

Some Properties of Spatial Scan Statistic Bernoulli Model: Example Simulation for Small and Large Data Using SaTScan

Titin Siswantining^{1*}, Asep Saefuddin², A.N.Khairil², N. Nunung², M. Wayan².

¹Department of Mathematics, Faculty of Mathematics and Natural Science, University of Indonesia, Depok 16424, Indonesia. titin@sci.ui.ac.id

²Department of Statistics, Faculty of Mathematics and Natural Science, Bogor Agricultural University, Bogor, Indonesia

Abstract: Spatial Scan Statistic using SaTScan is used to detect randomness in the cluster. This paper aims to demonstrate some of the properties owned by the spatial scan statistic of Bernoulli model, which exhibit analytically unbiased, have minimum variance and consistent. Simulations using the SaTScan showed that the amount of data give different results on the cluster properties of randomness. Analysis with small data provides a more random cluster, or unfavorable results of the analysis. It is showed that the greater amount of data provide the LLR enlarge; p-value smaller; RR decreases and Smaller biased.

Key words: Spatial Scan Statistic, Bernoulli model, SaTScan, Unbiased, minimum variance, consistent.

I. Introduction

Spatial scan statistic is used to detect clusters in the point process [1]. Suppose N is a spatial point process with $N(A)$ is a random number of points in set $A \subset G$. Moving windows in the study area is defined as a set of zones Z , $Z \subset G$. Z are used interchangeably to mark the subset G and the set of parameters defined in the zone. The spatial scan statistic utilize a test statistic based on the ratio maximum likelihood [1,2,3].

Spatial scan statistic has two models are the Poisson and Bernoulli. This paper discusses only a Bernoulli model. In a Bernoulli model, considering only the size of N such that $N(A)$ is an integer for all subsets $A \subset G$. Each unit size in accordance with the entity or individuals who have a statement of conditions, e.g. with or without disease, or included in a particular species or not, poor or not. It is expressed as individual dots, and the location of individuals at the point. There is exactly one zone $Z \subset G$ such that each individual has a chance of p in the zone, while the chances of individuals outside the zone is q , in the model. Null hypothesis $H_0: p=q$, and the alternative hypothesis $H_1: p>q$, $Z \in \mathcal{Z}$. Under the null hypothesis H_0 , $N(A) \sim \text{Bin}(M(A), p)$ for all sets or the set A . Under the alternative hypothesis H_1 , $N(A) \sim \text{Bin}(M(A), p)$ for all sets $A \subset G$, and $N(A) \sim \text{Bin}(M(A), q)$ for all subsets $A \subset Z^c$ [1].

Some statistical methods for analyzing cluster of spatial point process describe it merely, means that only can detect the location of the cluster without including inference, or inference without the ability to detect the location of the cluster. Important characteristics of the spatial statistical tests were both, such that if the null hypothesis is rejected in the locations of specific areas that cause rejection.

As described before, it is necessary to see how the statistical properties of the Spatial Scan Statistic Bernoulli model analyzed by SaTScan, especially the statistic derived from the direct estimation (DE), which is usually obtained from the survey. Statistical properties belonging to the unbiased, minimum variance and consistently demonstrated in this paper. Spatial scan statistic, especially the consistent nature of the spatial scan statistic is always obtained when the size of a large example [1,4], while in practice due to technical or economic reason, a large sample size is difficult to obtain, so it needs a way out to overcome this.

1.1. Statistical Properties of Bernoulli Model

Suppose that the random variable has a binomial distribution, has a probability density function (pdf) the following:

$$f(x) = \binom{n}{x} p^x (1-p)^{n-x}, x = 0, 1, \dots, n. \quad (1)$$

According to Hogg and Craig [3], unbiased estimator with minimum variance and a unique (unique MVUE = Minimum Variance Unbiased Estimate) to be searched by using the sufficient statistic. To find an estimator is a sufficient statistic or not, can be seen from the following manner.

Exponential and $\ln f(x)$ is:

$$f(x) = \exp \left[x \ln \frac{p}{1-p} + \ln \binom{n}{x} + n \ln(1-p) \right] \tag{2}$$

General form of the exponential class is:

$$f(x) = \exp [a(p)K(x) + S(x) + b(p)]$$

An exponential class of (2), selected with $a(p) = \ln \frac{p}{1-p}$, $K(x) = x$, $S(x) = \ln \binom{n}{x}$, and

$b(p) = n \ln(1-p)$. According to Hogg and Craig [2005], found that probability density function (pdf) of the exponential class meets

- a. the set $\{x | x = 0, 1, \dots, n\}$ does not depend on p ,
- b. $a(p)$ function which is continuous and not constant at $0 < p < 1$,
- c. $K(x)$ function is not constant at $x = 0, 1, \dots, n$ probability

Because it satisfies the above three properties, the density function (pdf) satisfies to the regular case of the exponential class. According to Hogg and Craig [5], therefore pdfs satisfy to the regular of exponential class case, so $K(X) = X$ is complete and sufficient statistic for p .

However, please note that the

$$E(X) = np \text{ is bias.}$$

With algebraic manipulation, is selected such that

$$E\left(\frac{X}{n}\right) = \frac{1}{n} E(X) = \frac{1}{n} np = p. \text{ Proved.}$$

It is mean that $\frac{X}{n}$ is an unbiased estimator for p . Unbiased properties $\frac{X}{n}$ was proven.

Because X is complete sufficient statistic, as well as $\frac{X}{n}$ is an unbiased for p , then the Lehmann-

Scheffe theorem [4], ensured that $\frac{X}{n}$ is an unbiased estimator and variance the minimum as well as unique to p

Lehmann and Scheffe Theorem.

Supposed that X_1, \dots, X_n , n is an integer constant, which have a random sample pdf or PMF $f(x; \theta)$, $\theta \in \Omega$. For example $Y_1 = u_1(X_1, \dots, X_n)$ is a sufficient statistic for θ , and suppose that a complete family $\{f_{Y_1}(y_1; \theta); \theta \in \Omega\}$. And suppose that a complete family Y_1 is a unique MVUE of θ .

Furthermore, note that

$$\text{var}\left(\frac{X}{n}\right) = \frac{1}{n^2} \text{var}(X) = \frac{p(1-p)}{n} \tag{3}$$

Thus, the Chebyshev's inequality is obtained that

$$\Pr \left[\left| \frac{X}{n} - p \right| < k \sqrt{\frac{p(1-p)}{n}} \right] \geq 1 - \frac{1}{k^2} \tag{4}$$

Chosed that $k = \varepsilon \sqrt{\frac{n}{p(1-p)}}$ for any $\varepsilon > 0$, inequality (4) is transformed into

$$\Pr \left(\left| \frac{X}{n} - p \right| < \varepsilon \right) \geq 1 - \frac{p(1-p)}{\varepsilon^2 n} \tag{5}$$

Taking the limit in inequality (5),

$$\lim_{n \rightarrow \infty} \Pr \left(\left| \frac{X}{n} - p \right| < \varepsilon \right) \geq 1 - \lim_{n \rightarrow \infty} \frac{p(1-p)}{\varepsilon^2 n} = 1 \tag{6}$$

However, because the probably is not greater than one, then

$$\lim_{n \rightarrow \infty} \Pr \left(\left| \frac{X}{n} - p \right| < \varepsilon \right) = 1. \tag{7}$$

Found that $\frac{X}{n}$ consistent estimator for p . Thus, if X has a binomial distribution, so we obtained conclusion that $\frac{X}{n}$ is an unbiased estimator with minimum variance, unique and consistent for p .

1.2. Statistical properties of SaTScan

In the case of scan statistic Bernoulli model, the random variable is $n(G)$ and sample size is $N(G)$, so the conditions are right ($p=q$), then its pdf is

$$f(n_G) = \binom{N(G)}{n_G} p^{n_G} (1-p)^{N(G)-n_G}.$$

With the previous explanation, and in a similar way, it was found that $\frac{n_G}{N(G)}$ is an unbiased estimator with minimum variance, unique and consistent for p . If n_z is the number of observation points in the zone Z and if H_1 so, with a similar explanation, it is obtained that $\frac{n_z}{N(Z)}$ is an unbiased estimator with minimum variance, unique and consistent for p , and $\frac{n_G - n_z}{N(G) - N(Z)}$ is an unbiased estimator with minimum variance, unique and consistent for q .

Noted that if $H_0 (p=q)$ true, then $\frac{n_z}{N(Z)}$ and $\frac{n_G - n_z}{N(G) - N(Z)}$ would be consistent to the same value, p .

Therefore, the triangle inequality,

$$\begin{aligned} \varepsilon &\leq \left| \frac{n_z}{N(Z)} - \frac{n_G - n_z}{N(G) - N(Z)} \right| \\ &\leq \left| \frac{n_z}{N(Z)} - p + p - \frac{n_G - n_z}{N(G) - N(Z)} \right| \\ &\leq \left| \frac{n_z}{N(Z)} - p \right| + \left| \frac{n_G - n_z}{N(G) - N(Z)} - p \right| \end{aligned}$$

So that

$$\begin{aligned} \Pr \left[\left| \frac{n_z}{N(Z)} - \frac{n_G - n_z}{N(G) - N(Z)} \right| \geq \varepsilon \right] &\leq \Pr \left[\left| \frac{n_z}{N(Z)} - p \right| + \left| \frac{n_G - n_z}{N(G) - N(Z)} - p \right| \geq \varepsilon \right] \\ &\leq \Pr \left[\left| \frac{n_z}{N(Z)} - p \right| \geq \frac{\varepsilon}{2} \text{ or } \left| \frac{n_G - n_z}{N(G) - N(Z)} - p \right| \geq \frac{\varepsilon}{2} \right] \\ &\leq \Pr \left[\left| \frac{n_z}{N(Z)} - p \right| \geq \frac{\varepsilon}{2} \right] + \Pr \left[\left| \frac{n_G - n_z}{N(G) - N(Z)} - p \right| \geq \frac{\varepsilon}{2} \right] \end{aligned}$$

The result obtained

$$\begin{aligned} \lim_{n \rightarrow \infty} \Pr \left[\left| \frac{n_z}{N(Z)} - \frac{n_G - n_z}{N(G) - N(Z)} \right| \geq \varepsilon \right] \\ \leq \lim_{n \rightarrow \infty} \Pr \left[\left| \frac{n_z}{N(Z)} - p \right| \geq \frac{\varepsilon}{2} \right] + \lim_{n \rightarrow \infty} \Pr \left[\left| \frac{n_G - n_z}{N(G) - N(Z)} - p \right| \geq \frac{\varepsilon}{2} \right] = 0 \end{aligned}$$

However, because the chances of non-negative value, then

$$\lim_{n \rightarrow \infty} \Pr \left[\left| \frac{n_z}{N(Z)} - \frac{n_G - n_z}{N(G) - N(Z)} \right| \geq \varepsilon \right] = 0.$$

Or simply, the values of $\frac{n_z}{N(Z)}$ and $\frac{n_G - n_z}{N(G) - N(Z)}$ much closed to consistency for a large samples.

From this description shows that it is true that the scan statistic have three properties: unbiased, minimum variance, and consistent. Spatial scan statistic requires large sample sizes especially for the consistently properties.

II. Simulation Data By Spatial Scan Statistic

To see how the statistical properties owned by the scan statistic, then performed a data simulation scan statistic and case study analysis. Simulations are carried out by Ugarte [6], which is simulated by dividing the area into 5 zones which percentile to the 10th, 25th, 50th, 75th, and percentile to the 90th. However, in this paper, from 35 area or village observed, divided into four zones, namely $p = 0.2; 0.4; 0.8$ and $p = 0.95$. At each village the size of the example let us assume that at 16, so that the paper is distributed Binomial many cases by the number of cases $y_i = n * p_i$ are many cases in the area to- i is the sample size multiplied by the proportion in the i -th village. Results of the spatial scan statistic showed that the proportion of the village with other village is very significant, to the second cluster. It can be seen in Table 1.

Table 1. Summary analyzed result by Scan statistic simulation of 35 villages.

	Cluster 1 (MLC)	Cluster 2
O/E Ratio	1.44	1.44
Risk Ratio (RR)	1.53	1.48
Log Likelihood Ratio (LLR)	16.561325	7.891521
p-value	0.00000034	0.0087
villages	31, 33, 25, 29	22, 30, 32, 35

The observation to the ratio of the expected value is 1.44, means that this estimator have a small bias values. The ratio of observed to expected value equal to 1, means that the estimator be the cluster size of the risk than the risk out of the cluster. The greater value means greater risk in different clusters with the risk outside the cluster. This shows that the estimator is closed to unbiased properties.

The ratio of observed to expected value of the same magnitude, both in cluster 1 (Most Likely Cluster, MLC) and in cluster 2, with a value of RR in MLC was higher than in cluster 2. Very small value of significance it is meaning between the groups with other groups is significantly different proportions.

In order to view further the statistic properties, so we simulate using the large amounts of data. For large amounts of data, simulation results are based on synthetic estimator, i.e. estimator obtained from DE, coupled with the estimation from the data village potentials (Podes) as an auxiliary variable. Simulations summary of the 247 villages synthetic estimator is presented in Table 2. While the comparison of small data with large data, is summarized in Table 3.

Table 2. Simulations summary of the 247 villages synthetic estimator

The synthetic estimator 247 villages	
LLR	413.540690
RR	2.53
O/E ratio	1.53
p-value	< 0.000000000000000010
Cluster 1	223, 224, 222, 213, 221, 214, 187, 220, 212, 168, 169, 186, 225, 208, 167, 166, 211, 215, 183, 165, 216, 182, 163, 184, 209, 176, 188, 164, 175, 218, 181, 173, 219, 172, 243, 185, 162, 217, 180, 244, 210, 174, 179, 239, 246, 177, 242, 178, 161, 247, 245, 171, 198, 238, 205, 204, 241, 61, 199, 158, 240, 170, 237, 206, 194, 13, 235, 59, 197, 160, 193, 157, 234, 154, 236, 232, 195, 60, 62, 200, 231, 155, 230, 233, 85, 203, 53, 58, 189, 228, 63, 159, 229, 156, 227, 151, 57, 84, 192, 226, 153, 82, 201, 70, 207, 202
Cluster 2	136, 137, 138, 147, 146, 149, 134, 140, 135, 145, 139, 141, 144, 133, 150 (O/E = 1.27; RR = 1.3)

Table 3. The simulation summary scan statistic of small and large data

Cluster 1				Cluster 2			
LLR	O/E ratio	RR	p-value	LLR	O/E ratio	RR	p-value

Some properties Of Spatial scan Statistic Bernoulli model: Example simulation For small And Large

Small Data	16.561325	1.44	1.53	0.00000034	7.891521	1.44	1.48	0.0087
Large Data	413.54069	1.53	2.53	$< 10^{-16}$	9.237041	1.27	1.30	0.021

In small data, LLR value is smaller than the large amount of data, so that the p-value becomes small, and RR is large, meaning that the more significant differences between clusters in to a cluster outside. With simulated data, the spatial position cannot be known; so as to simulate the statistical properties cannot be determined.

III. A Case Study On Poverty In Jember Region Indonesia

For data analysis used a case study on poverty in Jember, Indonesia. Summary of the analysis is shown in Table 4.

Table 4. Summary of SaTScan results to DE of 35 villages.

	LLR	RR	O/E ratio	p-value	Villages
Cluster 1	23.332	2.062	1.711	0.001	19, 8, 34, 5, 4, 26, 11
Cluster 2	16.266788	4.738	4.376	0.001	28, 20

The result of the analysis for a small amount of data is still rather difficult to interpret. Therefore, an analysis was performed for a large amount of data, by using a synthetic estimator. Synthetic estimator obtained from estimating DE with information obtained by borrowing strength from other areas. The information used in this issue is the information on Village Potential (Podes) derived from the Central Bureau of Statistic (BPS) Indonesia [7,8,9]. The analyzed result by using SaTScan, spatial scan statistic software, for the 247 villages is shown in Table 5.

If the two tables are compared, it appears that the greater amount of data will provide:

1. LLR enlarge
2. p-value smaller
3. RR decreases
4. Smaller bias.

The greater of LLR value gives a smaller p-value.

By using the case studies, the spatial position can be known with certainty. Therefore, the statistical properties can be determined. It appears that when a large number of data, then the ratio of O/E becomes smaller, close to value 1. This means that unbiased properties can be demonstrated.

Table 5. Summary of SaTScan results for the 247 synthetic estimator

	LLR	RR	O/E ratio	p-value	Regions
Cluster 1	65.647489	1.45	1.19	< 0.000000000000000010	234, 13, 235, 85, 170, 233, 63, 236, 62, 237, 238, 84, 239, 171, 241, 174, 231, 229, 242, 230, 240, 82, 61, 243, 60, 57, 69, 173, 81, 245, 175, 228, 244, 172, 176, 77, 227, 80, 70, 232, 226, 83, 68, 246, 75, 177, 164, 181, 182, 165, 167, 166, 76, 247, 58, 59, 180, 159, 79, 64, 91, 168, 65, 67, 169, 160, 51, 92, 74, 186, 66, 221, 153, 161, 178, 90, 71, 183, 72, 220, 154, 53, 73, 187, 184, 93, 222, 50, 163, 162, 87, 78, 155, 224, 52, 47, 89, 152, 179, 49, 212, 158, 37, 38, 208, 188, 213, 88, 185, 143, 157, 193, 223, 209, 86, 189, 211, 46, 150, 156
Cluster 2	17.698694	1.74	1.72	0.0000087	95, 94, 96, 101
Cluster 3	15.319119	1.92	1.91	0.000076	126, 127

IV. Conclusion And Suggestions

Spatial scan statistic of Bernoulli model have some properties are unbiased, minimum variance and consistent. These properties are equalities that should be owned by an estimator. Based on the analytic seen that

the spatial scan statistic will have these properties only if a large number of samples, particularly on the consistency properties.

Scan statistic requires large data, whereas in reality the available data is small. Analysis with small data provides a more random cluster, or an unfavorable result for the statistical analysis. Therefore, we need other methods that can provide a solution to the problem of small data. One proposed is to use the method of Small Area Estimation (SAE).

SAE has been known to overcome the parameter estimation for small data [10]. SAE is necessary to join into the scan statistic in order to obtain a better estimate, compared with the investigation without joined with SAE, which requires large sample sizes, where large sample sizes are very difficult to find. In theory, SAE has been able to handle the problem of small sample size. According to Rao [10], SAE also has a minimum variance. Therefore, it is possible to replace the role of DE with SAE estimators.

References

- [1]. Kulldorff, M. A. Spatial Scan Statistic. *Commun. Statist.-Theory Meth.* 26(6), 1997, 1481 – 1496.
- [2]. Duczmal, L., M. Kulldorff, L. Huang. Evaluation of spatial scan statistics for irregularly shaped disease clusters. To appear in *Journal of Computation and Graphical Statistics*. 2006.
- [3]. Tango, T. A Spatial Scan Statistic with a Restricted Likelihood Ratio. *Japanese Journal of Biometrics*, Vol. 29, No. 2, 2008, 75 – 95.
- [4]. Cucala, L., C. Demattei, P. Lopes, & A. Ribeiro. Spatial scan statistics for case event data based on connected components. *Biometrics*, 2009, 1 – 17.
- [5]. Hogg, R.V. & A.T. Craig. *Introduction to Mathematical Statistics*. 5th ed. (Prentice Hall, London, 2005)
- [6]. Ugarte, M.D., T. Goicoa, A.F. Militino. Empirical Bayes and Fully Bayes procedures to detect high-risk areas in disease mapping. *Computational Statistics and Data Analysis*, 53, 2009, 2938 – 2949.
- [7]. Badan Pusat Statistik (BPS, Central Bureau of Statistic) Indonesia. *Basic Poverty Measurement and Diagnostics Course*. (BPS Pubs, Jakarta, Indonesia, 2002).
- [8]. Badan Pusat Statistik (BPS, Central Bureau of Statistic) Indonesia. *Data and information poverty*. 2008. (BPS Pubs, Jakarta, Indonesia, 2008).
- [9]. Badan Pusat Statistik (BPS, Central Bureau of Statistic) Indonesia. *Indonesian Statistic 2011*. (BPS Pubs, Jakarta, Indonesia, 2012).
- [10]. Rao, J.N.K. *Small Area Estimation*. (Wiley-Interscience, USA, 2003)