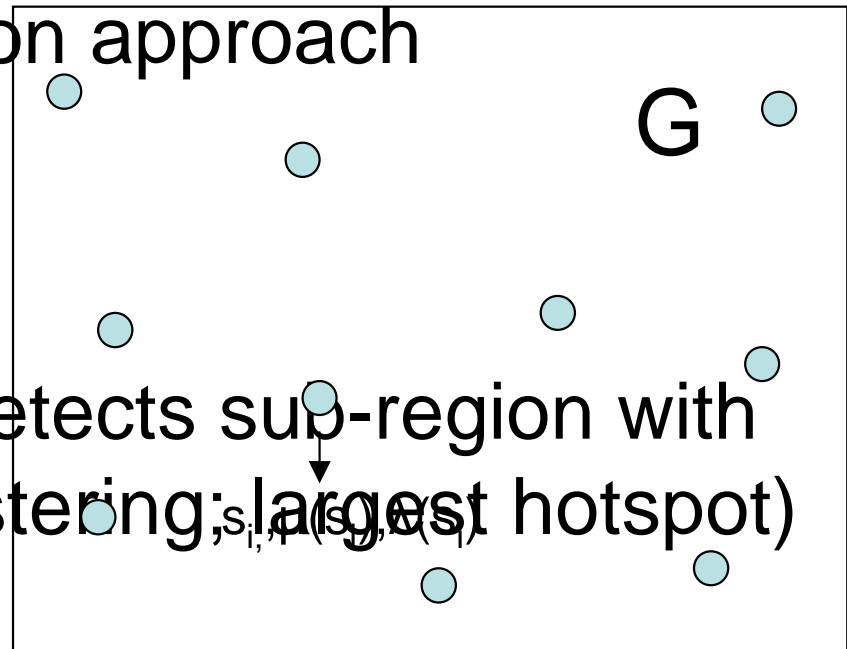
---

- Introduction: Spatial Scan statistic
- Bayesian formulation
- Adjusting for over-dispersion
  - Doubly stochastic/Cox process
- Data analysis: SIDs data
- Ongoing work
  - Incorporating spatial correlation



# Background

- $\mu(s)=v(s)$ ;  $\lambda(s)=\lambda \rightarrow$  homogeneous Poisson
- Global tests for spatial homogeneity
  - E.g. Ripley's K-function approach



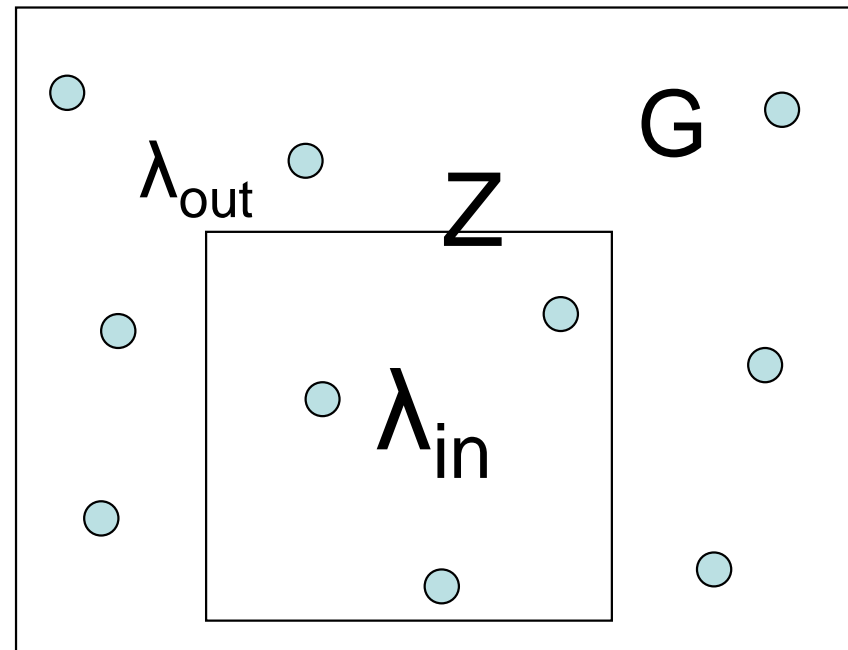
- Scan statistic:
  - Additional nuance (Detects sub-region with maximum spatial clustering; largest hotspot)



# Spatial scan statistic

- Geometric shape  $\rightarrow$  sub regions
  - Rectangle, Square, Circle, Ellipsoid,...
- Scan  $G$  using all possible sub regions

- LLRT: inside-outside
- Max discrepancy region
- Rand test  $\rightarrow$  p-value
- Inhomogeneity
  - Incorporate in  $\mu(s)$





# Model: Poisson case

---

$N$  points  $\{x_i\}$ : inhomogeneous Poisson process

$\mu(x)$ : Baseline measure;  $\lambda(x)$ : Intensity function

$$L(\lambda; \{x_i\}) \propto \left( \prod_i \lambda(x_i) \mu(x_i) \right) \exp\left(-\int_G \lambda(x) d\mu(x)\right)$$

Lattice Data :  $K$  regions;  $n_i$  = number of points in region  $i$

$x_i$  : Centroid of region  $i$

$\mu(x_i)$  : Adjusts inhomogeneity ; e.g., Population at risk

$$L(\lambda; \{x_i\}) \propto \left( \prod_i (\lambda(x_i) \mu(x_i))^{n_i} \right) \exp\left(-\sum_{i=1}^K \lambda(x_i) \mu(x_i)\right)$$



# Scan statistic: Max log(LRT)

**For region Z**

$$\lambda(x) = \lambda_{in} \quad x \in Z; \quad \lambda(x) = \lambda_{out} \quad x \in G - Z$$

$$Z_0 : \lambda_{in} = \lambda_{out} \quad \text{vs} \quad Z : \lambda_{in} > \lambda_{out}$$

**log(LRT) :**

$$D_Z = (n_Z \log(n_Z / \mu_Z) + n_{G-Z} \log(n_{G-Z} / \mu_{G-Z}) - n_G \log(n_G / \mu_G)) \mathbb{1}\left(\frac{n_Z}{\mu_Z} > \frac{n_{G-Z}}{\mu_{G-Z}}\right)$$

**Statistical discrepancy**

KL:  $(n_Z / N, n_{G-Z} / N)$  and  $(\mu_Z / \mu_G, \mu_{G-Z} / \mu_G)$  (Agarwal et al, 2006)

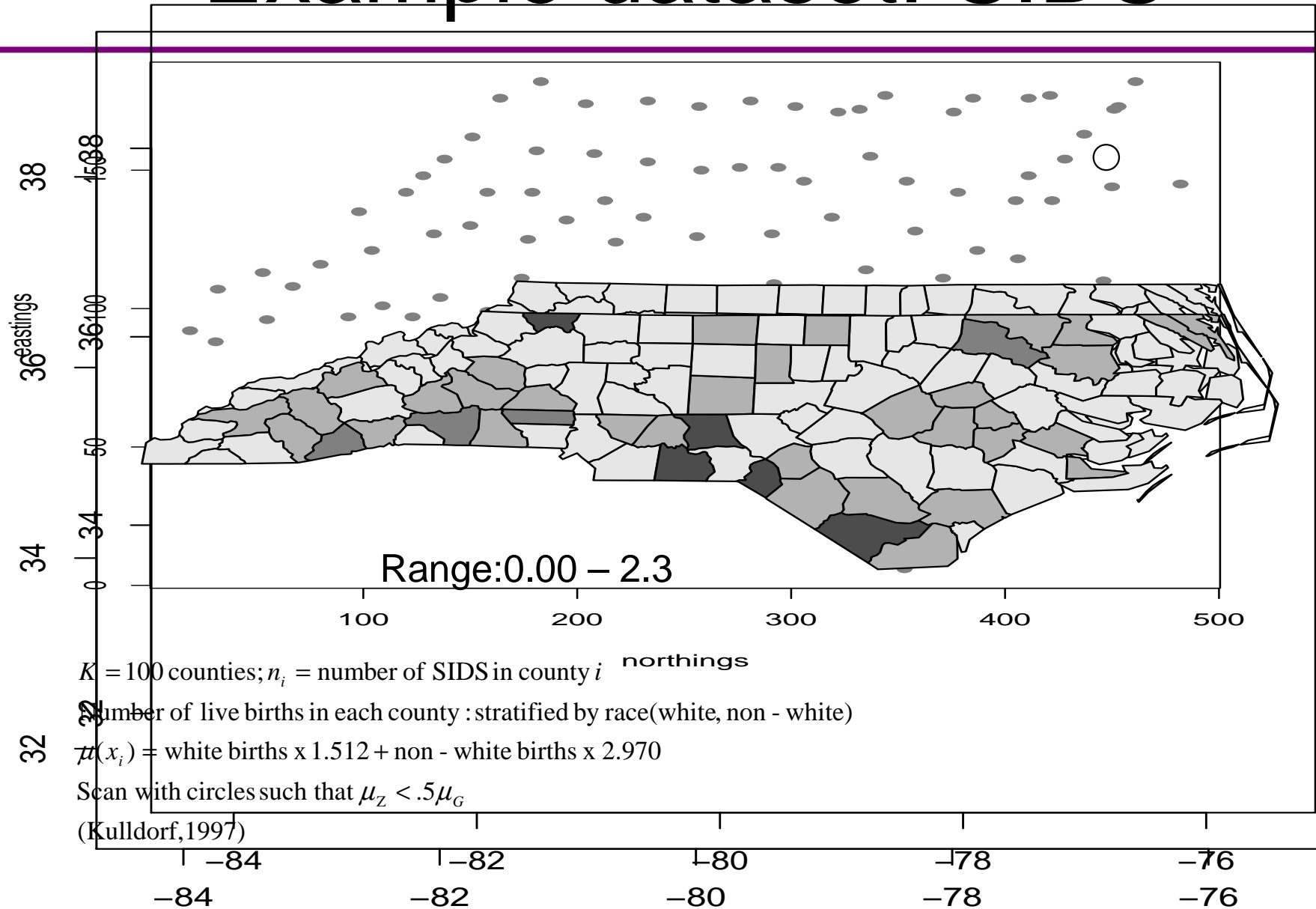
Max discrepancy region :  $Z^* = \arg \max_Z D_Z$

**Significance testing :**

p - value based on Monte - Carlo simulations (expensive)



# Example dataset: SIDS





# Bayesian approach

---

Model space:  $M = \{Z_0\} \cup \{Z_i : i = 1, \dots, N_{reg}\}$

$\pi^M$  : Prior distribution on M

$f(D | Z_i)$  = Marginal likelihood under  $Z_i$

Posterior :  $f(Z_i | D) \propto f(D | Z_i)\pi^M(Z_i)$

$$f(D | Z_i) = \int f(D | Z_i, \theta_{Z_i})\pi(\theta_{Z_i})d\theta_{Z_i}$$

Detect "hotspots" using a threshold on  $f(Z_i | D)$

Threshold : Expert decides; decision theoretic approach





# Poisson case

$$Z_i : \theta_{Z_i} = (\theta_{in} = \log(\lambda_{in}), \theta_{out} = \log(\lambda_{out}))$$

$$\pi(\theta_{in}, \theta_{out}) \sim N((\mu_0, \mu_0), (\delta, \delta), \rho = 0) 1(\theta_{in} \geq \theta_{out})$$

$$\delta = 100 * \text{Var}(n_i / \mu_i) * m^2; \text{ where } m = \text{mean}(n_i / \mu_i);$$

(proper but vague)

$$\pi^M(Z | P) = P 1(Z = Z_0) + (1 - P) \sum_{i=1}^{N_{reg}} 1(Z = Z_i) / N_{reg}$$

$f(D | Z)$ : Tierney - Kadane approximation

$$f(P) = (\alpha + 1) P^\alpha \text{ (Scott and Berger, 2006)}$$

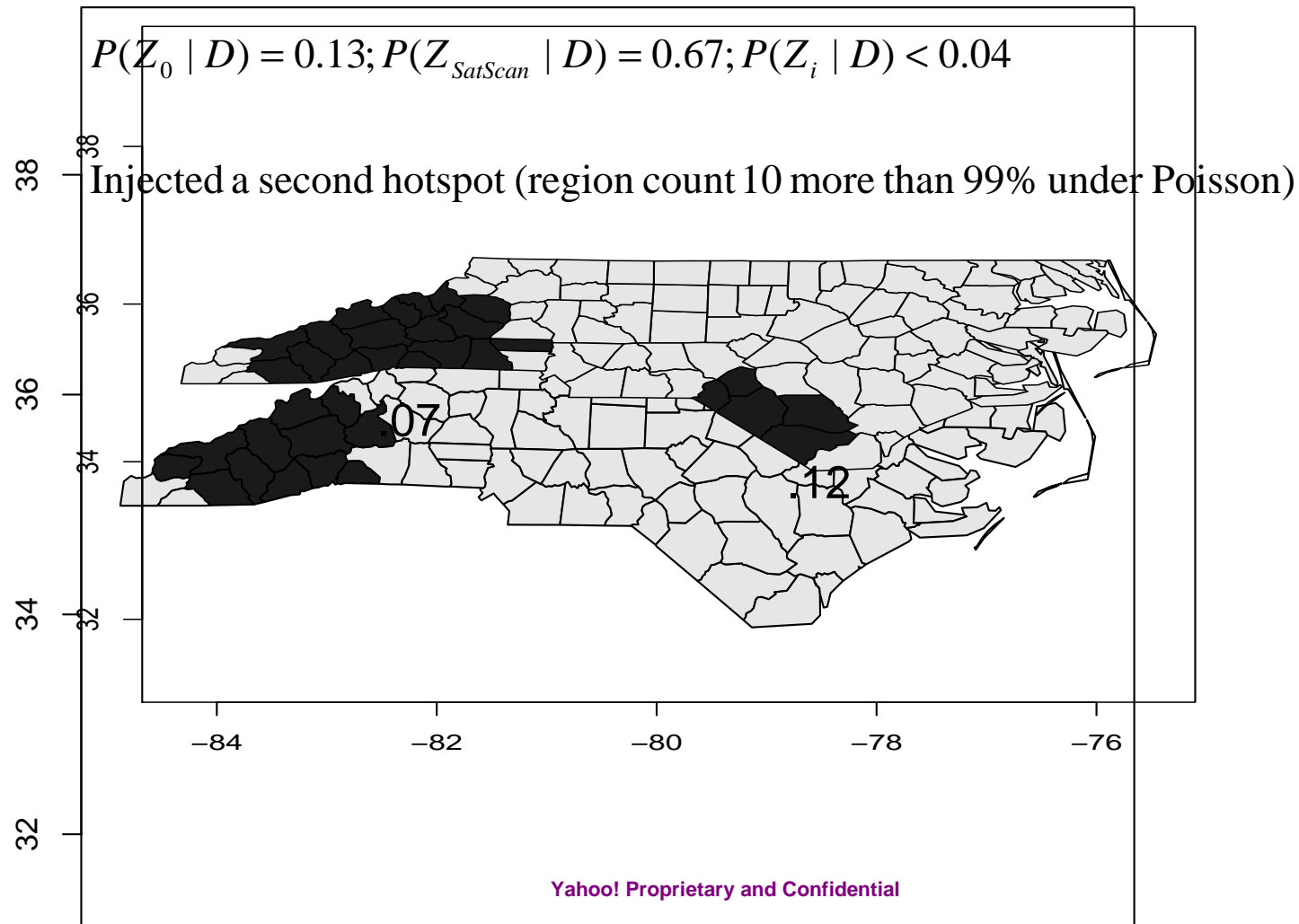
Choose  $\alpha$  such that median = .9



# Analysis for Poisson

Detect one hotspot, exactly the one by SatScan

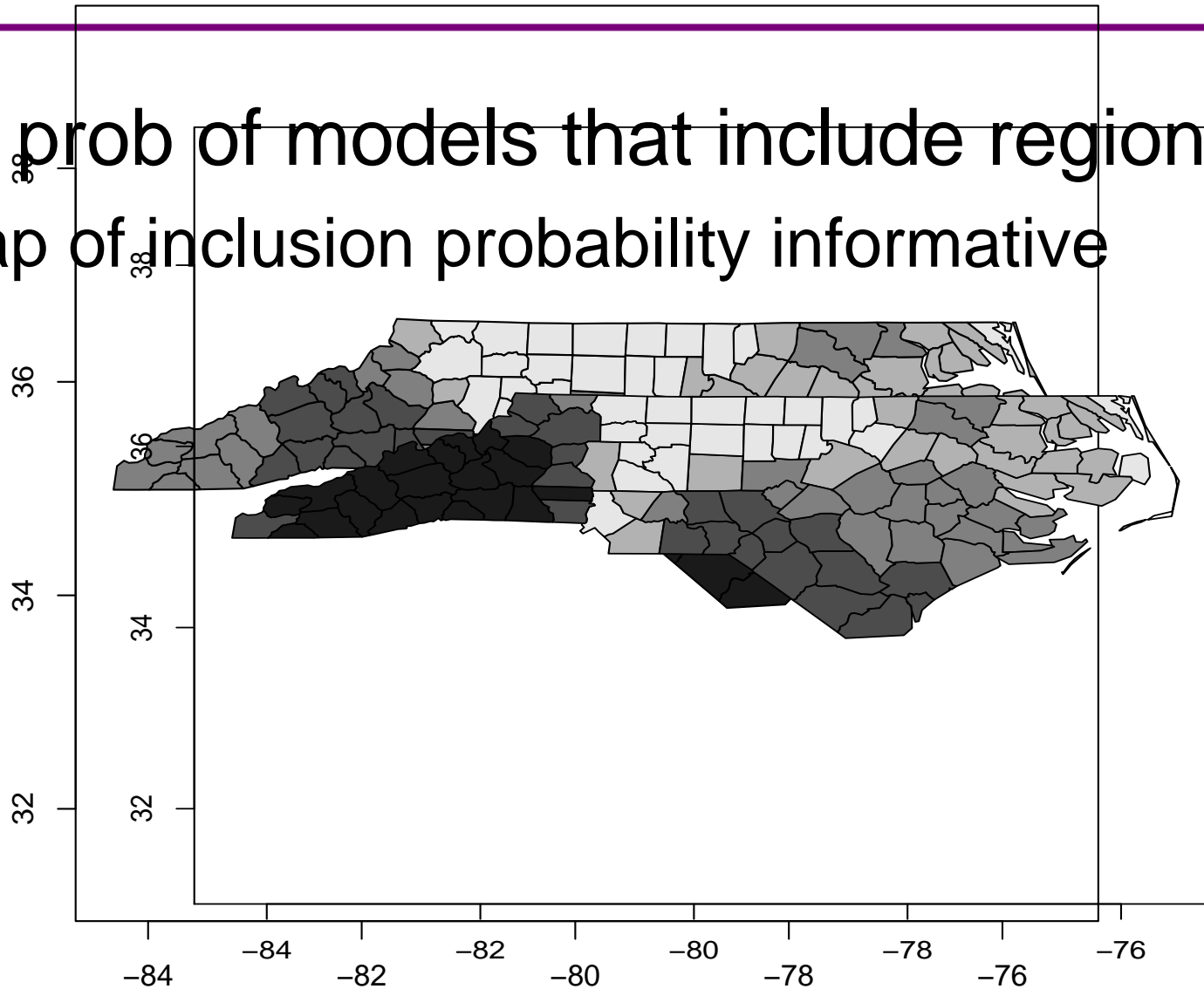
99% for region = 112; observed = 139





# Inclusion probability

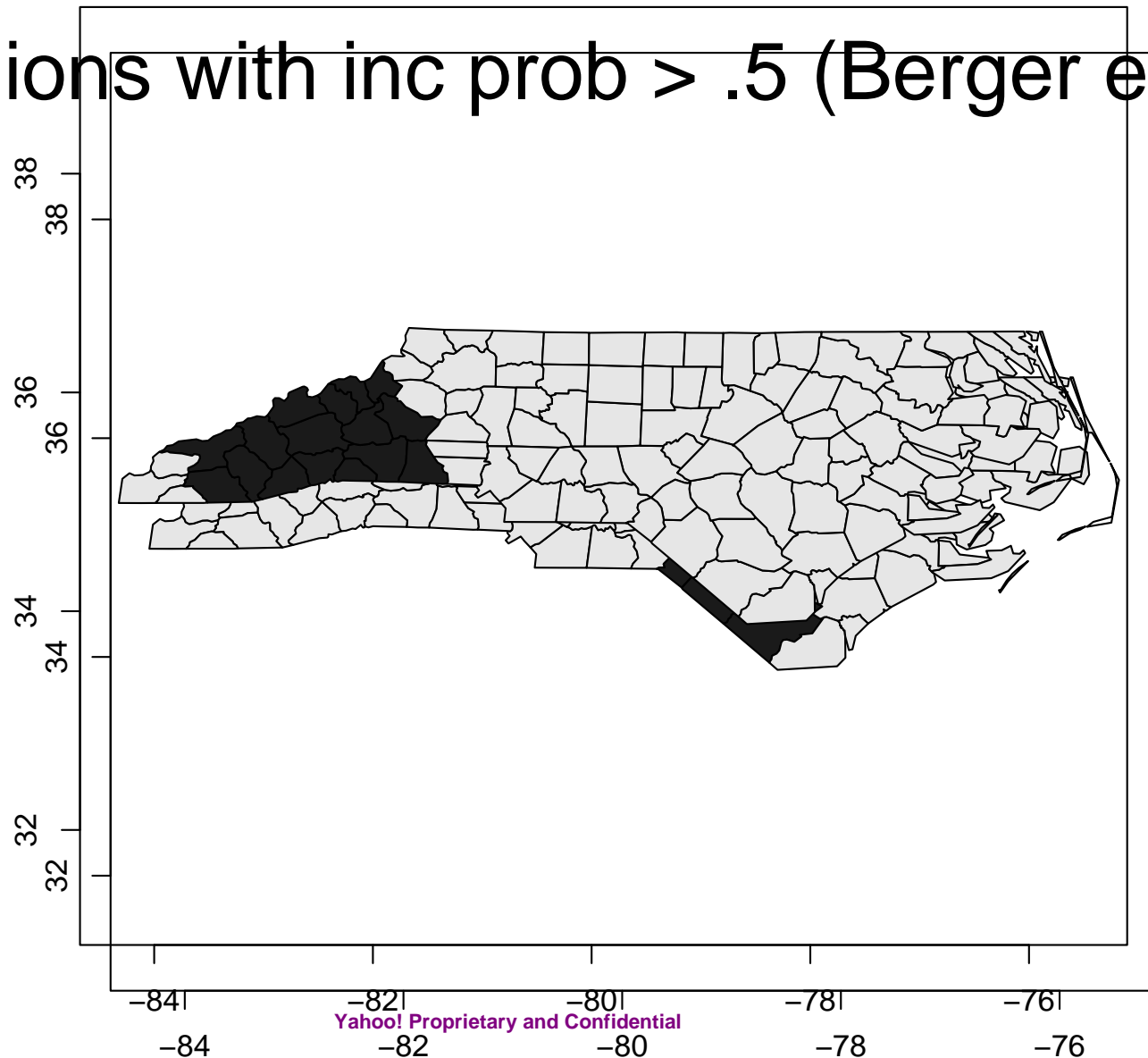
- Post prob of models that include region
  - Map of inclusion probability informative





# Median probability model

- All regions with inc prob  $> .5$  (Berger et al)





# Adjusting for over-dispersion

---

- Over-dispersion in data is routine
  - Adjust to reduce false +ve's
- SIDS data; negative binomial better  
 $\text{mean}(\mu)/\gamma = 15.03/16.44 = .94$

## Cox/doubly stochastic process

Conditional on an error process  $\nu(x)$ , inhomogeneous Poisson

$$\lambda(x_i) = \lambda_{in} \nu(x_i) \quad x_i \in Z$$

$$\lambda(x_i) = \lambda_{out} \nu(x_i) \quad x_i \in G - Z$$

$\nu(x_i)$  iid  $\text{Gamma}(\gamma_Z, \gamma_Z)$  accounts for overdispersion

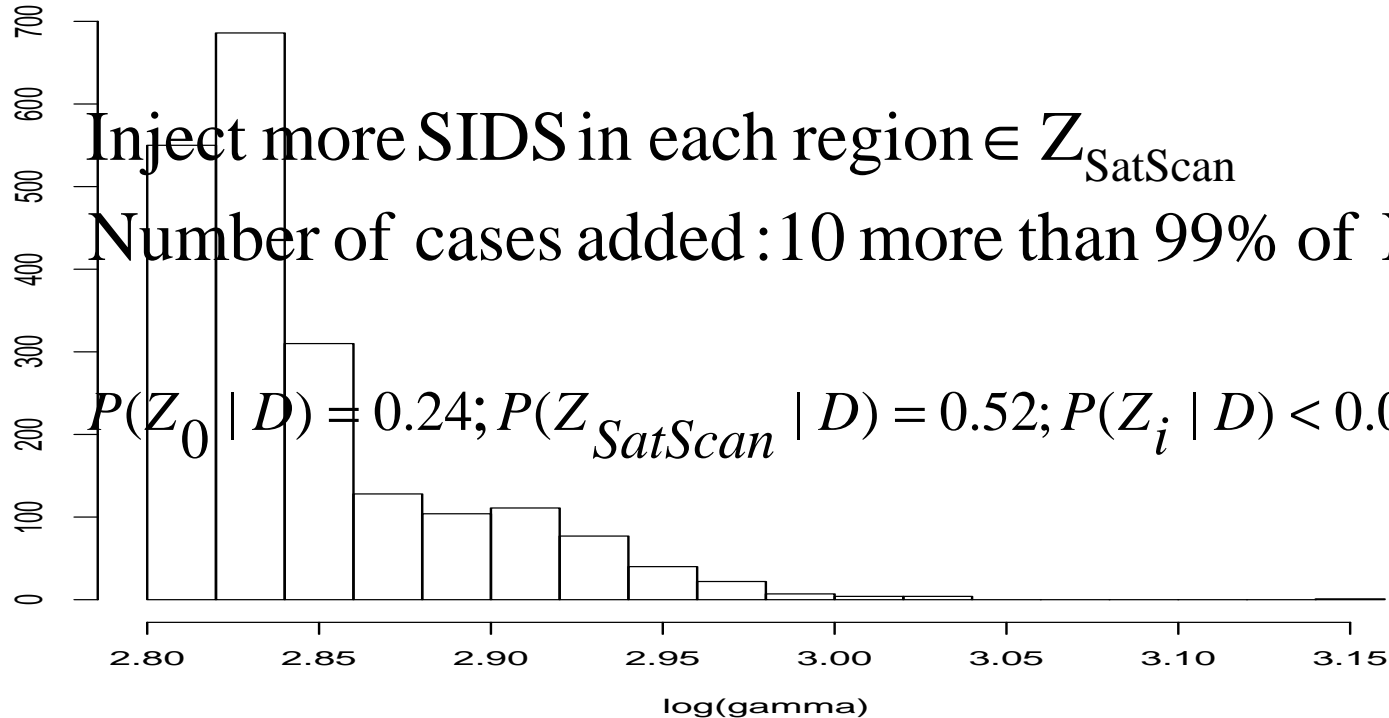


# Results with negative binomial

$$P(Z_0 | D) = 0.84; P(Z_{SatScan} | D) = 0.02; P(Z_i | D) < 0.01$$

Hotspot does not exist anymore!!

99% for SatScan region under NB = 153 > Obs = 139



$$P(Z_0 | D) = 0.24; P(Z_{SatScan} | D) = 0.52; P(Z_i | D) < 0.05$$



# Neil et al (NIPS 2005)

$$n_z \sim NB(\text{mean} = m * n_G / \mu_G; \text{var} = 1 / \mu_z);$$

$$n_{G-z} \sim NB(\text{mean} = n_G / \mu_G; \text{var} = 1 / \mu_{G-z})$$

Biased towards regions with larger measure;

Model	Poisson	NB	NLM
we estimate appropriate parameter from data			
No global model; hard to incorporate other features like spatial correlation. Trivial in our framework			
Null	0.13	0.84	0.74
<b>Fact</b> : The measure for SatScan region 20th percentile			
SatScan	0.67	0.02	0.07
penultimate	0.04	0.01	0.04



# Ongoing Work

---

- Incorporating spatial association
  - Put a spatial prior on the error process
- Loss functions to rank discrepant regions
  - Important for multiple hotspots
- Large scale power tests