

# CORRECTED MORAN'S I STATISTIC

by  
Jian Ying

A dissertation submitted to the faculty of  
The University of Utah  
in partial fulfillment of the requirements for the degree of

Master of Science  
in  
Computing

School of Computing  
The University of Utah  
March 2019

Copyright © Jian Ying 2019  
All Rights Reserved

The University of Utah Graduate School

STATEMENT OF DISSERTATION APPROVAL

The dissertation of Jian Ying  
has been approved by the following supervisory committee members:

Jeff M. Phillips , Chair(s) March 2019  
Date Approved

Aditya Bhaskara , Member March 2019  
Date Approved

Shandian Zhe , Member March 2019  
Date Approved

by Ross Whitaker , Chair/Dean of  
the Department/College/School of Computing  
and by David B. Kieda , Dean of The Graduate School.

## **ABSTRACT**

Moran's I is a statistic that measures spatial correlation of  $n$  spatial data points. It has been widely used in identifying spatial patterns. In this thesis, a method that estimates Moran's I by sampling and the bounds of error of the estimate is described. A corrected version of Moran's I that reduces bias of the estimate when the associated value can only be measured with noise is proposed and the properties are proved analytically and demonstrated by simulation.

# CONTENTS

<b>ABSTRACT</b> .....	<b>iii</b>
<b>LIST OF FIGURES</b> .....	<b>v</b>
<b>ACKNOWLEDGEMENTS</b> .....	<b>vi</b>
<b>CHAPTERS</b>	
<b>1. INTRODUCTION</b> .....	<b>1</b>
<b>2. ESTIMATE MORAN'S I BY SAMPLING</b> .....	<b>3</b>
2.1 Understanding the Error from Sampling .....	3
2.1.1 Modeling Data from Sampling .....	3
2.1.2 Algorithm to estimate $Z_w$ .....	4
2.1.3 Importance Sampling of Pairs .....	4
2.1.4 Approximation Bounds .....	5
2.1.4.1 Chebyshev's Inequality .....	5
2.1.4.2 Closeness in $\hat{Z}_w$ .....	5
2.1.4.3 Closeness in $\hat{Z}_2$ .....	6
2.1.4.4 Putting them together .....	6
2.1.5 Simulation .....	8
<b>3. CORRECTED MORAN'S I</b> .....	<b>10</b>
3.1 Corrected Moran's I when there is noise .....	10
3.1.1 Approximation Bounds .....	10
3.1.1.1 Closeness in numerators .....	11
3.1.1.2 Closeness in denominators .....	13
3.1.1.3 Putting these together .....	14
3.1.2 Simulations .....	17
3.1.2.1 scenario 1 .....	17
3.1.2.2 scenario 2 .....	17
3.2 Application in Non-Gaussian Distribution .....	19
<b>4. CONCLUSION</b> .....	<b>20</b>
<b>5. REFERENCE</b> .....	<b>21</b>

## LIST OF FIGURES

2.1	The effect of sample size on accuracy of Moran's I estimation . . . . .	9
3.1	The effect of sample size on observed and corrected I in Gaussian distribution	18
3.2	The effect of sample size on observed and corrected I in Bernoulli distribution . . . . .	19

## ACKNOWLEDGEMENTS

I would like to express the deepest appreciation to my committee chair Professor Jeff Phillips. Without his guidance and persistent help this dissertation would not have been possible. It took him a lot of extra effort to guide me, a part-time student without strong background in computer science, to finish the long journey. I really appreciate.

I would like to thank my committee members, Professor Aditya Bhaskara and Professor Shandian Zhe, for the useful comments, remarks and engagement through the learning process of this master thesis.

Finally, I want to thank my parents, Guangtai Ying and Huaizhi Song, and my sister Jing Ying, for their constant support. Special thanks to my wife, Minmin, for her love and care. Last but not least, to my son Kevin and my daughter Katherine, thank you for being good kids and you make my life much more fulfilled.

# CHAPTER 1

## INTRODUCTION

Moran's I is a statistic that measures spatial correlation of  $n$  spatial data points  $P = \{p_1, p_2, \dots, p_n\}$  [1], [2]. Each point  $p_i$  is associated with a value  $x_i \in \mathbb{R}$ . Let  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$  be the mean of all  $x_i$  values, and let  $\tilde{z}_i = x_i - \bar{x}$  be a centred value. Each pair of points  $p_i, p_j \in P$  also have an associated similarity weight  $w_{i,j} \in [0, 1]$ ; we often represent all of these pairs in a symmetric matrix  $W$ . For instance,  $w_{i,j}$  may be 1 if two points  $p_i, p_j$  are within a fixed distance (or the spatial regions they represent are adjacent) and 0 otherwise. Or we may define  $w_{i,j} = \exp(-\|p_i - p_j\|^2/2\sigma^2)$  for spatial points  $p_i, p_j \in \mathbb{R}^d$ . Then the Moran's I is defined as

$$I_{\text{true}} = \frac{1}{\sum_{i=1}^n \sum_{j=1}^n w_{i,j}} \frac{\sum_{i=1}^n \sum_{j=1}^n \tilde{z}_i w_{i,j} \tilde{z}_j}{\frac{1}{n} \sum_{i=1}^n \tilde{z}_i^2}. \quad (1.1)$$

Moran's I statistic tests whether there are some relationships between location and associated values. A statistically significant positive statistic indicates that the associated values of nearby locations are more clustered than randomly distributed. It has been widely used to identify the spatial pattern of ACS estimates, facilitating the analysis of regional convergence/inequality and corresponding policy and decision making. For example, it was used to study the relationship between lithium concentration in drinking water and suicide rate [3]. It was used in dialectology to measure the spatial variation of language [4]. It was also used to measure spatial genetic structure in plant populations [5]. Particularly, it is very useful in social study. For example, it was used to investigate criminal rate [6]. However, when the observation is uncertain, the value of the statistic and inference are no longer stable. No work to date has assessed the impacts of observation uncertainty on spatial autocorrelation identification, nor developed new spatial autocorrelation statistics that are robust to the input uncertainty. This research seeks to fill these gaps. Chapter 2 describes estimating Moran's I by sampling and the bounds of error of the estimate. Chapter 3 describes a corrected version of Moran's I that reduces bias of the estimate when



the value  $x_i$  can only be measured with noise.

## CHAPTER 2

### ESTIMATE MORAN'S I BY SAMPLING

#### 2.1 Understanding the Error from Sampling

In this section we discuss implications of error in the estimate of  $I_{\text{true}}$  when the data is generated by randomly sampling from a source. First, we identify the model variance implied by sampling and how this relates to the Gaussian distribution assumed above. Next, we show how to estimate  $I_{\text{true}}$  using only a random sample of pairs of data points; an important consideration for very large datasets. Finally, we show that this estimate can be improved by using importance sampling to select certain pairs with higher probability, then adjusting their weight in the final estimation to make it unbiased.

##### 2.1.1 Modeling Data from Sampling

Consider the case where  $S = \{p_1, p_2, \dots, p_n\}$  is finite but large point set. The relationships measured by  $w_{i,j}$  is fixed and easy to calculate, but the value  $x_i$  is hard to obtain. For example,  $x_i$  may be constructed as either an aggregate or average from a population associated with  $p_i$ . So it will be a saving of resources if we can estimate Moran's I with needed accuracy by sampling from  $S$ . Next we show that is doable.

First we define two sets  $\mathbb{A} = \{a(i, j) = \tilde{z}_i \tilde{z}_j\}$  and  $\mathbb{B} = \{b(i) = \tilde{z}_i^2\}$ . Then we define two new random variables:  $A$  represents the distribution of random sample from set  $\mathbb{A}$  and  $B$  represents the distribution of random sample from set  $\mathbb{B}$ . Then using  $W = \sum_{i=1}^n \sum_{j=1}^n w_{i,j}$ , define

$$Z_w = \frac{1}{W} \sum_{i=1}^n \sum_{j=1}^n \tilde{z}_i w_{i,j} \tilde{z}_j = \frac{1}{W} \sum_{i=1}^n \sum_{j=1}^n w_{i,j} a(i, j),$$
$$Z_2 = \frac{1}{n} \sum_{i=1}^n \tilde{z}_i^2 = \frac{1}{n} \sum_{i=1}^n b_i.$$

We can see that  $Z_w$  is the weighted mean of the random variable  $A$  with weight given by  $w_{i,j}$  and  $Z_2$  is the mean of the random variable  $B$ . Also observe that  $I_{\text{true}} = \frac{Z_w}{Z_2}$ . Note that if all of the weights  $w_{i,j}$  are known, then  $W$  is a known constant.

We can then estimate  $I_{\text{true}}$  in two steps. First, to estimate  $Z_2$  we simply sample  $\tilde{z}_i$ s uniformly and apply Chebyshev's Inequality to get bounds. To estimate  $Z_w$  we can uniformly sample pairs points  $(i,j)$  and get the values of  $\tilde{z}_i$  and  $\tilde{z}_j$  as well as  $w_{i,j}$ . However,  $Z_w$  can also be estimated as weighted mean of  $\tilde{z}_i\tilde{z}_j$  with weight given by  $w_{i,j}$ . We can improve performance of the estimation by importance sampling upon uniform sampling. Then we can apply Chebyshev's Inequality again to give bounds on this estimate.

### 2.1.2 Algorithm to estimate $Z_w$

Now we give an algorithm to estimate  $Z_w$ . Assuming that we have a fixed set of size  $n$  for which we want to calculate the Moran's I and that we have calculated the weight  $w_{i,j}$  and  $w_{i,\cdot} = \sum_{j=1}^n w_{i,j}$ . We denote the distribution with probability mass function (pmf)  $P(X = i) = w_{i,\cdot}$  as  $F$  and the distribution with pmf  $P(X = j) = w_{i,j}/w_{i,\cdot}$  as  $F_i$ . We want to sample  $n_1$  pairs of points by importance sampling to estimate  $Z_w$ . The Algorithm 1, below, implements the importance sampling we want.

---

#### Algorithm 1: Importance sampling to estimate $Z_w$

---

- 1: set  $S = 0$
  - 2: **for**  $i \leftarrow 1$  to  $n_1$  **do**
  - 3:   generate a random integer in  $i = [1, n]$  with a distribution  $F$ .
  - 4:   generate a random integer in  $j = [1, n]$  with a distribution  $F_i$ .
  - 5:   estimate  $\tilde{z}_i$  and  $\tilde{z}_j$
  - 6:    $S = S + \tilde{z}_i\tilde{z}_j$
  - 7: **return**  $S/n_1$
- 

### 2.1.3 Importance Sampling of Pairs

We can improve upon the above result by applying importance sampling. We have the importance sampling property:

**Property 2.1** Suppose we want to find  $\mu = E(f(X)) = \int_D f(x) p(x) dx$  where  $p$  is a probability density function on  $D \subseteq \mathbb{R}^d$  and  $f$  is the integrand. We take  $p(x) = 0$  for all  $x \notin D$ . If  $q$  is a positive probability density function on  $\mathbb{R}^d$ , then  $\mu = E_p(f(X)) = E_q\left(\frac{f(x)p(x)}{q(x)}\right)$  where  $E_p(\cdot)$  denotes expectation for  $X \sim p$  and  $E_q(\cdot)$  denotes expectation for  $X \sim q$ . When the distribution of  $X$  is discrete, simply change the integration to summation and density function to mass function and the theorem still holds.

Now we apply property 2.1 above to estimate  $Z_w$ .

We have  $Z_w = \frac{1}{W} \sum_{i=1}^n \sum_{j=1}^n \tilde{z}_i w_{i,j} \tilde{z}_j = \frac{n^2}{W} \sum_{i=1}^n \sum_{j=1}^n \tilde{z}_i w_{i,j} \tilde{z}_j \frac{1}{n^2} \equiv \frac{n^2}{W} \mu$  where  $\mu = \sum_{i=1}^n \sum_{j=1}^n \tilde{z}_i w_{i,j} \tilde{z}_j \frac{1}{n^2}$  is the expectation of  $f(i, j) \equiv \tilde{z}_i w_{i,j} \tilde{z}_j$ , where  $f(i, j) \sim p$  and  $p(i, j) \equiv \frac{1}{n^2}$  is a mass function.

Now let  $q(i, j) = \frac{w_{i,j}}{W}$  and  $q(i, j)$  is also a mass function. Then by theorem 2.1, we have

$$Z_w = \frac{n^2}{W} \mu = \frac{n^2}{W} E_q\left(\frac{f(i, j) p(i, j)}{q(i, j)}\right) = \frac{n^2}{W} E_q\left(\frac{(\tilde{z}_i w_{i,j} \tilde{z}_j) (\frac{1}{n^2})}{\frac{w_{i,j}}{W}}\right) = \frac{n^2}{W} E_q\left((\tilde{z}_i \tilde{z}_j) \left(\frac{W}{n^2}\right)\right) = E_q((\tilde{z}_i \tilde{z}_j)).$$

Therefore the importance sampling estimate of  $Z_w$  is  $\hat{Z}_w = \frac{1}{m} \sum_{i=1}^m (\tilde{z}_{i1} \tilde{z}_{i2})$ , where  $\tilde{z}_{i1}$  and  $\tilde{z}_{i2}$  are the two values of the sampled  $i$ 'th pair and the probability of the  $i$ 'th pair to be sampled is  $\frac{w_{i1,i2}}{W}$ .

## 2.1.4 Approximation Bounds

Next, we prove that  $\hat{I}$  estimated by sampling described above is close to  $I_{\text{true}}$  with high probability for large enough sample sizes. We will achieve this by proving that numerators  $\hat{Z}_w$  and  $Z_w$  are close, then showing the denominators  $\hat{Z}_2$  and  $Z_2$  are also close. Then we will combine these facts together to provide an overall bound.

### 2.1.4.1 Chebyshev's Inequality

We first represent Chebyshev's Inequality that will be used in our proof.

- Chebyshev's Inequality: if  $X$  is an arbitrary random variable and  $t > 0$ , then

$$\Pr(|X - E(X)| \geq t) \leq \frac{\text{Var}(X)}{t^2}.$$

Now consider a set  $X_1, X_2, \dots, X_n$  of  $n$  uncorrelated random variables (that is,  $\text{Cov}(X_i, X_j) = 0$ ) and their sum  $S_n = \sum_{i=1}^n X_i$ . Since  $\text{Cov}(X_i, X_j) = 0$ , we have  $\text{Var}(\sum_{i=1}^n X_i) = \sum_{i=1}^n \text{Var}(X_i)$ . Then for any  $t > 0$ , Chebyshev's Inequality implies:

$$\Pr(|S_n - E(S_n)| \geq t) \leq \frac{\text{Var}(S_n)}{t^2} = \frac{\sum_{i=1}^n \text{Var}(X_i)}{t^2}.$$

Furthermore, let  $\sigma^2 = \frac{1}{n} \sum_{i=1}^n \text{Var}(X_i)$  then  $\Pr(\frac{1}{n} |S_n - E(S_n)| \geq d) \leq \frac{\sigma^2}{nd^2}$ .

### 2.1.4.2 Closeness in $\hat{Z}_w$ .

**Lemma 2.1.1.** For any parameter  $t > 0$ , when  $n > \frac{C_1}{t^2 d}$  then  $\Pr[|\hat{Z}_w - Z_w| > t Z_w] \leq d$ , where  $C_1$  is a constant that is larger than  $\text{Var}(A) / Z_w^2$ .

*Proof.* We assume that  $\text{Var}(A) = V_A$ ,  $\text{mean}(A) = \mu_A$ . To estimate  $Z_w$ , we apply importance sampling to get a sample of size  $n$  from A. We have proved that  $\hat{Z}_w = \frac{1}{n} \sum_{i=1}^m (A_i)$  is an unbiased estimate of  $Z_w$ . Also We have  $\text{Var}(\hat{Z}_w) = V_A/n$ . Now we apply Chebyshev's Inequality to get bounds of estimates of  $\hat{Z}_w$ .

$$\Pr(|\hat{Z}_w - Z_w| \geq tZ_w) \leq \frac{\text{Var}(\hat{Z}_w)}{t^2 Z_w^2} = \frac{V_A}{nt^2 Z_w^2}.$$

So if  $n > \frac{V_A}{dt^2 Z_w^2}$  then  $\Pr(|\hat{Z}_w - Z_w| \geq tZ_w) \leq d$ .

In a special case that we know that  $V_A < C_1 Z_w^2$ , then if  $n > \frac{C_1}{dt^2}$ , then  $\Pr[|\hat{Z}_w - Z_w| > tZ_w] \leq d$ .

□

### 2.1.4.3 Closeness in $\hat{Z}_2$ .

**Lemma 2.1.2.** For any parameter  $t > 0$ , when  $n > \frac{C_2}{t^2 d}$  then  $\Pr[|\hat{Z}_2 - Z_2| > tZ_2] \leq d$ , where  $C_2$  is a constant that is larger than  $\text{Var}(B)/Z_2^2$ .

*Proof.* We assume that  $\text{Var}(B) = V_B$  and  $\text{mean}(B) = \mu_B$ . To estimate  $Z_2$ , we uniformly sample a sample of size  $n$  from B. Obviously  $\hat{Z}_2 = \frac{1}{n} \sum_{i=1}^m (B_i)$  is an unbiased estimate of  $Z_2$ .

Also We have  $\text{Var}(\hat{Z}_2) = V_B/n$  Now we apply Chebyshev's Inequality to get bounds of estimates of

$$\Pr(|\hat{Z}_2 - Z_2| \geq tZ_2) \leq \frac{\text{Var}(\hat{Z}_2)}{t^2 Z_2^2} = \frac{V_B}{nt^2 Z_2^2}.$$

So if  $n > \frac{V_B}{dt^2 Z_2^2}$  then  $\Pr(|\hat{Z}_2 - Z_2| \geq tZ_2) \leq d$ . In a special case that we know that  $V_B < C_2 Z_2^2$ , then if  $n > \frac{C_2}{dt^2}$  then  $\Pr[|\hat{Z}_2 - Z_2| > tZ_2] \leq d$ . □

### 2.1.4.4 Putting them together

Finally we put these results about the numerator and denominator stability together.

**Theorem 2.1.1.** Consider a dataset  $P$  of size  $m$  where each value point has an associated value  $z_i$ . Define a random variable  $A(i) = z_{i_1} z_{i_2}$  where  $z_{i_1}$  and  $z_{i_2}$  are a pair of values and  $B(i) = z_{i_1}^2$ . If we apply importance sampling with  $q(i) = w(i) = w_{i_1, i_2}$  to sample  $n_1$  pairs of points to estimate  $Z_w$  and uniformly obtain another sample of size  $n_2$  to estimate  $Z_2$ , Then, for any parameter  $\alpha \in (0, 1)$ , if  $n_1 > \frac{18C_1}{\delta\alpha^2}$  and  $n_2 > \frac{18C_2}{\delta\alpha^2}$ , with probability at least  $1 - \delta$ , we have

$$1 - \alpha \leq \frac{\hat{I}}{I_{\text{true}}} \leq 1 + \alpha.$$

*Proof.* Overall, when all the conditions mentioned in the proof above hold, which are  $n_1 > \frac{C_1}{dt^2}$  and  $n_2 > \frac{C_2}{dt^2}$ , then

$$\Pr(\hat{Z}_2/Z_2 \leq 1 - t \text{ or } \hat{Z}_2/Z_2 \geq 1 + t) = \Pr(|\hat{Z}_2 - Z_2| \geq tZ_2) \leq d$$

and

$$\Pr(\hat{Z}_w/Z_w \leq 1 - t \text{ or } \hat{Z}_w/Z_w \geq 1 + t) = \Pr(|\hat{Z}_w - Z_w| \geq tZ_w) \leq d$$

hold and then

$$\Pr((\hat{Z}_w/Z_w)/(\hat{Z}_2/Z_2) < (1 - t)/(1 + t) \text{ or } (\hat{Z}_w/Z_w)/(\hat{Z}_2/Z_2) < (1 + t)/(1 - t)) \leq (2d).$$

And since  $(\hat{Z}_w/Z_w)/(\hat{Z}_2/Z_2) = \hat{I}/I_{\text{true}}$  so

$$\Pr(\hat{I}/I_{\text{true}} < (1 - t)/(1 + t) \text{ or } \hat{I}/I_{\text{true}} > (1 + t)/(1 - t)) \leq (2d).$$

So

$$\Pr((1 - t)/(1 + t) \leq \hat{I}/I_{\text{true}} \leq (1 + t)/(1 - t)) \geq (1 - 2d).$$

Now, we rewrite the conditions and inequality in terms of  $\alpha$  and  $\delta$ .

Let  $t = \frac{\alpha}{2 + \alpha}$  then

$$\begin{aligned} (1 + t)/(1 - t) &= (1 + \frac{\alpha}{2 + \alpha}) / (1 - \frac{\alpha}{2 + \alpha}) \\ &= (2 + \alpha + \alpha) / (2 + \alpha - \alpha) \\ &= (2 + 2\alpha) / 2 \\ &= (1 + \alpha) \end{aligned}$$

and

$$\begin{aligned} (1 - t)/(1 + t) &= (1 - \frac{\alpha}{2 + \alpha}) / (1 + \frac{\alpha}{2 + \alpha}) \\ &= (2 + \alpha - \alpha) / (2 + \alpha + \alpha) \\ &= 2 / (2 + 2\alpha) \\ &= 1 / (1 + \alpha) \\ &= 1 - \alpha / (1 + \alpha) \\ &> 1 - \alpha \end{aligned}$$

Also let  $d = \frac{1}{2}\delta$  then

$$1 - 2d = 1 - \delta$$

So

$$\Pr(1 - \alpha \leq I_{\text{corr}}/I_{\text{true}} \leq 1 + \alpha) \geq (1 - \delta).$$

Now plug  $t = \frac{\alpha}{2 + \alpha}$  and  $\varepsilon = \frac{1}{2}\delta$  in to the conditions we mentioned in the proof of Closeness in denominators and denominators. We get the conditions in terms of  $\alpha$  and  $\delta$  and they are

- $n_1 > \frac{C_1}{\left(\frac{\alpha}{2+\alpha}\right)^2 \frac{1}{2}\delta} = \frac{2(2+\alpha)^2 C_1}{\alpha^2 \delta}$
- $n_2 > \frac{C_2}{\left(\frac{\alpha}{2+\alpha}\right)^2 \frac{1}{2}\delta} = \frac{2(2+\alpha)^2 C_2}{\alpha^2 \delta}$

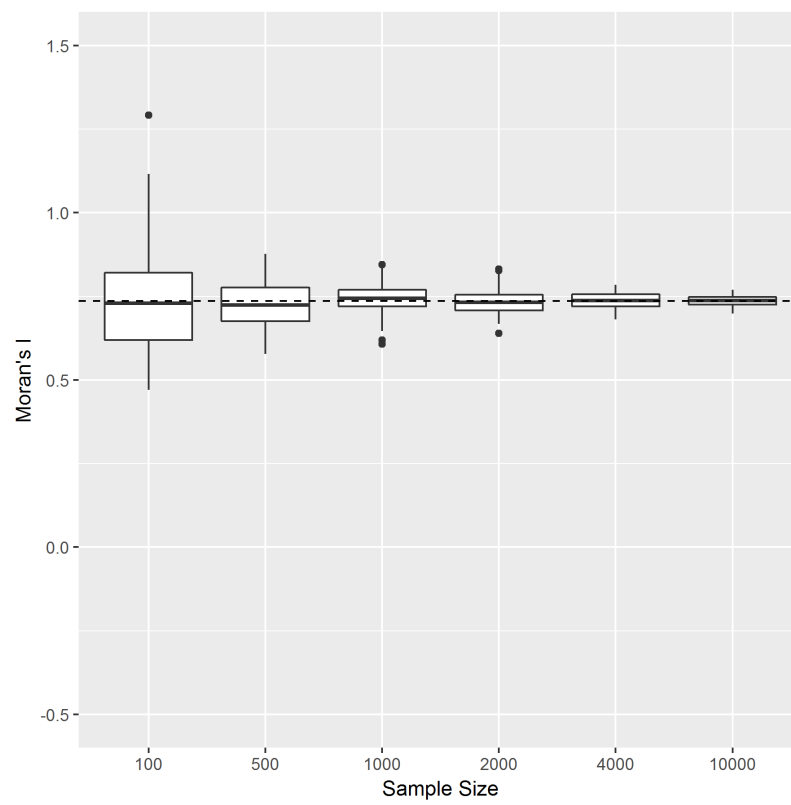
Since  $0 < \alpha < 1$ , conditions below is looser than conditions above.

- $n_1 > \frac{18C_1}{\alpha^2 \delta}$
- $n_2 > \frac{18C_2}{\alpha^2 \delta}$

That finishes the proof of the theorem. □

### 2.1.5 Simulation

A simulation study was performed to investigate how sample size of the sampling algorithm affect the accuracy of Moran's I estimation. The true attribute value  $x_i$  for each point  $p_i$  on an 50 by 50 grid is generated using R package *gstat*. The weight matrix  $W$  is defined so  $w_{i,j}$  is 1 when  $p_i$  and  $p_j$  are adjacent on the grid, and 0 otherwise. True Moran's I was calculated using formula 1.1 as the reference. The sampling algorithms was applied to the generated data for 100 times for each of tested sample sizes. The Moran's I estimation for different sample sizes are summarized by box plot on figure 2.1. It is clearly shown that with sample size increases, the estimation of I approaches the true I.



**Figure 2.1.** The effect of sample size on accuracy of Moran's I estimation



## CHAPTER 3

### CORRECTED MORAN'S I

#### 3.1 Corrected Moran's I when there is noise

Very often, the true value associated with point  $p_i$  (denoted  $\tilde{z}_i$ ) is not observed. Instead what we observe is  $z_i = \tilde{z}_i + \varepsilon_i$ , where  $\varepsilon_i$  is an unknown noise term, and we assume  $\varepsilon_i \perp \varepsilon_j$  (they are independent). We will also assume that  $\varepsilon_i \sim N(0, v_i)$ . It is now typically for uncertain data to provide some information from which we can at least construct a good estimate of this. When  $z_i$  (instead of  $\tilde{z}_i$ ) is used to calculate Moran's I, the estimated Moran's I will be smaller than the true one. This is similar to what happens in the calculation of the correlation coefficient. Let  $\beta = \tilde{\beta} + \varepsilon_\beta$  and  $\theta = \tilde{\theta} + \varepsilon_\theta$  where  $\tilde{\beta}$  and  $\tilde{\theta}$  are the true values and  $\beta$  and  $\theta$  are estimate of  $\tilde{\beta}$  and  $\tilde{\theta}$ , respectively, where  $\varepsilon_\beta$  and  $\varepsilon_\theta$  are the associated measurement errors. Then the observed correlation  $\text{corr}(\beta, \theta)$  is smaller than the true correlation  $\text{corr}(\tilde{\beta}, \tilde{\theta})$ . This phenomenon is called *disattenuation* and can be corrected by the equation:  $\text{corr}(\tilde{\beta}, \tilde{\theta}) = \text{corr}(\beta, \theta) / \sqrt{R_\beta R_\theta}$  where  $R_\beta = (\text{Var}(\beta) - \text{Var}(\varepsilon_\beta)) / \text{Var}(\beta)$  and  $R_\theta = (\text{Var}(\theta) - \text{Var}(\varepsilon_\theta)) / \text{Var}(\theta)$ .

Inspired by the correction of disattenuation in correlation coefficient, we propose to correct the Moran's I by

$$I_{\text{corr}} = \frac{1}{\sum_i \sum_j \omega_{i,j} \frac{1}{n} \sum_i z_i^2 - \sum_i v_i} \sum_i \sum_j z_i \omega_{i,j} z_j. \quad (3.1)$$

##### 3.1.1 Approximation Bounds

Next, we prove that  $I_{\text{corr}}$  is close to  $I_{\text{true}}$  with high probability for large enough dataset sizes. We will achieve this by proving that numerators  $N \equiv \sum_i \sum_j z_i \omega_{i,j} z_j$  and  $N_0 \equiv \sum_i \sum_j \tilde{x}_i \omega_{i,j} \tilde{x}_j$  are close, then showing the denominators  $D \equiv \frac{1}{n} \sum_i z_i^2 - \sum_i v_i$  and  $D_0 \equiv \frac{1}{n} \sum_i \tilde{z}_i^2$  are also close. Then we will combine these facts together to provide an overall

bound. We assume that the noise is bounded such that  $v_i \leq v = c_1 D_0$  ( $c_1$  is a constant). By definition,  $D_0$  is the variance of the true value among the locations, and  $v$  is the bound of noise. So  $c_1$  can be seen as the ratio of noise to signal. Since rescaling the weight  $w_{ij}$  doesn't change  $I_{corr}$ , without losing generality, we can assume  $0 \leq \omega_{i,j} \leq 1$ , and in this paper we only consider the case that the weight matrix is sparse such that  $\sum_i \omega_{i,j} = \sum_j \omega_{i,j} < w \forall i$  but not too sparse such that  $W = \sum_i \sum_j \omega_{i,j} > c_2 n$ .

We first list a few standard properties about random variables  $X$  and  $Y$  we will use.

(P1)  $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$

(P2) If  $X \perp Y$  then  $E(XY) = E(X)E(Y)$ .

that is, if  $X$  and  $Y$  are independent, then we can decompose the product of their expected values.

(P3) If  $X \perp Y$  then  $\text{Var}(XY) = E(X)^2 \text{Var}(Y) + E(Y)^2 \text{Var}(X) + \text{Var}(X)\text{Var}(Y)$ , and further if  $E(X) = 0$  and  $E(Y) = 0$  then  $\text{Var}(XY) = \text{Var}(X)\text{Var}(Y)$ .

(P4) Chebyshev's Inequality: if  $X$  is an arbitrary random variable and  $t > 0$ , then

$$\Pr(|X - E(X)| \geq t) \leq \frac{\text{Var}(X)}{t^2}.$$

Now consider a set  $X_1, X_2, \dots, X_n$  of  $n$  uncorrelated random variables (that is,  $\text{Cov}(X_i, X_j) = 0$ ) and their sum  $S_n = \sum_{i=1}^n X_i$ . Since  $\text{Cov}(X_i, X_j) = 0$ , we have  $\text{Var}(\sum_{i=1}^n X_i) = \sum_{i=1}^n \text{Var}(X_i)$ . Then for any  $t > 0$ , Chebyshev's Inequality implies:

$$\Pr(|S_n - E(S_n)| \geq t) \leq \frac{\text{Var}(S_n)}{t^2} = \frac{\sum_{i=1}^n \text{Var}(X_i)}{t^2}.$$

Furthermore, let  $\sigma^2 = \frac{1}{n} \sum_{i=1}^n \text{Var}(X_i)$  then  $\Pr(\frac{1}{n}|S_n - E(S_n)| \geq d) \leq \frac{\sigma^2}{nd^2}$ .

### 3.1.1.1 Closeness in numerators.

**Lemma 3.1.1.** For any parameter  $t > 0$ , when  $n > \frac{w(2+c_1)c_1}{t^2(1_{true}c_2)^2d}$  then  $\Pr[|N - N_0| > tN_0] \leq d$

*Proof.* We first expand the definition of  $N$

$$N = \sum_i \sum_j z_i \omega_{i,j} z_j = \sum_i \sum_j (\tilde{z}_i + \varepsilon_i) \omega_{i,j} (\tilde{z}_j + \varepsilon_j) = \sum_i \sum_j (\tilde{z}_i \tilde{z}_j + \tilde{z}_i \varepsilon_j + \tilde{z}_j \varepsilon_i + \varepsilon_i \varepsilon_j) \omega_{ij}.$$

Since  $E[\varepsilon_i] = 0$  and  $\varepsilon_i \perp \varepsilon_j$ , then  $E(\tilde{z}_i \varepsilon_j) = E(\tilde{z}_j \varepsilon_i) = E(\varepsilon_i \varepsilon_j) = 0$ , hence  $E(N) = E(N_0)$ .

Next we can argue that  $\varepsilon_i$  and  $\varepsilon_i\varepsilon_j$  are uncorrelated since

$$\text{Cov}(\varepsilon_i, \varepsilon_i\varepsilon_j) = \text{E}(\varepsilon_i\varepsilon_i\varepsilon_j) - \text{E}(\varepsilon_i)\text{E}(\varepsilon_i\varepsilon_j) = \text{E}(\varepsilon_i^2)\text{E}(\varepsilon_j) - \text{E}(\varepsilon_i)\text{E}(\varepsilon_i\varepsilon_j) = 0 - 0 = 0.$$

Moreover we have  $\varepsilon_i\varepsilon_k$  and  $\varepsilon_\ell\varepsilon_j$  are uncorrelated for any  $i \neq k \neq j \neq \ell$  because

$$\text{Cov}(\varepsilon_i\varepsilon_k, \varepsilon_i\varepsilon_j) = \text{E}(\varepsilon_i\varepsilon_k\varepsilon_i\varepsilon_j) - \text{E}(\varepsilon_i\varepsilon_k)\text{E}(\varepsilon_i\varepsilon_j) = \text{E}(\varepsilon_i^2)\text{E}(\varepsilon_k)\text{E}(\varepsilon_j) - \text{E}(\varepsilon_i)\text{E}(\varepsilon_k)\text{E}(\varepsilon_i\varepsilon_j) = 0 - 0 = 0.$$

So all the terms in  $N$  are uncorrelated from each other. Also  $\text{Var}(\tilde{z}_i\varepsilon_j) = \tilde{z}_i^2v_i$  and  $\text{Var}(\varepsilon_i\varepsilon_j) =$

$\text{Var}(\varepsilon_i)\text{Var}(\varepsilon_j) = v_iv_j$  due to that  $\varepsilon_i \perp \varepsilon_j$  and  $\text{E}(\varepsilon_i) = \text{E}(\varepsilon_j) = 0$ . So  $\text{Var}(N) = \sum_i \sum_j (\tilde{z}_i^2v_j + \tilde{z}_j^2v_i + v_iv_j)\omega_{ij}^2$

According to Chebyshev's Inequality,

$$\Pr(|N - N_0| > tN_0) \leq \frac{\sum_i \sum_j (\tilde{z}_i^2v_j + \tilde{z}_j^2v_i + v_iv_j)\omega_{ij}^2}{t^2N_0^2}.$$

In a special case where  $v_i \leq v = c_1D_0$  ( $c_1$  is a constant),  $0 \leq \omega_{i,j} \leq 1$ , and the weight matrix is sparse such that  $\sum_i \omega_{i,j} = \sum_j \omega_{i,j} < w \forall i$  but not too sparse such that  $\sum_i \sum_j \omega_{i,j} = W > c_2n$  then

$$\begin{aligned} \Pr(|N - N_0| > tN_0) &\leq \frac{\sum_i \sum_j (\tilde{z}_i^2v_j + \tilde{z}_j^2v_i + v_iv_j)\omega_{ij}^2}{t^2N_0^2} \leq \frac{\sum_i \sum_j (\tilde{z}_i^2 + \tilde{z}_j^2 + v)v\omega_{ij}}{t^2N_0^2} \\ &= \frac{(\sum_i \sum_j \tilde{z}_i^2\omega_{ij} + \sum_i \sum_j \tilde{z}_j^2\omega_{ij} + \sum_i \sum_j v\omega_{ij})v}{t^2N_0^2} \\ &\leq \frac{(\sum_i \tilde{z}_i^2w + \sum_j \tilde{z}_j^2w + \sum_i vw)v}{t^2N_0^2} \\ &= \frac{(nD_0w + nD_0w + nvw)v}{t^2N_0^2} \\ &= \frac{(nD_0w + nD_0w + nc_1D_0w)c_1D_0}{t^2N_0^2} \\ &= \frac{nD_0^2w(2 + c_1)c_1}{t^2N_0^2} \\ &= \frac{nD_0^2w(2 + c_1)c_1}{t^2(I_{\text{true}}D_0W)^2} \\ &= \frac{nw(2 + c_1)c_1}{t^2(I_{\text{true}}W)^2} \\ &\leq \frac{nw(2 + c_1)c_1}{t^2(I_{\text{true}}c_2n)^2} \\ &= \frac{nw(2 + c_1)c_1}{t^2(I_{\text{true}}c_2n)^2} \\ &= \frac{w(2 + c_1)c_1}{t^2(I_{\text{true}}c_2)^2n} \\ &\equiv p_1 \end{aligned}$$

In the case that  $n$  is a big number,  $p_1$  is a small number. Specifically, when  $n > \frac{w(2+c_1)c_1}{t^2(1_{true}c_2)^2d}$ , then  $p_1 < d$

### 3.1.1.2 Closeness in denominators.

**Lemma 3.1.2.** For any parameter  $t > 0$ , when  $n > \frac{(4+2c_1)c_1}{dt^2}$  then  $\Pr(|D - D_0| > tD_0) < d$

*Proof.* Recall

$$D \equiv \frac{1}{n} \sum_i z_i^2 - \sum_i v_i = \frac{1}{n} \sum_i ((\tilde{z}_i + \varepsilon_i)^2 - v_i) = \frac{1}{n} \sum_i (\tilde{z}_i^2 + 2\tilde{z}_i\varepsilon_i + \varepsilon_i^2 - v_i).$$

Moreover, since  $E(\tilde{z}_i\varepsilon_i) = 0$  and  $\text{Var}(\tilde{z}_i\varepsilon_i) = \tilde{z}_i^2 v_i$  and because  $\varepsilon_i \sim N(0, v_i)$  then hence  $\frac{\varepsilon_i^2}{v_i} \sim \chi_1^2$  and  $E(\frac{\varepsilon_i^2}{v_i}) = 1$ . This implies  $E(\varepsilon_i^2 - v_i) = 0$ , that  $\text{Var}(\frac{\varepsilon_i^2}{v_i}) = 2$  and  $\text{Var}(\varepsilon_i) = 2v_i^2$ .

Also

$$\text{Var}(2\tilde{z}_i\varepsilon_i + \varepsilon_i^2) = \text{Var}(2\tilde{z}_i\varepsilon_i) + 2\text{Cov}(2\tilde{z}_i\varepsilon_i, \varepsilon_i^2) + \text{Var}(\varepsilon_i^2)$$

But

$$\text{Cov}(2\tilde{z}_i\varepsilon_i, \varepsilon_i^2) = E(2\tilde{z}_i\varepsilon_i\varepsilon_i^2) - E(2\tilde{z}_i\varepsilon_i)E(\varepsilon_i^2) = 2\tilde{z}_iE(\varepsilon_i^3) - 2\tilde{z}_iE(\varepsilon_i)E(\varepsilon_i^2) = 2\tilde{z}_i \cdot 0 - 2\tilde{z}_i \cdot 0E(\varepsilon_i^2) = 0$$

The reason that  $E(\varepsilon_i^3) = 0$  is that  $\varepsilon_i \sim N(0, V_i)$  which is a distribution symmetric about 0.

So

$$\text{Var}(2\tilde{z}_i\varepsilon_i + \varepsilon_i^2) = \text{Var}(2\tilde{z}_i\varepsilon_i) + 2\text{Cov}(2\tilde{z}_i\varepsilon_i, \varepsilon_i^2) + \text{Var}(\varepsilon_i^2) = 4\tilde{z}_i^2 v_i + 2v_i^2$$

and  $E(D) = E(D_0)$  due to that  $E(\tilde{z}_i\varepsilon_i) = 0$  and  $E(\varepsilon_i^2 - v_i) = 0$ .

In summary,  $E(D) = E(D_0)$  and  $\text{Var}(D) = \frac{1}{n^2} \sum_i (4\tilde{z}_i^2 v_i + 2v_i^2)$ . Now according to Chebyshev's Inequality,

$$\begin{aligned}
\Pr(|D - D_0| > tD_0) &\leq \frac{\frac{1}{n^2} \sum_i (4\tilde{z}_i^2 v_i + 2v_i^2)}{t^2 D_0^2} \\
&= \frac{\sum_i (4\tilde{z}_i^2 v_i + 2v_i^2)}{n^2 t^2 D_0^2} \\
&\leq \frac{\sum_i (4\tilde{z}_i^2 + 2v) v}{n^2 t^2 D_0^2} \\
&= \frac{(4nD_0 + 2nv)v}{n^2 t^2 D_0^2} \\
&= \frac{(4nD_0 + 2nc_1 D_0)c_1 D_0}{n^2 t^2 D_0^2} \\
&= \frac{(4 + 2c_1)c_1}{nt^2} \\
&\equiv p_2
\end{aligned}$$

and when  $n$  is a big number  $p_2$  is small. In particular, when  $n > \frac{(4+2c_1)c_1}{dt^2}$  then  $p_2 < d$ .  $\square$

### 3.1.1.3 Putting these together.

Finally we put these results about the numerator and denominator stability together.

**Theorem 3.1.1.** Consider a dataset  $P$  of size  $n$  where each value  $z_i = \tilde{z}_i + \varepsilon_i$  has independent Gaussian noise  $\varepsilon_i \sim \mathcal{N}(0, v_i)$ , that  $v_i \leq v = c_1 D_0$  ( $c_1$  is a constant),  $0 \leq \omega_{i,j} \leq 1$ , and the weight matrix is sparse such that  $\sum_i \omega_{i,j} = \sum_j \omega_{i,j} < w \forall i$  but not too sparse such that  $W = \sum_i \sum_j \omega_{i,j} > c_2 n$ . Then, for any parameter  $\alpha \in (0, 1)$ , if  $n > \max(18w \cdot \frac{(2+c_1)c_1}{(c_2)^2} \cdot \frac{1}{(I_{true})^2} \cdot \frac{1}{\alpha^2} \cdot \frac{1}{\delta}, 36(2+c_1)c_1 \cdot \frac{1}{\alpha^2} \cdot \frac{1}{\delta})$ , with probability at least  $1 - \delta$ , we have

$$1 - \alpha \leq \frac{I_{corr}}{I_{true}} \leq 1 + \alpha.$$

*Proof.* Overall, when all the conditions mentioned in the proof above hold, which are  $n > \frac{w(2+c_1)c_1}{t^2(I_{true}c_2)^2d}$  and  $n > \frac{(4+2c_1)c_1}{dt^2}$ , then

$$\Pr(D/D_0 \leq 1 - t \text{ or } D/D_0 \geq 1 + t) = \Pr(|D - D_0| \geq tD_0) \leq d$$

and

$$\Pr(N/N_0 \leq 1 - t \text{ or } N/N_0 \geq 1 + t) = \Pr(|N - N_0| \geq tN_0) \leq d$$

hold and then

$$\Pr((N/N_0)/(D/D_0) < (1-t)/(1+t) \text{ or } (N/N_0)/(D/D_0) < (1+t)/(1-t)) \leq (2d).$$

And since  $(N/N_0)/(D/D_0) = I_{\text{corr}}/I_{\text{true}}$  so

$$\Pr(I_{\text{corr}}/I_{\text{true}} < (1-t)/(1+t) \text{ or } I_{\text{corr}}/I_{\text{true}} > (1+t)/(1-t)) \leq (2d).$$

So

$$\Pr((1-t)/(1+t) \leq I_{\text{corr}}/I_{\text{true}} \leq (1+t)/(1-t)) \geq (1-2d).$$

Now, we rewrite the conditions and inequality in terms of  $\alpha$  and  $\delta$ .

Let  $t = \frac{\alpha}{2+\alpha}$  then

$$\begin{aligned} (1+t)/(1-t) &= (1 + \frac{\alpha}{2+\alpha}) / (1 - \frac{\alpha}{2+\alpha}) \\ &= (2+\alpha+\alpha) / (2+\alpha-\alpha) \\ &= (2+2\alpha) / 2 \\ &= (1+\alpha) \end{aligned}$$

and

$$\begin{aligned} (1-t)/(1+t) &= (1 - \frac{\alpha}{2+\alpha}) / (1 + \frac{\alpha}{2+\alpha}) \\ &= (2+\alpha-\alpha) / (2+\alpha+\alpha) \\ &= 2 / (2+2\alpha) \\ &= 1 / (1+\alpha) \\ &= 1 - \alpha / (1+\alpha) \\ &> 1 - \alpha \end{aligned}$$

Also let  $d = \frac{1}{2}\delta$  then

$$1 - 2d = 1 - \delta$$

So

$$\Pr(1 - \alpha \leq I_{\text{corr}}/I_{\text{true}} \leq 1 + \alpha) \geq (1 - \delta).$$

Now plug  $t = \frac{\alpha}{2+\alpha}$  and  $d = \frac{1}{2}\delta$  in to the conditions we mentioned in the proof of Closeness in denominators and denominators. We get the conditions in terms of  $\alpha$  and  $\delta$  and they are

- $n > \frac{w(2+c_1)c_1}{(\frac{\alpha}{2+\alpha})^2(I_{\text{true}}c_2)^2\frac{1}{2}\delta} = \frac{2w(2+\alpha)^2(2+c_1)c_1}{\alpha^2(I_{\text{true}}c_2)^2\delta}$

- $n > \frac{2(2+c_1)c_1}{(\frac{\alpha}{2+\alpha})^2 \frac{1}{2}\delta} = \frac{4(2+\alpha)^2(2+c_1)c_1}{\alpha^2\delta}$

Since  $0 < \alpha < 1$ , conditions below is looses than conditions above.

- $n > \frac{18w(2+c_1)c_1}{\alpha^2(I_{true}c_2)^2\delta} = 18w \cdot \frac{(2+c_1)c_1}{(c_2)^2} \cdot \frac{1}{(I_{true})^2} \cdot \frac{1}{\alpha^2} \cdot \frac{1}{\delta}$

- $n > \frac{36(2+c_1)c_1}{\alpha^2\delta} = 36(2+c_1)c_1 \cdot \frac{1}{\alpha^2} \cdot \frac{1}{\delta}$

That finishes the proof of the theorem. □

Now we convert the relative bound to additive bound.

**Theorem 3.1.2.** Consider a dataset  $P$  of size  $n$  where each value  $z_i = \tilde{z}_i + \varepsilon_i$  has independent Gaussian noise  $\varepsilon_i \sim \mathcal{N}(0, v_i)$ , that  $v_i \leq v = c_1 D_0$  ( $c_1$  is a constant),  $0 \leq \omega_{i,j} \leq 1$ , and the weight matrix is sparse such that  $\sum_i \omega_{i,j} = \sum_j \omega_{i,j} < w \forall i$  but not too sparse such that  $W = \sum_i \sum_j \omega_{i,j} > c_2 n$ . Then, for any parameter  $\beta \in (0, 1)$ , if  $n > \max(18w \cdot \frac{(2+c_1)c_1}{(c_2)^2} \cdot \frac{1}{\beta^2} \cdot \frac{1}{\delta}, 36(2+c_1)c_1 \cdot \frac{1}{\beta^2} \cdot \frac{1}{\delta})$ , with probability at least  $1 - \delta$ , we have

$$I_{true} - \beta \leq I_{corr} \leq I_{true} + \beta.$$

*Proof.* According to theorem 1.1, when  $n > \max(18w \cdot \frac{(2+c_1)c_1}{(c_2)^2} \cdot \frac{1}{(I_{true})^2} \cdot \frac{1}{\alpha^2} \cdot \frac{1}{\delta}, 36(2+c_1)c_1 \cdot \frac{1}{\alpha^2} \cdot \frac{1}{\delta})$ , with probability at least  $1 - \delta$ , we have

(E1)

$$1 - \alpha \leq \frac{I_{corr}}{I_{true}} \leq 1 + \alpha.$$

That is equivalent to

(E2)

$$I_{true} - \alpha I_{true} \leq I_{corr} \leq I_{true} + \alpha I_{true}.$$

Now let  $\beta = \alpha I_{true}$  then we have  $\alpha = \frac{\beta}{I_{true}}$  and E2 becomes

(E3)

$$I_{true} - \beta \leq I_{corr} \leq I_{true} + \beta.$$

Plug  $\beta = \alpha I_{true}$  and  $\alpha = \frac{\beta}{I_{true}}$  in  $n > \max(18w \cdot \frac{(2+c_1)c_1}{(c_2)^2} \cdot \frac{1}{(I_{true})^2} \cdot \frac{1}{\alpha^2} \cdot \frac{1}{\delta}, 36(2+c_1)c_1 \cdot \frac{1}{\alpha^2} \cdot \frac{1}{\delta})$  and because  $-1 \leq I_{true} \leq 1$ , we can easily get that  $n > \max(18w \cdot \frac{(2+c_1)c_1}{(c_2)^2} \cdot \frac{1}{\beta^2} \cdot \frac{1}{\delta}, 36(2+c_1)c_1 \cdot \frac{1}{\beta^2} \cdot \frac{1}{\delta})$

$c_1)c_1 \cdot \frac{1}{\beta^2} \cdot \frac{1}{\delta}$ ) is tighter than  $n > \max(18w \cdot \frac{(2+c_1)c_1}{(c_2)^2} \cdot \frac{1}{(I_{true})^2} \cdot \frac{1}{\alpha^2} \cdot \frac{1}{\delta}, 36(2+c_1)c_1 \cdot \frac{1}{\alpha^2} \cdot \frac{1}{\delta})$ .

That finishes the proof.

□

### 3.1.2 Simulations

#### 3.1.2.1 scenario 1

A simulation study was performed to investigate if the formula correct the estimate of Moran's I. The true attribute value  $x_i$  for each point  $p_i$  on an 20 by 20 grid is generated using R package *gstat*. The weight matrix  $W$  is defined so  $w_{i,j}$  is  $\exp(-\|p_i - p_j\|^2/2(s^2))$ , where  $\|p_i - p_j\|$  is the Euclidean distance between  $p_i$  and  $p_j$  and  $s$  is a pre-specified scale constant. The observed attribute value  $y_i$  is generated by adding a noise  $\epsilon_i$  to  $x_i$ , and  $\epsilon_i \sim N(0, v)$  where  $v = f \cdot \text{var}(x)$  and  $f$  varies from 0.2 to 1 in the simulation. The true Moran's I is calculated by applying formula 1.1 on  $x_i$ , and the observed Moran's I is calculated by applying formula 1.1 on  $y_i$ , and corrected Moran's I is calculated by applying formula 3.1 on  $y_i$ . The simulation is repeated for 100 times for setting of  $f$  and mean observed I and corrected I are calculated and presented in table 3.1.

f	0.2	0.4	0.6	0.8	1
true I	0.630	0.630	0.630	0.630	0.630
mean observed I	0.523	0.447	0.392	0.351	0.314
mean corrected I	0.628	0.629	0.630	0.633	0.633

**Table 3.1.** Experiments for Gaussian noise

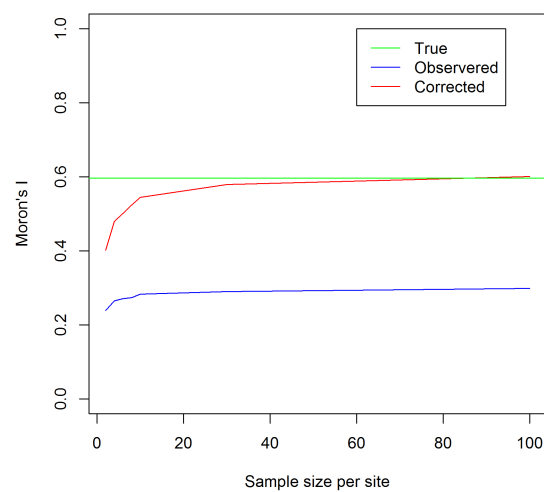
Conclusion: The proposed formula corrected the estimate of Moran's I of Gaussian distribution.

#### 3.1.2.2 scenario 2

It is common that there are multiple source of variation. Suppose at site  $i$  at time point  $t$ , there is a population  $Y_{i,t} \sim G(\theta_{i,t}, \Theta_i)$  where  $\theta_{i,t}$  is the population mean at site  $i$  at time point  $t$  and  $\Theta_i$  is the other parameters of the distribution. We also suppose  $\theta_{i,t}$  itself is also a random variable and that  $\theta_{i,t} \sim N(\theta_i, \sigma_i)$  where  $\theta_i$  are long term mean of  $\theta_{i,t}$ . Now we consider the situation that  $G(\theta_{i,t}, \Theta_i)$  is Gaussian distribution, that is  $Y_{i,t} \sim N(\theta_{i,t}, \Theta_i)$ , and we are interested in the spatial correlation of  $\theta_i$  among all the sites. We need to estimate



$\theta_i$ . There are two source of variation of the estimation of  $\theta_i$ . One is the sampling error, say we can sample a sample of size  $n_i$  at time point  $t$  to estimate  $\theta_{i,t}$ , the other is the variation along the time, say we can only do sampling at  $m$  time point. However, very often we can make  $n_i$  so big that the sampling error is ignorable. Then what we have are  $\theta_{i,t}$  which follows Gaussian distribution and the correction mentioned above can also get unbiased estimation of true Moran's I. A simulation is performed to demonstrate this. The true attribute value  $x_i$  for each point  $p_i$  on an 20 by 20 grid is generated using R package *gstat*. A noise  $\varepsilon_i$  with known variance  $v_1$  is added to the true value. A sample of size  $n$  was drew from Gaussian distribution  $N(x_i + \varepsilon_i, v_2)$  from each site and mean  $m_i$  of the samples were calculated. Observed and corrected Moran's I were calculated using these  $m_i$ . This process was performed for different values of  $n$ . The true Moran's I was also calculated using  $x_i$ . The results are shown on Figure 3.1.

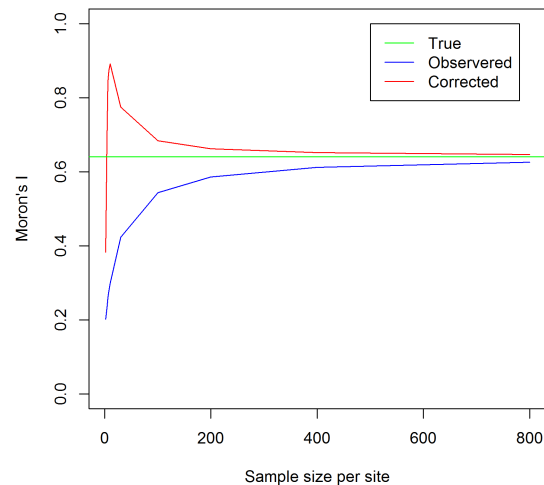


**Figure 3.1.** The effect of sample size on observed and corrected I in Gaussian distribution

Conclusion: The corrected Moran's I estimate approaches the true value as the sample size increases while the original version of Moran's I is always biased even the sample size is large .

### 3.2 Application in Non-Gaussian Distribution

So far we assume that the observed attribute value  $y_i$  follows Gaussian distribution. Now we investigate if the correction we proposed can also be applied to Non-Gaussian distribution, say Bernoulli distribution, by simulation. In this case the attribute values we are interested in is the probability of event  $q_i$ . In the simulation,  $q_i$  are restricted in the range  $[0.1, 0.5]$ . Trial number  $N_i$  in each location varies in the simulation. Event count  $n_i$  is generated follow binomial distribution with proportion  $q_i$  and size  $N_i$  and observed proportion  $q_i^*$  is then calculated by dividing  $n_i$  by  $N_i$ . The variance of  $q_i^*$  is calculated as  $v_i = q_i^* \cdot (1 - q_i^*) / N_i$ . Using  $q_i$ ,  $q_i^*$  and  $v_i$  as input, true Moran's I, observed Moran's I and corrected Moran's I were calculated similar as above and presented on Figure 3.1.



**Figure 3.2.** The effect of sample size on observed and corrected I in Bernoulli distribution

Conclusion: The proposed formula also correct the estimate of Moran's I of Bernoulli distribution.

## CHAPTER 4

### CONCLUSION

During the last decades, the speed of accumulating information has progressed beyond people's imagination. Information is not knowledge however. The first step to gain knowledge from information is to recognize patterns of the information. One of the important patterns is spatial correlation. Moran's I is a commonly used statistic to assess spatial correlation. We provided a sampling method to estimate this statistic and proved bounds of the estimate. We also proposed a corrected formula to estimate Moran's I when there are noise on the attribute values. We proved closeness of the estimate to the true value and demonstrated the closeness property by simulation in 2D space. The future direction can be exploring the performance on the corrected formula on higher dimensional space.

## CHAPTER 5

### REFERENCE

- [1] P. A. Moran, "Notes on continuous stochastic phenomena," *Biometrika*, vol. 37, no. 1/2, pp. 17–23, 1950.
- [2] H. Li, C. A. Calder, and N. Cressie, "Beyond moran's i: Testing for spatial dependence based on the spatial autoregressive model," *Geographical Analysis*, vol. 39, no. 4, pp. 357–375, 2007.
- [3] M. Helbich, M. Leitner, and N. D. Kapusta, "Geospatial examination of lithium in drinking water and suicide mortality," *International journal of health geographics*, vol. 11, no. 1, p. 19, 2012.
- [4] J. Grieve, "A regional analysis of contraction rate in written standard american english," *International Journal of Corpus Linguistics*, vol. 16, no. 4, pp. 514–546, 2011.
- [5] M. Maki and M. Masuda, "Spatial autocorrelation of genotypes in a gynodioecious population of *chionographis japonica* var. *kurohimensis* (liliaceae)," *International journal of plant sciences*, vol. 154, no. 4, pp. 467–472, 1993.
- [6] F. C. Mencken and C. Barnett, "Murder, nonnegligent manslaughter, and spatial autocorrelation in mid-south counties," *Journal of Quantitative Criminology*, vol. 15, no. 4, pp. 407–422, 1999.