

# Structured Correlation Detection with Application to Colocalization Analysis in Dual-Channel Fluorescence Microscopic Imaging

Shulei Wang<sup>\*</sup>, Jianqing Fan<sup>†</sup>, Ginger Pocock<sup>§</sup> and Ming Yuan<sup>\*,‡</sup>

Morgridge Institute for Research<sup>\*</sup> and University of Wisconsin-Madison<sup>\*,§</sup>

and

Princeton University<sup>†</sup>

(April 11, 2016)

## Abstract

Motivated by the problem of colocalization analysis in fluorescence microscopic imaging, we study in this paper structured detection of correlated regions between two random processes observed on a common domain. We argue that although intuitive, direct use of the maximum log-likelihood statistic suffers from potential bias and substantially reduced power, and introduce a simple size-based normalization to overcome this problem. We show that scanning with the proposed size-corrected likelihood ratio statistics leads to optimal correlation detection over a large collection of structured correlation detection problems.

---

<sup>\*</sup>Supported in part by NSF FRG Grant DMS-1265202 and NIH Grant 1U54AI117924-01.

<sup>†</sup>Supported in part by NSF Grants DMS-1206464 and DMS-1406266 and NIH grants R01-GM072611-11.

<sup>‡</sup>Ming Yuan wishes to thank Paul Ahlquist, Kevin Eliceiri and Nathan Sherer for introducing him to colocalization analysis in microscopic imaging, and Richard Samworth for helpful discussions and careful reading of an earlier draft that has led to much improved presentation. Address for correspondence: Department of Statistics, University of Wisconsin-Madison, 1300 University Avenue, Madison, WI 53706.

# 1 Introduction

Most if not all biological processes are characterized by complex interactions among biomolecules such as proteins. A common way to decipher such interactions is through multichannel fluorescence microscopic imaging where each molecule is labeled with fluorescence of a unique emission wavelength, and their biological interactions can be identified with correlation between the expression of fluorescent proteins in certain compartments. Although an ad hoc approach, visual inspection of the overlaid image from both channels is arguably the most common way to determine colocalization in multichannel fluorescence microscopy. Potential pitfalls of this naïve strategy, however, are also well-documented as merged images are heavily influenced by factors such as bleed-through, cross-talk, and relative intensities between different channels. See, e.g., Bolte and Cordelières (2006) and Comeau et al. (2006).

Since the pioneering work of Manders and his collaborators in early 1990s, quantitative methods have also been introduced to colocalization analysis. See, e.g., Manders et al. (1992) and Manders et al. (1993). These approaches typically proceed by first manually selecting a region where the two molecules are likely to colocalize. The degree of colocalization is then measured through various notions of correlation coefficient, most notably Pearson’s correlation coefficient or Manders’ correlation coefficient, computed specific to the chosen area. See Manders et al. (1993), Costes et al. (2004), Parmryd, Adler et al. (2008), Hecce et al. (2013) among others. Obviously, the performance of these approaches depends critically on the manually-selected region of interest, which not only makes the analysis subjective but also creates a bottleneck for high-throughput microscopic imaging processes. Moreover, even if the region is selected in a principled way, colocalization may not be directly inferred from the value of the correlation coefficient computed on the region because the value of the coefficient itself does not translate into statistical significance. This problem could be alleviated using permutation tests, as suggested by Costes et al. (2004). In doing so, however, one still neglects the fact that the region of interest is selected based upon its plausibility of colocalization, and the resulting p-value may appear significant merely because of our failure to adjust for the selection bias. The present work is motivated by this need for an automated and statistically valid way to detect colocalization.

Colocalization analysis can naturally be formulated as an example of a broad class of problems that we shall refer to as structured correlation detection where we observe multiple

collections of random variables on a common domain and want to determine if there is a subset of these variables that are correlated. These types of problem arise naturally in many different fields. For example, in finance, detecting time periods where two common stocks show unusual correlation is essential to the so-called pairs trading strategy (see, e.g. Vidyamurthy, 2004). Other potential examples of structured correlation detection problems can also be found in Chen and Gupta (1997), Robinson et al. (2008), Wieda et al. (2011), and Rodionov (2015), among many others. To fix ideas, in what follows, we shall focus our discussion in the context of colocalization analysis. More specifically, denote by  $\mathbb{I}$  the index set of all pixels in the field of view. In a typical two or three dimensional image,  $\mathbb{I}$  could be a lattice of the corresponding dimension. In practice, it is also possible that  $\mathbb{I}$  is a certain subset of a lattice. For example, when investigating intercellular activities,  $\mathbb{I}$  only includes pixels that correspond to the interior of a cell, or a compartment (e.g. nucleus) of a cell. For each location  $i \in \mathbb{I}$ , let  $X_i$  and  $Y_i$  be the intensities measured at the two channels respectively, as illustrated in the left panel of Figure 1.

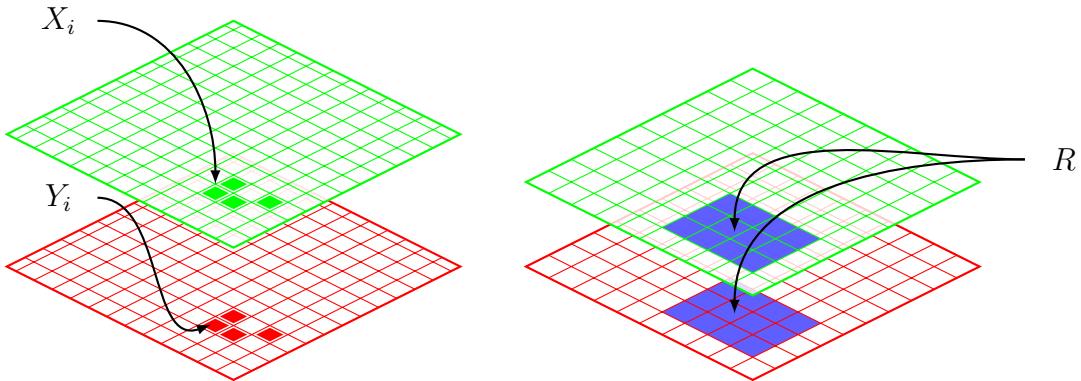


Figure 1: Pixel view of dual channel images.

In the absence of colocalization, we assume that  $X_i$  and  $Y_i$  are uncorrelated and can be modeled as

$$\begin{bmatrix} X_i \\ Y_i \end{bmatrix} \sim N \left( \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \right). \quad (1)$$

where the marginal means  $\mu_1$  and  $\mu_2$  and variances  $\sigma_1^2$  and  $\sigma_2^2$  may be unknown. In the presence of colocalization,  $X_i$  and  $Y_i$  are correlated, and we therefore treat them as observations

from a correlated bivariate normal distribution

$$\begin{bmatrix} X_i \\ Y_i \end{bmatrix} \sim N \left( \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix} \right). \quad (2)$$

When colocalization occurs, it typically does not occur at isolated locations. As a result, the region  $R$  of colocalization is more structured than an arbitrary subset of  $\mathbb{I}$ . For example, colocalization can frequently be observed on a contiguous region  $R$ , as illustrated in the right panel of Figure 1. Let  $\mathcal{R}$  be a library containing all possible regions where correlation may be present. For example,  $\mathcal{R}$  could be the collection of all ellipses or polygons on a two dimensional lattice ( $\mathbb{I}$ ). The primary goal of correlation detection in general, and colocalization analysis in particular, is to check if there is an unknown region  $R \in \mathcal{R}$  such that (1) holds for all  $i \in \mathbb{I} \setminus R$ , and (2) holds for all  $i \in R$  and for some  $\rho \neq 0$ .

The fact that we do not know on which region  $R \in \mathcal{R}$  correlation is present naturally brings about the issue of structured multiple testing. Various aspects of structured multiple testing has been studied in recent years. See, e.g., Desolneux et al. (2003), Pacifico et al. (2004), Arias-Castro et al. (2005), Hall and Jin (2010), Walther (2010), Arias-Castro et al. (2011), Fan et al. (2012), Chan and Walther (2013), and Cai and Yuan (2014), among many others. The problem of colocalization analysis, however, is unique in at least two critical aspects. First, most if not all existing works focus exclusively on signals at the mean level. Our interest here, on the other hand, is on the correlation coefficient. Not only do we want to detect signals in terms of the correlation, but also we want to do so in the presence of unknown marginal means and variances as nuisance parameters. In addition, prior work typically deals with situations where  $\mathbb{I}$  is one-dimensional and  $\mathcal{R}$  is a collection of segments, which is amenable to statistical analysis and sometimes also allows for fast computation. In the case of colocalization analysis, however, the index set  $\mathbb{I}$  is multidimensional and the set  $\mathcal{R}$  usually contains more complex geometric shapes. To address both challenges, we develop in this article a general methodology for correlation detection on a general domain that can be readily applied for colocalization analysis.

Our method is motivated by an observation that, for a fairly general family of  $\mathcal{R}$ , the likelihood ratio statistics exhibit a subtle dependence on the size of a candidate region. As a result, their direct use for correlation detection may lead to nontrivial bias, and substan-

tially reduced power. To overcome this problem, we introduce a size-corrected likelihood ratio statistic and show that scanning with the corrected likelihood ratio statistic yields optimal correlation detection for a large family of  $\mathcal{R}$  in the sense that it can detect elevated correlation at a level no other detectors could improve significantly. We show that the corrected likelihood ratio statistic based scan can also be computed efficiently for a large collection of geometric shapes in arbitrary dimension, characterized by their covering numbers under a suitable semimetric. This includes among others, convex polygons or ellipses, arguably two of the most commonly encountered shapes in practice.

The rest of the paper is organized as follows. In the next section, we introduce a size-corrected likelihood ratio statistic for a general index set  $\mathbb{I}$  and collection  $\mathcal{R}$ , and discuss how it can be used to detect colocalization. We shall also investigate efficient implementation as well as theoretical properties of the proposed method. Section 3 gives several concrete examples of  $\mathbb{I}$  and  $\mathcal{R}$  and show how the general methodology can be applied to these specific situations. Numerical experiments are presented in Section 5 to further illustrate the merits of the proposed methods. Proofs are given in Section 6, with auxiliary results relegated to the Appendix.

## 2 Structured Correlation Detection

In a general correlation detection problem,  $\mathbb{I}$  can be an arbitrary index set and  $\mathcal{R} \subset 2^{\mathbb{I}}$  is a given collection of subsets of  $\mathbb{I}$ , and we are interested in testing the null hypothesis  $H_0$  that (1) holds for all  $i \in \mathbb{I}$  against a composite alternative  $H_a$  that (2) holds for all  $i \in R$  whereas (1) holds for all  $i \notin R$ , for some  $R \in \mathcal{R}$ . We shall argue in this section that the usual maximum log-likelihood ratio statistic may not be suitable for correlation detection, and introduce a size-based correction to address the problem.

### 2.1 Likelihood ratio statistics

A natural test statistic for our purpose is the scan, or maximum log-likelihood ratio statistic:

$$L^* = \max_{R \in \mathcal{R}} L_R,$$

where  $L_R$  is the log-likelihood ratio statistic for testing  $H_0$ :

$$L_R = -(|R| - 2) \log(1 - r_R^2). \quad (3)$$

Here  $|R|$  is the cardinality of  $R$  and  $r_R$  is Pearson correlation within  $R$ :

$$r_R = \frac{\sum_{i \in R} (X_i - \bar{X}_R)(Y_i - \bar{Y}_R)}{\sqrt{\sum_{i \in R} (X_i - \bar{X}_R)^2 \sum_{i \in R} (Y_i - \bar{Y}_R)^2}}$$

where

$$\bar{X}_R = \frac{1}{|R|} \sum_{i \in R} X_i, \quad \text{and} \quad \bar{Y}_R = \frac{1}{|R|} \sum_{i \in R} Y_i.$$

It is worth noting that strictly speaking,  $L_R$  defined by (3) is not the genuine likelihood ratio statistic, which would replace the factor  $|R| - 2$  on the right hand side of (3) by  $|R|$ . Our modification accounts for the correct degrees of freedom so that, for a fixed uncorrelated region  $R$ ,

$$L_R \approx (|R| - 2) \frac{r_R^2}{1 - r_R^2} \sim t_{|R|-2}^2.$$

See, e.g., Muirhead (2008). Obviously, when  $|R|$  is large,  $L_R$  approximately follows a  $\chi_1^2$  distribution and the effect of such correction becomes negligible.

The use of scan or maximum log-likelihood ratio statistics for detecting spatial clusters or signals is very common across a multitude of fields. See, e.g., Fan (1996), Fan et al. (2001) and Glaz et al. (2001) and references therein. Their popularity is also justified as it is well known that scan statistics are minimax optimal if  $|R|$  is small when compared with  $|\mathbb{I}|$ . See, e.g., Lepski and Tsybakov (2000), Dümbgen and Spokoiny (2001), and Dümbgen and Walther (2008). But we show here that such a strategy may not be effective for correlation detection unless  $|R|$  is very small. In particular, we show that, in the absence of a correlated region, the magnitude of  $L_R$  depends critically on its size  $|R|$ , and therefore, the maximum of  $L_R$ 's over regions of different sizes is typically dominated by those evaluated on smaller regions. As a result, direct use of  $L^*$  for correlation detection could be substantially conservative in detecting larger correlated regions.

We now examine the behavior of the maximum of  $L_R$  for  $R \in \mathcal{R}$  of a particular size. Note that it is possible that there is no element in  $\mathcal{R}$  that is of a particular size. To avoid

lengthy discussion to account for such trivial situations, we shall consider instead the subset

$$\mathcal{R}(A) = \{R \in \mathcal{R} : |R| \in (A/2, A]\},$$

for some positive  $A$ . In other words,  $\mathcal{R}(A)$  is the collection of all possible correlated regions of size between  $A/2$  and  $A$ . The factor of  $1/2$  is chosen arbitrarily and can be replaced by any constant in  $(0, 1)$ . Basically  $\mathcal{R}(A)$  includes elements of  $\mathcal{R}$  that, roughly speaking, are of size  $A$ . It is clear that

$$L^* = \max_A \left\{ \max_{R \in \mathcal{R}(A)} L_R \right\}.$$

We shall argue that  $\max_{R \in \mathcal{R}(A)} L_R$  may have different magnitudes for different  $A$ s under the null hypothesis. In particular, we shall show that for a large collection of  $\mathcal{R}(A)$ ,  $\max_{R \in \mathcal{R}(A)} L_R$  can be characterized precisely.

Obviously, the behavior of  $\max_{R \in \mathcal{R}(A)} L_R$  depends on the complexity of  $\mathcal{R}(A)$ . More specifically, we shall first assume that the possible correlated regions are indeed more structured than arbitrary subsets of  $\mathbb{I}$  in that there exist constants  $c_1, c_2 > 0$  independent of  $A$  and  $n := |\mathbb{I}|$  such that

$$|\mathcal{R}(A)| \leq c_1 n A^{c_2}. \tag{4}$$

In other words, (4) dictates that  $|\mathcal{R}(A)|$  increases with  $A$  only polynomially. This is to be contrasted with the completely unstructured setting where  $\mathcal{R} = 2^{\mathbb{I}}$ , the collection of all subsets of  $\mathbb{I}$ , and the number of all subsets of  $\mathbb{I}$  of size  $A$  is of the order  $n^A$ , which depends on  $A$  exponentially. Condition (4) essentially requires that  $\mathcal{R}$  is a much smaller subset of  $2^{\mathbb{I}}$  and therefore indeed imposes structures on the possible regions of correlation.

Naïve counting of the size of  $\mathcal{R}(A)$  as above, however, may not reflect its real complexity. To this end, we also need to characterize the dissimilarity of elements of  $\mathcal{R}(A)$ . For any two sets  $R_1, R_2 \in 2^{\mathbb{I}}$ , write

$$d(R_1, R_2) = 1 - \frac{|R_1 \cap R_2|}{\sqrt{|R_1||R_2|}}.$$

It is easy to see that  $d(\cdot, \cdot)$  is a semimetric on  $2^{\mathbb{I}}$ . We now consider the covering number of sets of a particular size in  $\mathcal{R}$  under  $d$ . Let  $N(A, \epsilon)$  be the smallest integer such that there is

a subset, denoted by  $\mathcal{R}_{\text{app}}(A, \epsilon)$ , of  $\mathcal{R}$  with

$$|\mathcal{R}_{\text{app}}(A, \epsilon)| = N(A, \epsilon)$$

and

$$\sup_{R_1 \in \mathcal{R}(A)} \inf_{R_2 \in \mathcal{R}_{\text{app}}(A, \epsilon)} d(R_1, R_2) \leq \epsilon.$$

It is worth emphasizing that we require the covering set  $\mathcal{R}_{\text{app}}(A, \epsilon) \subset \mathcal{R}$ . It is clear that  $N(A, \epsilon)$  is a decreasing function of  $\epsilon$  and  $N(A, 0) = |\mathcal{R}(A)|$ . We shall also adopt the convention that  $N(A, 1)$  represents the largest number of non-overlapping elements from  $\mathcal{R}(A)$ . Clearly, without any structural assumption, we can always divide  $\mathbb{I}$  into  $n/A$  subsets of size  $A$ . We shall assume that the collection  $\mathcal{R}(A)$  is actually rich enough that

$$N(A, 1) \geq c_3 \frac{n}{A}, \tag{5}$$

for some constant  $c_3 > 0$ . Conversely, we shall assume also that there are not too many “distinct” sets in  $\mathcal{R}(A)$  in that there are certain constants  $c_4, c_5, c_6 > 0$  independent of  $A$  and  $N$  such that

$$N(A, \epsilon) \leq c_4 \frac{n}{A} \left( \log \frac{n}{A} \right)^{c_5} \left( \frac{1}{\epsilon} \right)^{c_6}. \tag{6}$$

Conditions (4), (5) and (6) are fairly general and hold for many common choices of  $\mathcal{R}$ . Consider, for example, the case when  $\mathbb{I} = \{1, 2, \dots, n\}$  is a one-dimensional sequence and

$$\mathcal{R} = \{(a, b] : 0 \leq a < b \leq n\}$$

is the collection of all possible segments on  $\mathbb{I}$ . It is clear that there are at most  $n - \ell$  segments of length  $\ell$  for any  $\ell \in (A/2, A]$ , which means

$$|\mathcal{R}(A)| \leq \frac{1}{2}nA.$$

In addition, for any  $A$ , there are at least  $\lfloor n/A \rfloor$  distinct segments

$$\{((i-1)A, iA] : i = 1, \dots, \lfloor n/A \rfloor\},$$



of length  $A$ , implying that (5) also holds. On the other hand, it is not hard to see that the collection of all segments starting at  $(i - 1)\epsilon A/2$  ( $i = 1, 2, \dots$ ) of length between  $A/2$  and  $A$  can approximate any segment of length between  $A/2$  and  $A$  with approximation error  $\epsilon$ . Therefore,

$$N(A, \epsilon) \leq \left( \frac{A/2}{\epsilon A/2} \right) \left( \frac{n}{\epsilon A/2} \right) = 2 \frac{n}{A} \left( \frac{1}{\epsilon} \right)^2,$$

so that (6) also holds. In the next section, we shall consider more complex examples motivated by colocalization analysis and show that these conditions are expected to hold in fairly general settings.

We now show that if  $\mathcal{R}(A)$  satisfies these conditions,  $\max_{R \in \mathcal{R}(A)} L_R$  concentrates sharply around  $2 \log(n/A)$ .

**Theorem 1.** *Suppose that (1) holds for all  $i \in \mathbb{I}$ . Assume also that (4) and (6) hold. Then*

$$\max_{R \in \mathcal{R}(A)} L_R \leq 2 \log(n/A) + O_p(\log \log(n/A)), \quad \text{as } n \rightarrow \infty. \quad (7)$$

*If in addition, (5) holds, then*

$$\max_{R \in \mathcal{R}(A)} L_R = 2 \log(n/A) + O_p(\log \log(n/A)), \quad \text{as } n \rightarrow \infty. \quad (8)$$

We adopted a generic chaining (see, e.g., Talagrand, 2000) argument for the proof of Theorem 1. The analysis itself may be of independent interest and can be applied to other similar problems such as deriving asymptotic bounds for likelihood ratio statistics in structured detection of mean shifts.

## 2.2 Size-corrected likelihood ratio statistics

An immediate consequence of Theorem 1 is that the value of  $L^*$  alone may not be a good measure of the evidence of correlation. It also depends critically on the size of  $R$  for which  $L_R$  is maximized. As such, when using  $L^*$  as a test statistic, the critical value is largely driven by  $\max_{R \in \mathcal{R}(A)} L_R$  corresponding to smaller  $A$ 's. Therefore, a test based on  $L^*$  could be too conservative when correlation is present on a region with a large cardinality. Motivated by this observation, we now consider normalizing  $\max_{R \in \mathcal{R}(A)} L_R$ , leading to a size-corrected

log-likelihood ratio statistic:

$$\begin{aligned} T^* &= \max_A \left\{ \frac{1}{\log \log(n/A)} \left[ \max_{R \in \mathcal{R}: |R|=A} L_R - 2 \log(n/A) \right] \right\} \\ &= \max_{R \in \mathcal{R}} \left\{ \frac{1}{\log \log(n/|R|)} \left[ L_R - 2 \log \left( \frac{n}{|R|} \right) \right] \right\}. \end{aligned}$$

For brevity, we shall hereafter assume that  $\max_{R \in \mathcal{R}} |R| \leq n/4$ . In general, we can always replace  $\log x$  by  $\log_+(x) := \log(\max\{x, 1\})$  to avoid the trivial cases where the logarithms may not be well defined.

It is clear that under the null hypothesis, the distribution of  $T^*$  is invariant to the nuisance parameters and therefore can be readily evaluated through Monte Carlo simulation. More specifically, one can simulate  $(X_i^*, Y_i^*)^\top \sim N(0, I_2)$  independently for  $i \in \mathbb{I}$ , and compute  $T^*$  for the simulated data. The distribution of  $T^*$  can be approximated by the empirical distribution of the test statistics estimated by repeating this process. Denote by  $q_\alpha$  the  $(1 - \alpha)$ -quantile of  $T^*$  under the null hypothesis. We shall then proceed to reject  $H_0$  if and only if  $T^* > q_\alpha$ . This clearly is an  $\alpha$ -level test by construction. We shall show that in Section 4, it is also a powerful test for detecting correlation.

One of the potential challenges for scan statistics is computation. To compute  $T^*$ , we need to enumerate all elements in  $\mathcal{R}$ , which could be quite burdensome. A key insight obtained from studying  $T^*$  however suggests an alternative to  $T^*$  that is more amenable for computation. More specifically, it is noted that although numerous, regions of large size, namely  $\mathcal{R}(A)$  with a large  $A$ , may have fewer “distinct” elements. As such, we do not need to evaluate  $L_R$  on each  $R \in \mathcal{R}(A)$  but rather on a smaller covering set  $\mathcal{R}_k(A)$ .

With slight abuse of notation, write

$$\mathcal{R}_k = \{R \in \mathcal{R} : |R| \in (2^{-k}n, 2^{-(k-1)}n]\}, \quad k = 2, \dots, \lfloor \log_2 n \rfloor + 1.$$

It is clear that  $T^* = \max_k T_k^*$  where

$$T_k^* = \max_{R \in \mathcal{R}_k} \left\{ \frac{1}{\log \log(n/|R|)} \left[ L_R - 2 \log \left( \frac{n}{|R|} \right) \right] \right\}.$$

It turns out that for

$$k \leq k_* := \lfloor \log_2 n - 2 \log_2 \log n \rfloor,$$

we can approximate  $T_k^*$  very well by scanning through only a small number of  $R$ s from  $\mathcal{R}_k$ . In particular, let  $\tilde{\mathcal{R}}_k$  be a  $1/(4k^2)$  covering set of  $\mathcal{R}_k$  with

$$|\tilde{\mathcal{R}}_k| = N \left( 2^{-(k-1)}n, \frac{1}{4k^2} \right).$$

We shall proceed to approximate  $T_k^*$  by

$$\tilde{T}_k^* = \max_{R \in \tilde{\mathcal{R}}_k} \left\{ \frac{1}{\log \log(n/|R|)} \left[ L_R - 2 \log \left( \frac{n}{|R|} \right) \right] \right\},$$

when  $k \leq k_*$ . Denote by

$$\tilde{T}^* = \max_k \tilde{T}_k^*,$$

where, with slight abuse of notation,  $\tilde{T}_k^* = T_k^*$  for  $k > k_*$ . Instead of using  $T^*$ , we shall now consider  $\tilde{T}^*$  as our test statistic. As before, we can compute the  $1 - \alpha$  quantile  $\tilde{q}_\alpha$  of  $\tilde{T}^*$  under the null hypothesis by Monte Carlo method and proceed to reject  $H_0$  if and only if  $\tilde{T}^* > \tilde{q}_\alpha$ .

Compared with  $T^*$ , the new statistic  $\tilde{T}^*$  is much more computationally friendly. More specifically, under the complexity condition (6), it amounts to computing the corrected likelihood ratio statistic on a total of

$$\begin{aligned} & \sum_{k \leq k_*} N \left( 2^{-(k-1)}n, \frac{1}{4k^2} \right) + \sum_{k > k_*} N(2^{-(k-1)}n, 0) \\ & \leq c_4(\log 2)^{c_5} 4^{c_6} n (\log n)^{c_5+2c_6+1} + c_1 n (\log n)^{2c_2+1} \end{aligned}$$

sets. In other words, the number of size-corrected likelihood ratio statistics we need to evaluate in computing  $\tilde{T}^*$  is linear in  $n$ , up to a certain polynomial of logarithmic factor.

### 3 Correlation Detection on a Lattice

While a general methodology was presented for correlation detection under a generic domain in the previous section, we now examine more specific examples motivated by colocalization analysis in microscopic imaging, and discuss further the operating characteristics of the proposed approach. In particular, we shall focus on correlation detection in a two-dimensional

lattice where  $\mathbb{I} = \{(i, j) : 1 \leq i, j \leq m\}$  so that  $n = m^2$ , for concreteness, although the discussion can be extended straightforwardly to more general situations such as rectangular or higher order lattices.

Most of the imaging tools allow users to visually identify areas of colocalization using variants of the Lasso tool. This allows either a convex polygonal or ellipsoidal region to be selected. Motivated by this, we shall consider specifically in this section the detection of correlation on either an unknown convex polygonal or ellipsoidal region on a two-dimensional lattice. We show that in both cases, the collection  $\mathcal{R}$  of all possible correlated areas satisfies conditions (4), (5) and (6) and therefore the size-corrected scan statistic  $\tilde{T}^*$  can be efficiently computed.

### 3.1 Polygons

We first treat convex  $k$ -polygons. Any  $k$ -polygon can be indexed by its vertices  $\{(a_i, b_i) : 1 \leq i \leq k\}$ , and will therefore be denoted by  $K(\{(a_i, b_i) : 1 \leq i \leq k\})$ . For expositional ease, we focus on the case when the vertices are located on the lattice, although the general case can also be treated with further care. The convexity of a polygon allows us to define its center as  $(\bar{a}, \bar{b})$  where

$$\bar{a} = \frac{1}{k} \sum_{i=1}^k a_i, \quad \text{and} \quad \bar{b} = \frac{1}{k} \sum_{i=1}^k b_i.$$

Denote by

$$r_i = \sqrt{(a_i - \bar{a})^2 + (b_i - \bar{b})^2}$$

the distance from the  $i$ th vertex to the center. To fix ideas, we will focus attention on nearly regular polygons, where  $r_i$ s are of the same order. In this case, the collection of possible correlated regions is:

$$\mathcal{R}_{\text{polygon}}(k, M) = \left\{ K(\{(a_i, b_i) : 1 \leq i \leq k\}) : \max_i r_i / \min_i r_i \leq M \right\}.$$

Recall that

$$\mathcal{R}_{\text{polygon}}(A; k, M) = \{R \in \mathcal{R}_{\text{polygon}}(k, M) : |R| \in (A/2, A]\}.$$

The following result states that (4) holds for  $\mathcal{R}_{\text{polygon}}(k, M)$ .

**Proposition 1.** *There exists a constant  $c > 0$  depending on  $k$  and  $M$  only such that*

$$|\mathcal{R}_{\text{polygon}}(A; k, M)| \leq cnA^k.$$

We now verify (5) for  $\mathcal{R}_{\text{polygon}}(k, M)$ . To this end, we note that any convex  $k$ -polygon can be identified with a minimum bounding circle as shown in Figure 2. Clearly if two polygons intersect, so do their minimum bounding circles. This immediately implies that (5) holds, because we can always place  $\lfloor m/r \rfloor^2$  mutually exclusive circles of radius  $r$  over an  $m \times m$  lattice.

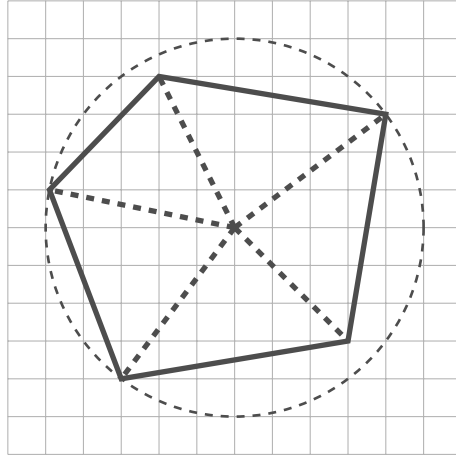


Figure 2: Convex polygon and its minimum bounding circle.

Finally, we show (6) also holds for  $\mathcal{R}_{\text{polygon}}(k, M)$  by constructing an explicit covering set. The idea is fairly simple – we apply a local perturbation to each vertex:

$$\pi_s(K(\{(a_i, b_i) : 1 \leq i \leq k\})) = K(\{(2^s \lfloor 2^{-s} a_i \rfloor, 2^s \lfloor 2^{-s} b_i \rfloor) : 1 \leq i \leq k\}).$$

It can be shown that

**Proposition 2.** *Let  $\pi_s$  be defined above. Then there exists an absolute constant  $c > 0$  such that*

$$\rho(K(\{(a_i, b_i) : 1 \leq i \leq k\}), \pi_s(K(\{(a_i, b_i) : 1 \leq i \leq k\}))) \geq 1 - c(\min_i r_i)^{-1} 2^s.$$

It is clear that, there exist constants  $0 < c_7 < c_8$  depending on  $k$  and  $M$  only such that

$$\mathcal{R}_{\text{polygon}}(A; k, M) \subset \{K \in \mathcal{R}_{\text{polygon}}(k, M) : c_7 A^{1/2} \leq r_i \leq c_8 A^{1/2}, i = 1, 2, \dots, k\}.$$

Therefore, by taking  $s = \log_2(\epsilon A^{1/2})$ , we get

$$N(A, \epsilon) \leq c_9 \frac{n}{A} \left( \log \left( \frac{n}{A} \right) \right)^{k-1} \left( \frac{1}{\epsilon} \right)^{2k+2}$$

In addition, this argument suggests a simple strategy by *digitalization* ( $\pi_s$ ) to construct a covering set for  $\mathcal{R}$ .

From this particular case, we can see the tremendous computational benefit of  $\tilde{T}^*$  over  $T^*$ . To evaluate  $T^*$ , we need to compute the size-corrected likelihood ratio statistics for a total of  $|\mathcal{R}| = O(n^k)$  possible regions. In contrast, computing  $\tilde{T}^*$  only involves  $O(n \text{polylog}(n))$  regions as shown in the previous section. Here  $\text{polylog}(\cdot)$  stands for a certain polynomial of  $\log(\cdot)$ .

## 3.2 Ellipses

Next, we consider the case when  $\mathcal{R}$  is a collection of ellipses on a two-dimensional lattice. Recall that any ellipse can be indexed by its center  $(\tau_1, \tau_2)^\top$ , and a positive definite matrix  $\Sigma \in \mathbb{R}^{2 \times 2}$ :

$$\mathcal{E}((\tau_1, \tau_2)^\top, \Sigma) = \left\{ (x_1, x_2)^\top \in \mathbb{R}^2 : (x_1 - \tau_1, x_2 - \tau_2) \Sigma^{-1} \begin{pmatrix} x_1 - \tau_1 \\ x_2 - \tau_2 \end{pmatrix} \leq 1 \right\}.$$

For brevity, we shall consider the case when  $\Sigma$  is well conditioned in that its condition number, that is the ratio between its eigenvalues, is bounded to avoid lengthy discussion about the effect of discretization. In this case,

$$\mathcal{R}_{\text{ellipse}} = \{ \mathcal{E}((\tau_1, \tau_2)^\top, \Sigma) \cap \mathbb{I} : 1 \leq \tau_1, \tau_2 \leq m, \Sigma \succ 0, \lambda_{\max}(\Sigma)/\lambda_{\min}(\Sigma) \leq M \}.$$

We first note that any ellipse can be identified with its circumscribing rectangle as shown in Figure 3. Therefore, immediately following the bound on the number of rectangles on a

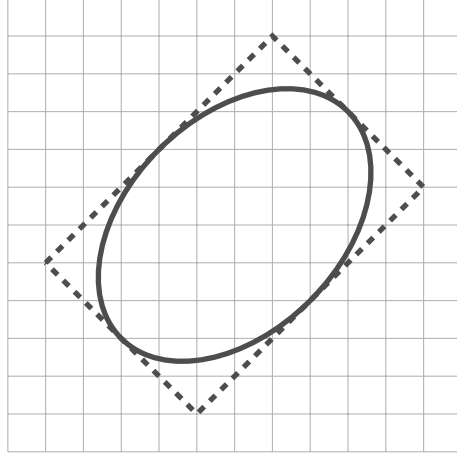


Figure 3: Circumscribing rectangle of an ellipse

lattice, for example by Proposition 1 with  $k = 4$ , we get

$$\mathcal{R}_{\text{ellipse}} \leq cnA^4,$$

for some constant  $c > 0$ . Similarly, if two ellipses intersect, then so do their minimum bounding rectangles. By the argument for polygons, we therefore know that (5) and (6) also hold for  $\mathcal{R}_{\text{ellipse}}$ .

## 4 Optimality

We now study the power of the proposed test  $T^*$  and its variant  $\tilde{T}^*$ . We shall first investigate the required strength of correlation so it can be detected using the proposed tests.

**Theorem 2.** *Assume that (4) and (6) hold. If there exists a correlated region  $R \in \mathcal{R}$ , with  $|R| \rightarrow \infty$ , such that (1) holds for  $i \notin R$  and (2) holds for  $i \in R$ , and*

$$|R| \log \left( \frac{1}{1 - \rho^2} \right) \geq (2 + \delta_n) \log \left( \frac{n}{|R|} \right) \quad (9)$$

*for some  $\delta_n > 0$  such that  $\delta_n \log^{1/2}(n/|R|) \rightarrow \infty$ , then  $T^* > q_\alpha$  and  $\tilde{T}^* > \tilde{q}_\alpha$  with probability tending to one as  $n \rightarrow \infty$ .*

Theorem 2 shows that whenever correlation on a region  $R$  satisfies (9), our tests will con-

sistently reject the null hypothesis and have power tending to one. The detection boundary of the proposed tests for a correlated region  $R$  can therefore be characterized by (9). More specifically, depending on the cardinality  $|R|$ , there are three different regimes.

- For large regions where  $|R| \asymp n$ , correlation is detectable if  $|R|\rho^2 \rightarrow \infty$ . Recall that, from Neyman-Pearson Lemma, even if the correlated region  $R$  is known in advance, we can only consistently detect it under the same requirement. Put differently, the proposed method is as powerful as if we knew the region in advance.
- For regions of intermediate sizes such that  $\log n \ll |R| \ll n$ , the detection boundary becomes  $\rho^2 \geq (2 + \delta_n)|R|^{-1} \log(n/|R|)$ , provided that  $\delta_n \sqrt{\log(n/|R|)} \rightarrow \infty$ . Here, we can see that weaker correlation can be detected over larger regions.
- And finally for small regions where  $|R| \ll \log(n)$ , detection is only possible for nearly perfect correlation in that  $\rho^2 \geq 1 - \exp(-(2 + \delta_n) \log(n)/|R|)$  where  $\delta_n \sqrt{\log n} \rightarrow \infty$ .

It turns out that the detection boundary achieved by  $T^*$  and  $\tilde{T}^*$  as shown in Theorem 2 is indeed sharply optimal.

**Theorem 3.** *Assume that (5) holds. For any  $\alpha$ -level test  $\Delta$ , there exists an instance where correlation occurs on some  $R \in \mathcal{R}$  obeying*

$$|R| \log \left( \frac{1}{1 - \rho^2} \right) \geq (2 - \delta_n) \log \left( \frac{n}{|R|} \right) \quad (10)$$

for a certain  $\delta_n > 0$  with  $\delta_n \log^{1/2}(n/|R|) \rightarrow \infty$ , such that the type II error of  $\Delta$  converges to  $1 - \alpha$  as  $n \rightarrow \infty$ . Moreover, if there exists some  $\alpha$ -level test  $\Delta$  for which the type II error converges to 0 as  $n \rightarrow \infty$  on any instance where correlation occurs on some  $R \in \mathcal{R}$  obeying

$$|R| \log \left( \frac{1}{1 - \rho^2} \right) \geq c_n \quad \text{and} \quad |R| \rightarrow \infty, \quad (11)$$

then it is necessary to have  $c_n \rightarrow \infty$  as  $n \rightarrow \infty$ .

In other words, Theorem 3 shows that any test is essentially powerless for detecting correlation with

$$|R| \log \left( \frac{1}{1 - \rho^2} \right) \leq (2 - \delta_n) \log \left( \frac{n}{|R|} \right)$$



for any  $\delta_n > 0$  such that  $\delta_n \log^{1/2}(n/|R|) \rightarrow \infty$ . Together with Theorem 2, we see that, when  $n/|R| \rightarrow \infty$ , the optimal detection boundary for colocalization for a general index set  $\mathbb{I}$  and a large collection of  $\mathcal{R}$ 's that satisfy certain complexity requirements is

$$|R| \log \left( \frac{1}{1 - \rho^2} \right) = 2 \log \left( \frac{n}{|R|} \right);$$

and the size-corrected scan statistic is sharply optimal.

The second statement of Theorem 3 deals with the case when  $\limsup n/|R|$  is finite. Together with Theorem 2, (11) implies that in this case, the correlated region can be detected if and only if

$$\rho^2 |R| \rightarrow \infty$$

and size-corrected scan statistic is again optimal.

To better appreciate the effect of the size of a correlated region on its detectability, it is instructive to consider the cases where  $|R| = n^\alpha$  for some  $0 < \alpha < 1$  or  $|R| = (\log n)^\alpha$  for some  $\alpha > 1$ . In the former case when  $|R| = n^\alpha$ , the detection boundary is

$$\rho^2 = 2(1 - \alpha)n^{-\alpha} \log n.$$

In the latter case when  $|R| = (\log n)^\alpha$ , the detection boundary is

$$\rho^2 = 2(\log n)^{\alpha-1}.$$

In both cases, it is clear that much weaker correlation can be detected on larger regions.

## 5 Numerical Experiments

We now conduct numerical experiments to further demonstrate the practical merits of the proposed methodology.

### 5.1 Simulation

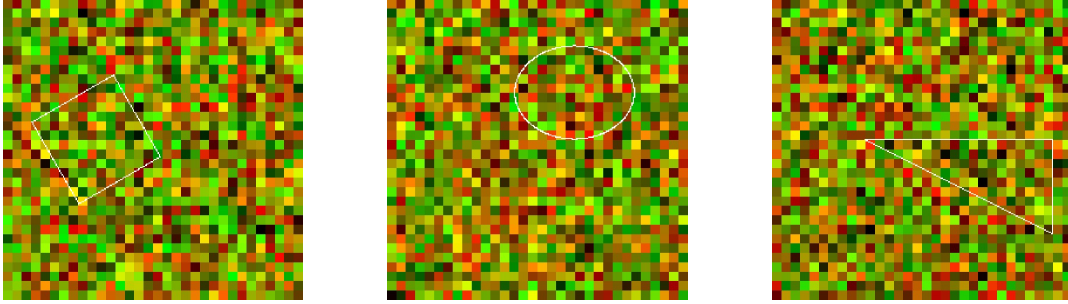
We begin with a series of four sets of simulation studies. To fix ideas, we focus on two dimensional lattices in our simulation studies. The first set of simulations was designed

to show the flexibility of the general method by considering a variety of different shapes of correlated regions, namely the choice of the library  $\mathcal{R}$ , including axis-aligned rectangles, triangles and axis-aligned ellipses. We compare the performance of size-corrected likelihood ratio statistic and the uncorrected likelihood ratio statistic to demonstrate the necessity and usefulness of the proposed correction. The second set was carried out to compare the full scan statistic  $T^*$  and the nearly linear time scan  $\tilde{T}^*$  and illustrate similar performance between the two methods yet considerable computation gain by using  $\tilde{T}^*$ . The third and fourth sets of simulation studies were conducted to confirm qualitatively our theoretical findings about the effect of the size  $|\mathbb{I}|$  of the lattice and the area  $A$  of correlated region on its detectability. In each case, we shall assume that only the shape of the correlated region is known and therefore  $\mathcal{R}$  is the collection of all regions of a particular shape. In addition, we simulate the null distribution and identify the upper 5% quantile of the null distribution based on 1000 Monte Carlo simulations. We reject the null hypothesis for a simulation run if the corresponding test statistic,  $T^*$ ,  $\tilde{T}^*$  or  $L^*$ , exceeds their respective upper quantile. This ensures that each test is at level 5%, up to Monte Carlo simulation error.

As argued in the previous sections, our methods can handle a variety of geometric shapes. We now demonstrate this versatility through simulation where we consider detecting a correlated region in the form of a triangle, an ellipse or a rectangle. In particular, we simulated data on a  $32 \times 32$  squared lattice. Correlation was imposed on a right triangle with side length 10, 20 and  $10\sqrt{5}$ , or an axis-aligned ellipse with short axis 4.94 and long axis 6.36, or a rectangle of size  $10 \times 10$ . The location of these correlated regions was selected uniformly over the lattice. Typical simulated examples of different correlation and shape are given in Figure 4.

To assess the power of  $T^*$ , we considered two relatively small values of correlation coefficient  $\rho$ : 0.2 and 0.4. For comparison purposes, we computed for each simulation run both  $T^*$  and the uncorrected maximum likelihood ratio statistic  $L^*$ . The experiment was repeated for 500 times for each combination of shape and correlation coefficient. The results are summarized in Table 1. These results not only show the general applicability of our methods but also demonstrate the improved power of the size correction we apply.

We now compare the full scan statistic  $T^*$  with its more computationally efficient variant  $\tilde{T}^*$ . We focus on the case when the correlated region is known to be an axis-aligned rectangle.



(a) Rectangle with  $\rho = 0.2$ . (b) Ellipse with  $\rho = 0.2$ . (c) Triangle with  $\rho = 0.2$ .

Figure 4: Simulated examples: overlaid images from green channel and red channel in a typical simulation run for different shapes of correlated regions.

Shape	Rectangle		Ellipse		Triangle	
$\rho$	0.2	0.4	0.2	0.4	0.2	0.4
$T^*$	0.16	0.42	0.25	0.6	0.21	0.58
$L^*$	0.04	0.20	0.03	0.51	0.03	0.26

Table 1: Power comparison between  $T^*$  and  $L^*$  for different combinations of shape and correlation coefficient.

The true correlated region is a randomly selected  $10 \times 10$  rectangle on a  $64 \times 64$  lattice. We consider a variety of different correlation coefficients 0.2, 0.4, 0.6, and 0.8. The performance and computing time of both tests are reported in Table 2, which is also based on 500 runs for each value of the correlation coefficient. It is clear from Table 2 that the two tests enjoy similar performance with  $T^*$  slightly more powerful. Yet  $\tilde{T}^*$  is much more efficient to evaluate as expected.

Correlation Coefficient		0.2	0.4	0.6	0.8
Power	$T^*$	0.108	0.228	0.502	0.708
	$\tilde{T}^*$	0.106	0.214	0.410	0.606
Time (ms)	$T^*$	444.084	447.236	452.634	453.064
	$\tilde{T}^*$	139.026	139.344	140.554	142.144

Table 2: Comparison between  $T^*$  and  $\tilde{T}^*$ .

We note that the computing gain of  $\tilde{T}^*$  over  $T^*$  becomes more significant for larger images. In particular, we ran similar scans over lattices of size  $256 \times 256$ ,  $256 \times 512$  and  $512 \times 512$ . The computing time for a typical dataset is presented in Table 3:

Size of Lattice	$256 \times 256$	$256 \times 512$	$512 \times 512$
Computing time of $T^*$ (s)	129.942	487.238	1934.996
Computing time of $\tilde{T}^*$ (s)	16.59	45.117	144.206

Table 3: Comparison of computing times for  $T^*$  and  $\tilde{T}^*$ .

We now evaluate the effect of the size of a correlated region on its detectability. In the light of the observations made in the previous set of experiments, we focus on using  $\tilde{T}^*$  to detect a correlated rectangle on a  $64 \times 64$  lattice. We consider four different sizes of the correlated rectangle:  $5 \times 5$ ,  $10 \times 10$ ,  $20 \times 20$  and  $40 \times 40$ . For each given size of the correlated region, we varied the correlation coefficient to capture the relationship between the power of our detection scheme and  $\rho$ . The results summarized in Figure 5 are again based on 500 runs for each combination of size and correlation coefficient of the correlated region. The observed effect of  $A$  on its detectability is consistent with the results established in Theorem 2 and Theorem 3: larger regions are easier to detect with the same correlation coefficient.

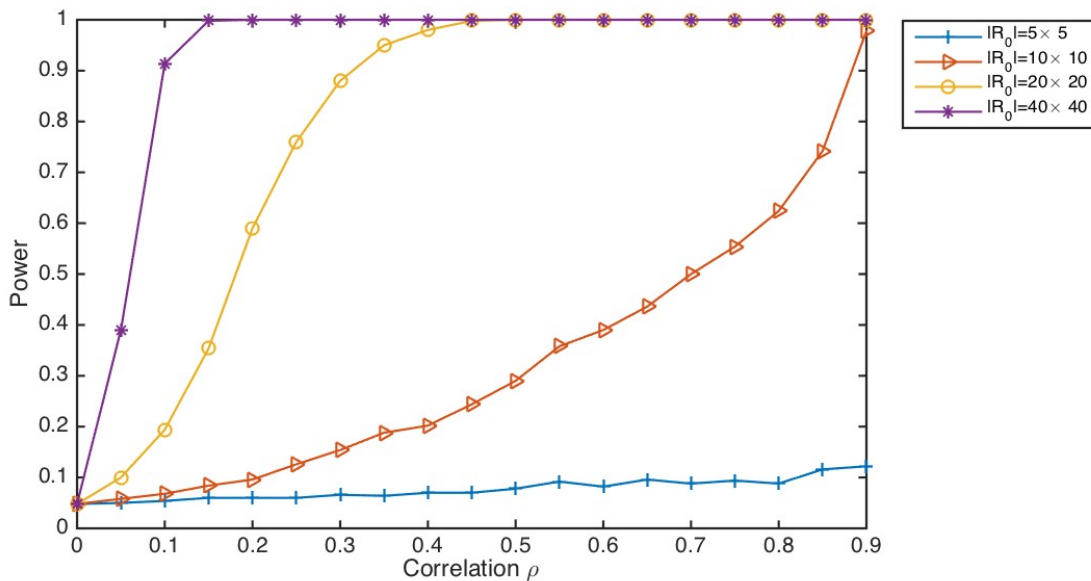


Figure 5: Power plot for detecting a correlated rectangle of different sizes on a  $64 \times 64$  lattice.

Our final set of simulations is designed to assess the effect of  $\mathbb{I}$ . To this end we consider identifying a  $10 \times 10$  correlated rectangle on a squared lattice of size  $32 \times 32$ ,  $64 \times 64$ , or  $128 \times 128$ . As in the previous example, we repeat the experiment 500 times for each

combination of  $\mathbb{I}$  and a variety of values of  $\rho$ . The results are presented in Figure 6. The observed effect of  $|\mathbb{I}|$  is again consistent with our theoretical developments: as the size of lattice increases, detection becomes harder for a region of the same size and correlation.

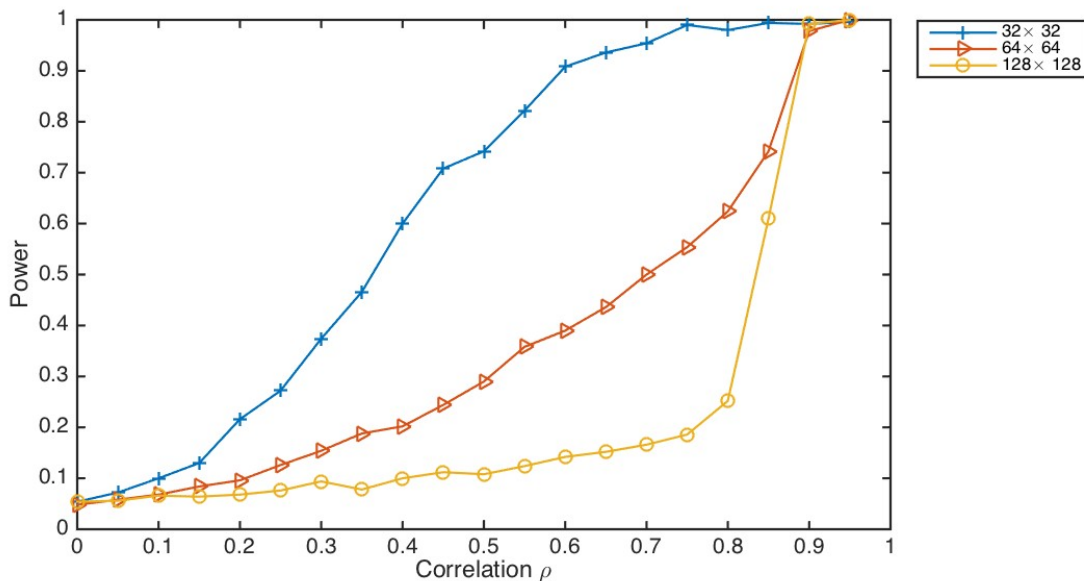


Figure 6: Power plot for detecting a  $10 \times 10$  correlated rectangle on squared lattices of different sizes.

## 5.2 Real data example

For illustration purposes, we now return to the image example we briefly mentioned in the introduction. This image originated from a concerted effort by several research groups to dissect the post-transcriptional process of human immunodeficiency virus type 1 (HIV-1) using imaging based approaches.

HIV uses the cell's mRNA nuclear export pathway to initiate the post-transcriptional stages of the viral life cycle. Nuclear export of a segment of the viral gRNA bearing the Rev Response Element (RRE) is essential to HIV gene regulation, viral replication and disease. It is well established that this process is regulated by a viral Rev trafficking protein that binds to the RRE and recruits the cellular CRM1 nuclear export receptor. In addition to nuclear export, Rev may also play a role in gRNA encapsidation and translation of gRNA. To gain insight into Rev's motility in the nucleus and cytoplasm to better understand Rev

trafficking dynamics, as well as Rev’s roles in viral gene expression and virus particle assembly live cell imaging was employed to monitor Rev and CRM1 behavior. It is expected that the REV/CRM1/RRE ribonuclear complex will have high colocalization in the cytoplasm during the viral life cycle. It has been shown that HIV-1 genomic RNAs (gRNAs) frequently exhibit “burst” nuclear export kinetics. These events are characterized by striking en masse evacuations of gRNAs from the nucleus to flood the cytoplasm in conjunction with the onset of viral late gene expression. Burst nuclear export is regulated through interactions between the viral protein Rev, cellular nuclear export factor CRM1, and the gRNA’s cis-acting RNA Rev response element (RRE). By monitoring mutant versions of the Rev protein unable to bind CRM1, export element deficient versions of the gRNA (RRE Minus), and lack of gRNA in the visual system, we can determine CRM1 trafficking behavior in the context of the virus.

A specific data example is given in Figure 7 where dual-channel images of a wild type cell and a mutant cell are presented side-by-side. In each image, CRM1 is represented by green and the gRNA by red. While the “burst” nuclear export is visible for the wild type cell, it is less evident for the mutant cell. To further quantify such differences between the two cell types, we applied our method to this particular example following standard steps to preprocess the image: applying Otsu’s method to each channel to remove background and then identifying spatial compartments where both channels are significantly expressed. On the post-processed images, we compute the test statistic  $\tilde{T}^*$  and evaluate its corresponding p-value again by simulating the null distribution through 1000 Monte Carlo experiments. For the wild type cell, we obtained  $\tilde{T}^* = 5.65 \times 10^3$  which is larger than any of the 1000 values from the Monte Carlo simulations under the null hypothesis, suggesting a p-value < 0.1%, up to a Monte Carlo simulation error. On the other hand, the test statistic for the mutant cell is 8.98, which corresponds to a p-value of 0.846.

## 6 Proofs of Main Results

We now present the proofs to our main results, namely Theorems 1, 2 and 3. Proofs of Propositions 1 and 2, as well as a number of auxiliary results, are relegated to the Appendix. To distinguish from the constants appeared in the previous sections, we shall use the capital letter  $C$  to denote a generic positive constant that may take different values at each appearance.

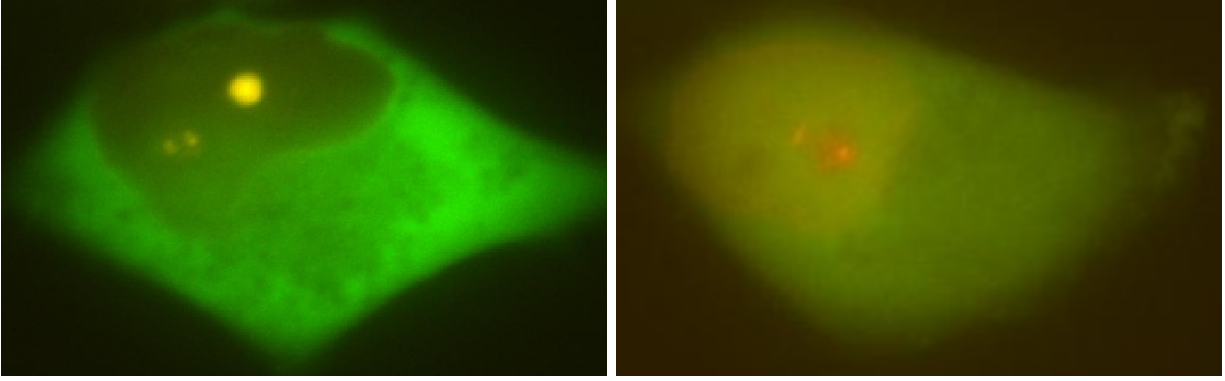


Figure 7: Colocalization between CRM1 and gRNA: comparison between the wild type (on the left) and mutant (on the right).

*Proof of Theorem 1.* We first prove the upper bound (7) under conditions (4) and (6). To this end, we shall establish a stronger result that there exists a constant  $C > 0$  such that for any  $0 < t < (\log n)^3$ .

$$\mathbb{P} \left\{ \max_{R \in \mathcal{R}(A)} L_R > 2 \log n + C(\log \log n + t) \right\} \leq \exp(-t). \quad (12)$$

It is clear that (7) follows immediately from (12).

We now proceed to prove (12). We shall consider the cases where  $A \leq (\log n)^5$  and  $A \geq (\log n)^5$  separately. First consider the situation when  $A \leq (\log n)^5$ . By Lemma 6, there exists a constant  $C > 0$  such that for any fixed  $R \in \mathcal{R}(A)$

$$\mathbb{P} \{L_R > x\} \leq C \exp(-x/2).$$

Applying union bound yields

$$\begin{aligned} \mathbb{P} \left\{ \max_{R \in \mathcal{R}(A)} L_R > x \right\} &\leq C |\mathcal{R}(A)| \exp(-x/2) \\ &\leq c_1 C n A^{c_2} \exp(-x/2) \\ &\leq c_1 C n (\log n)^{5c_2} \exp(-x/2), \end{aligned}$$

where the second inequality follows from (4). Equation (12) then follows by taking

$$x = 2 \log(c_1 C) + 2 \log n + 10c_2 \log \log n + 2t.$$

The treatment for  $A \geq (\log n)^5$  is more involved and we apply a chaining argument. Let  $\mathcal{R}_{\text{app}}(A, e^{-s})$  be an  $e^{-s}$  covering set of  $\mathcal{R}(A)$  so that

$$|\mathcal{R}_{\text{app}}(A, e^{-s})| = N(A, e^{-s}).$$

For any segment  $R \in \mathcal{R}(A)$ , denote by

$$\pi_s(R) = \operatorname{argmin}_{R' \in \mathcal{R}_{\text{app}}(A, e^{-s})} d(R, R').$$

Of course, the minimizer on the right hand side may not be uniquely defined, in which case, we take  $\pi_s(R)$  to be an arbitrarily chosen minimizer.

Write

$$L_R = \sum_{s=s_*}^{s^*-1} (L_{\pi_{s+1}(R)} - L_{\pi_s(R)}) + (L_R - L_{\pi_{s^*}(R)}) + L_{\pi_{s^*}(R)},$$

where  $s^* > s_* \geq \log \log(n/A)$  are to be specified later. It is clear that

$$\max_{R \in \mathcal{R}(A)} L_R \leq \sum_{s=s_*}^{s^*-1} \max_{R \in \mathcal{R}(A)} |L_{\pi_{s+1}(R)} - L_{\pi_s(R)}| + \max_{R \in \mathcal{R}(A)} |L_R - L_{\pi_{s^*}(R)}| + \max_{R \in \mathcal{R}(A)} |L_{\pi_{s^*}(R)}|. \quad (13)$$

We now bound the three terms on the right hand side of (13) separately.

By definition,

$$d(R, \pi_s(R)) \leq e^{-s}, \quad \text{and} \quad d(R, \pi_{s+1}(R)) \leq e^{-(s+1)}.$$

Hence there exists a constant  $C > 0$  such that

$$|\pi_s(R) \cap \pi_{s+1}(R)| \geq (1 - Ce^{-s})|R|, \quad \text{and} \quad d(\pi_s(R), \pi_{s+1}(R)) \leq Ce^{-s}.$$

Now by Lemma 7, for any fixed  $R \in \mathcal{R}(A)$ ,

$$|L_{\pi_s(R)} - L_{\pi_{s+1}(R)}| \leq C (e^{-s/2}x + |R|^{-1/2}x^{3/2})$$



with probability at least  $1 - Ce^{-x}$ . An application of the union bound yields

$$\begin{aligned}
& \mathbb{P} \left\{ \max_{R \in \mathcal{R}(A)} |L_{\pi_{s+1}(R)} - L_{\pi_s(R)}| > C \left( e^{-s/2}x + \sqrt{2}A^{-1/2}x^{3/2} \right) \right\} \\
& \leq CN(A, e^{-s})N(A, e^{-(s+1)})e^{-x} \\
& \leq C[N(A, e^{-(s+1)})]^2 e^{-x} \\
& \leq c_4^2 C \left( \frac{n}{A} \right)^2 \left( \log \frac{n}{A} \right)^{2c_5} e^{2c_6(s+1)} e^{-x},
\end{aligned}$$

where the last inequality follows from (6). In particular, taking

$$x = t + 2 \log s + \log(c_4^2 C) + 2 \log(n/A) + 2c_5 \log \log(n/A) + 2c_6(s+1)$$

yields, with probability at least  $1 - s^{-2}e^{-t}$ ,

$$\max_{R \in \mathcal{R}(A)} |L_{\pi_{s+1}(R)} - L_{\pi_s(R)}| \leq C \left( (s+t+\log(n/A))e^{-s/2} + A^{-1/2}(s+t+\log(n/A))^{3/2} \right).$$

Here we used the fact that  $s \geq s_* \geq \log \log(n/A)$ . Now applying the union bound over all  $s_* \leq s < s^*$ , we get, with probability at least  $1 - s_*^{-1}e^{-t} \geq 1 - e^{-t}$ ,

$$\begin{aligned}
\sum_{s=s_*}^{s^*-1} \max_{R \in \mathcal{R}(A)} |L_{\pi_{s+1}(R)} - L_{\pi_s(R)}| & \leq C \sum_{s=s_*}^{s^*-1} \left( (s+t+\log(n/A))e^{-s/2} + A^{-1/2}(s+t+\log(n/A))^{3/2} \right) \\
& \leq C \left( s_* e^{-s_*/2} + A^{-1/2}(s^*)^{5/2} \right) \\
& \quad + C \left( e^{-s_*/2}(t+\log(n/A)) + A^{-1/2}s^*(t+\log(n/A))^{3/2} \right).
\end{aligned}$$

To bound the second term on the right hand side of (13), we again apply Lemma 7. For any fixed  $R \in \mathcal{R}(A)$ , we get

$$\mathbb{P} \left\{ |L_R - L_{\pi_{s^*}(R)}| \geq C \left( e^{-s^*/2}x + \sqrt{2}A^{-1/2}x^{3/2} \right) \right\} \leq Ce^{-x}.$$

Another application of the union bound yields,

$$\begin{aligned}
\max_{R \in \mathcal{R}(A)} |L_R - L_{\pi_{s^*}(R)}| & \leq C \left( e^{-s^*/2} \log |\mathcal{R}(A)| + A^{-1/2}(\log |\mathcal{R}(A)|)^{3/2} + e^{-s^*/2}t + A^{-1/2}t^{3/2} \right) \\
& \leq c_2 C \left( e^{-s^*/2} \log n + A^{-1/2}(\log n)^{3/2} + e^{-s^*/2}t + A^{-1/2}t^{3/2} \right),
\end{aligned}$$

with probability at least  $1 - Ce^{-t}$ , where we used (4) in the last inequality.

Finally, for the third term on the right hand side of (13), we have

$$\mathbb{P} \left\{ \max_{R \in \mathcal{R}_{\text{app}}(A, e^{-s_*})} |L_R| \geq x \right\} \leq CN(A, e^{-s_*}) e^{-x/2} \leq c_4 C \left( \frac{n}{A} \right) \left( \log \frac{n}{A} \right)^{c_5} e^{c_6 s_*} e^{-x/2}.$$

Taking

$$x = 2 \log(c_4 C) + 2 \log \frac{n}{A} + c_5 \log \log \frac{n}{A} + 2c_6 s_* + t$$

yields, with probability at least  $1 - Ce^{-t}$ ,

$$\max_{R \in \mathcal{R}_{\text{app}}(A, e^{-s_*})} |L_R| \leq 2 \log(c_4 C) + 2 \log \frac{n}{A} + c_5 \log \log \frac{n}{A} + 2c_6 s_* + t.$$

In summary, we get, with probability at least  $1 - Ce^{-t}$ ,

$$\begin{aligned} \max_{R \in \mathcal{R}(A)} L_R &\leq C \left( s_* e^{-s_*/2} + A^{-1/2} (s_*)^{5/2} + e^{-s_*/2} t + A^{-1/2} s_* t^{3/2} + e^{-s_*/2} \log n \right. \\ &\quad \left. + A^{-1/2} (\log n)^{3/2} + e^{-s_*/2} t + A^{-1/2} t^{3/2} + e^{-s_*/2} \log \frac{n}{A} + A^{-1/2} s_* (\log(n/A))^{3/2} \right) \\ &\quad + 2 \log(c_4 C) + 2 \log \frac{n}{A} + c_5 \log \log \frac{n}{A} + 2c_6 s_* + t. \end{aligned}$$

Recall that  $A \geq (\log n)^5$ . If we take  $s^* = 2 \log n$  and  $s_* = 2 \log \log(n/A)$ , then for any  $t \leq (\log n)^3$ , we can deduce from the above inequality that

$$\max_{R \in \mathcal{R}(A)} L_R \leq 2 \log \frac{n}{A} + C \left( \log \log \frac{n}{A} + t \right), \quad (14)$$

which implies (12).

We now prove (8) if in addition, (5) holds. In the light of (6), we can find a subset  $\tilde{\mathcal{R}}(A)$  of  $\mathcal{R}(A)$  such that for any  $R_1, R_2 \in \tilde{\mathcal{R}}(A)$ ,  $R_1 \cap R_2 = \emptyset$  and

$$|\tilde{\mathcal{R}}(A)| \geq c_3 \frac{n}{A}.$$

Obviously,

$$\max_{R \in \mathcal{R}(A)} L_R \geq \max_{R \in \tilde{\mathcal{R}}(A)} L_R.$$

If  $A \leq (\log n)^5$ , then

$$\begin{aligned}
\mathbb{P} \left\{ \max_{R \in \tilde{\mathcal{R}}(A)} L_R \leq x \right\} &= \prod_{R \in \tilde{\mathcal{R}}(A)} \mathbb{P}\{L_R \leq x\} \\
&\leq \prod_{R \in \tilde{\mathcal{R}}(A)} (1 - C|R|^{-1/2}e^{-x/2}) \\
&\leq [1 - CA^{-1/2}e^{-x/2}]^{c_3n/A} \\
&\leq [1 - C(\log n)^{-5/2}e^{-x/2}]^{c_3n/A},
\end{aligned}$$

where the first inequality follows from the lower bound given by Lemma 6. It can then be derived that

$$\max_{R \in \tilde{\mathcal{R}}(A)} L_R \geq 2 \log n + O_p(\log \log n). \quad (15)$$

Together with (7), (15) implies the desired claim when  $A \leq (\log n)^5$ .

Next we consider the case when  $A \geq (\log n)^5$ . We proceed in a similar fashion as before but rely on the following tail bound of  $L_R$ : if  $A \geq 24$ , then there exists a constant  $C > 0$  such that for any  $R \in \mathcal{R}(A)$  and  $0 < x < \sqrt{A}$ ,

$$\mathbb{P}\{L_R \geq x\} \leq Cx^{-1/2} \exp(-x/2). \quad (16)$$

If (16) holds, then

$$\mathbb{P} \left\{ \max_{R \in \tilde{\mathcal{R}}(A)} L_R \leq x \right\} \geq (1 - Cx^{-1/2}e^{-x/2})^{c_3n/A},$$

which yields

$$\max_{R \in \tilde{\mathcal{R}}(A)} L_R \geq \max_{R \in \tilde{\mathcal{R}}(A)} L_R \geq 2 \log(n/A) + O_p(\log \log(n/A)).$$

Together with (7), this concludes the proof.

It now remains to prove (16). Write

$$T_R = (|R| - 2)r_R^2.$$

Note that  $\log(1+x) > x - x^2/2$  for any  $x > 0$ . We get

$$L_R \geq (|R| - 2) \log(1 + r_R^2) \geq T_R^2 - \frac{T_R^4}{2(|R| - 2)} \geq T_R^2 - \frac{T_R^4}{A - 4},$$

for any  $A \geq 5$ , where in the last inequality we used the fact that  $|R| > A/2$  for any  $R \in \mathcal{R}(A)$ . This can be further lower bounded by  $T_R^2 - 3T_R^4/A$  for any  $A \geq 6$ . Thus, for any  $0 < x < A/24$ ,

$$\begin{aligned} \mathbb{P}\{L_R \geq x\} &\geq \mathbb{P}\left\{T_R^2 - \frac{3T_R^4}{A} \geq x\right\} \\ &\geq \mathbb{P}\left\{T_R^2 - \frac{3T_R^4}{A} \in [x, 2x)\right\} \\ &\geq \mathbb{P}\{T_R^2 \in [x + 12x^2/A, 2x + 3x^2/A)\} \\ &\geq \mathbb{P}\left\{T_R \in [\sqrt{x + 12x^2/A}, \sqrt{2x + 3x^2/A})\right\}. \end{aligned}$$

Because  $T_R \sim t_{|R|-2}$ , we have

$$\begin{aligned} &\mathbb{P}\left\{T_R \in [\sqrt{x + 12x^2/A}, \sqrt{2x + 3x^2/A})\right\} \\ &\geq C \int_{\sqrt{x+12x^2/A}}^{\sqrt{2x+3x^2/A}} \left(1 + \frac{u^2}{|R|-2}\right)^{-\frac{|R|-1}{2}} du \\ &\geq C \int_{\sqrt{x+12x^2/A}}^{\sqrt{2x+3x^2/A}} \exp\left[-\frac{|R|-1}{2} \log\left(1 + \frac{u^2}{|R|-2}\right)\right] du \\ &\geq C \int_{\sqrt{x+12x^2/A}}^{\sqrt{2x+3x^2/A}} \exp\left(-\frac{|R|-1}{2(|R|-2)} u^2\right) du, \end{aligned}$$

for some constant  $C > 0$ , where in the last inequality we used the fact that  $\log(1+x) \leq x$

for all  $x \geq 0$ . Thus,

$$\begin{aligned}
& \mathbb{P} \left\{ T_R \in [\sqrt{x + 12x^2/A}, \sqrt{2x + 3x^2/A}] \right\} \\
& \geq C(2x + 3x^2/A)^{-1/2} \int_{\sqrt{x+12x^2/A}}^{\sqrt{2x+3x^2/A}} u \exp \left( -\frac{|R| - 1}{2(|R| - 2)} u^2 \right) du \\
& = C(2x + 3x^2/A)^{-1/2} (1 - (|R| - 1)^{-1}) \left[ \exp \left( -\frac{|R| - 1}{2(|R| - 2)} (x + 12x^2/A) \right) \right. \\
& \quad \left. - \exp \left( -\frac{|R| - 1}{2(|R| - 2)} (2x + 3x^2/A) \right) \right].
\end{aligned}$$

Recall that  $0 < x < A/24$ . We get

$$\begin{aligned}
& \mathbb{P} \left\{ T_R \in [\sqrt{x + 12x^2/A}, \sqrt{2x + 3x^2/A}] \right\} \\
& \geq Cx^{-1/2} \exp \left( -\frac{|R| - 1}{2(|R| - 2)} (x + 12x^2/A) \right) \\
& \geq Cx^{-1/2} \exp \left( -\frac{A - 1}{2(A - 2)} (x + 12x^2/A) \right) \\
& \geq Cx^{-1/2} \exp(-x/2),
\end{aligned}$$

where in the last inequality, we used the fact that  $x \leq \sqrt{A}$ . The proof is then completed.  $\square$

*Proof of Theorem 2 (Consistency of  $T^*$ ).* We first show that the claim is true for  $T^*$ . To this end, we begin by arguing that  $q_\alpha = O(1)$ , and then show that under  $H_1$ ,  $T^* \rightarrow \infty$ . Note that

$$\begin{aligned}
T^* & = \max_{R \in \mathcal{R}} \left\{ \frac{1}{\log \log(n/|R|)} \left[ L_R - 2 \log \left( \frac{n}{|R|} \right) \right] \right\} \\
& = \max_{1 \leq k \leq \log n} \max_{R \in \mathcal{R}(e^{-k+1}n)} \left\{ \frac{1}{\log \log(n/|R|)} \left[ L_R - 2 \log \left( \frac{n}{|R|} \right) \right] \right\}.
\end{aligned}$$

As shown in the proof of Theorem 1, there exists a constant  $C > 0$  such that for any  $0 < t < (\log n)^3$ ,

$$\mathbb{P} \left\{ \max_{R \in \mathcal{R}(e^{-k+1}n)} L_R \geq 2k + C(\log k + t) \right\} \leq \exp(-t).$$

Taking  $t = x + \log(2k^2)$  yields

$$\mathbb{P} \left\{ \max_{R \in \mathcal{R}(e^{-k+1}n)} \left\{ \frac{1}{\log \log(n/|R|)} \left[ L_R - 2 \log \left( \frac{n}{|R|} \right) \right] \right\} \geq C(x+1) \right\} \leq \frac{1}{2k^2} \exp(-x).$$

Applying union bound over all  $k$ , we get

$$\mathbb{P} \{T^* \geq C(x+1)\} \leq \sum_{1 \leq k \leq \log n} \frac{1}{2k^2} \exp(-x) \leq \exp(-x),$$

which implies that  $q_\alpha \leq C(1 - \log(1 - \alpha))$ .

It now suffices to show that if (9) holds for some  $R \in \mathcal{R}$ , then  $T^* \rightarrow \infty$ . To this end, note that

$$T^* \geq \frac{1}{\log \log(n/|R|)} \left[ L_R - 2 \log \left( \frac{n}{|R|} \right) \right].$$

We treat the case  $|R| \geq \log n$  and  $|R| \leq \log n$  separately.

Consider first the situation when  $|R| \leq \log n$ . By Lemma 3,

$$\left( \frac{1 - r_R^2}{1 - \rho^2} \right) \sum_{i \in R} (Y_i - \bar{Y}_R)^2 \sim \chi_{|R|-2}^2.$$

Applying the  $\chi^2$  tail bounds of Laurent and Massart (2000), we get, with probability at least  $1 - 2e^{-x}$ ,

$$\left( \frac{1 - r_R^2}{1 - \rho^2} \right) \sum_{i \in R} (Y_i - \bar{Y}_R)^2 \leq (|R| - 2) + 2\sqrt{x(|R| - 2)} + 2x$$

and

$$\sum_{i \in R} (Y_i - \bar{Y}_R)^2 \geq (|R| - 1) - 2\sqrt{x(|R| - 1)}.$$

Under this event,

$$\frac{1 - r_R^2}{1 - \rho^2} \leq \frac{(|R| - 2) + 2\sqrt{x(|R| - 2)} + 2x}{(|R| - 1) - 2\sqrt{x(|R| - 1)}}.$$

Assuming that  $x = o(|R|)$ , this can be further simplified as

$$\frac{1 - r_R^2}{1 - \rho^2} \leq 1 + o\left(\sqrt{\frac{x}{|R|}}\right).$$

If in addition,  $x \rightarrow \infty$ , then

$$\begin{aligned} -(|R| - 2) \log(1 - r_R^2) &\geq -|R| \log(1 - \rho^2) + o\left(\sqrt{x|R|}\right) \\ &\geq 2 \log(n/|R|) + \delta_n \log(n/|R|) + o\left(\sqrt{x|R|}\right), \end{aligned}$$

which diverges with  $n$  because

$$\delta_n \log(n/|R|) \gg \sqrt{\log(n/|R|)} \gg |R| \gg \sqrt{x|R|}.$$

Since

$$\delta_n \log(n/|R|) \gg \sqrt{\log(n/|R|)} \gg \log \log(n/|R|),$$

this immediately suggests that

$$T^* \geq \frac{1}{\log \log(n/|R|)} \left[ L_R - 2 \log \left( \frac{n}{|R|} \right) \right] \rightarrow_p \infty.$$

Next consider the case when  $|R| \geq \log n$ . Assume without loss of generality that  $\rho > 0$ . The treatment for  $\rho < 0$  is identical. Following an argument similar to that for Lemma 4, we get

$$\sum_{i \in R} (X_i - \bar{X}_R)^2, \sum_{i \in R} (Y_i - \bar{Y}_R)^2 \leq (|R| - 1) + 2\sqrt{x(|R| - 1)} + 2x$$

and

$$\sum_{i \in R} (X_i - \bar{X}_R)(Y_i - \bar{Y}_R) \geq (|R| - 1)\rho - 2\sqrt{x(|R| - 1)} - 2x.$$

with probability at least  $1 - 6e^{-x}$ . Denote this event by  $\mathcal{E}(x)$ . We shall now proceed under  $\mathcal{E}(x)$  with an appropriately chosen  $x \rightarrow \infty$ .

$$r_R \geq \frac{\rho - 2\sqrt{x/(|R| - 1)} - 2[x/(|R| - 1)]}{1 + 2\sqrt{x/(|R| - 1)} + 2[x/(|R| - 1)]}. \quad (17)$$

It is not hard to see that under the condition (9),  $|R|\rho^2 \rightarrow \infty$ . Assuming that  $x \rightarrow \infty$  such that  $x = o(|R|\rho^2)$ , we get

$$r_R \geq \rho + o\left(\sqrt{\frac{x}{|R|}}\right)$$

Then,

$$L_R \geq -(|R| - 2) \log \left[ 1 - \left( \rho + o \left( \sqrt{\frac{x}{|R|}} \right) \right)^2 \right] \quad (18)$$

Recall that

$$-|R| \log(1 - \rho^2) \geq (2 + \delta_n) \log \left( \frac{n}{|R|} \right).$$

Denote by  $\rho_* > 0$  the solution to

$$-|R| \log(1 - \rho_*^2) = (2 + \delta_n) \log \left( \frac{n}{|R|} \right).$$

It is clear that  $\rho \geq \rho_*$ . Together with the fact that the right hand side of (18) is an increasing function of  $\rho$ , we get

$$\begin{aligned} L_R &\geq -(|R| - 2) \log \left[ 1 - \left( \rho_* + o \left( \sqrt{\frac{x}{|R|}} \right) \right)^2 \right] \\ &= -(|R| - 2) \log(1 - \rho_*^2) + o(\sqrt{x|R|}\rho_*) \\ &= (2 + \delta_n) \log \left( \frac{n}{|R|} \right) + o(\sqrt{x|R|}\rho_*^2). \end{aligned}$$

Note that

$$|R|\rho_*^2 \leq 2|R| \log(1 + \rho_*^2) \leq -2|R| \log(1 - \rho_*^2) = 2(2 + \delta_n) \log \left( \frac{n}{|R|} \right).$$

It is not hard to see that

$$\frac{\delta_n^2}{2 + \delta_n} \log \left( \frac{n}{|R|} \right) \rightarrow \infty$$

if (9) holds. Assuming that

$$x = o \left( \frac{\delta_n^2}{2 + \delta_n} \log \left( \frac{n}{|R|} \right) \right),$$

we get

$$T^* \geq (\log \log(n/|R|))^{-1} (L_R - 2 \log(n/|R|)) \rightarrow \infty.$$

This concludes the proof of consistency of  $T^*$  under (9). □



*Proof of Theorem 2 (Consistency of  $\tilde{T}^*$ ).* We now consider the computationally efficient test based on  $\tilde{T}^*$  is also consistent. As before, we begin by arguing that  $\tilde{q}_\alpha = O(1)$ , and then show that under  $H_1$ ,  $\tilde{T}^* \rightarrow \infty$ . To show that  $\tilde{q}_\alpha = O(1)$ , it suffices to note that

$$\begin{aligned} T^* &= \max_{R \in \mathcal{R}} \{(\log \log(n/|R|))^{-1} (L_R - 2 \log(n/|R|))\} \\ &\geq \max_{R \in \cup_k \tilde{\mathcal{R}}_k} \{(\log \log(n/|R|))^{-1} (L_R - 2 \log(n/|R|))\} \\ &= \tilde{T}^*. \end{aligned}$$

Therefore,  $\tilde{q}_\alpha \leq q_\alpha = O(1)$  following the argument before.

Next we show that under the alternative hypothesis where  $X_i$  and  $Y_i$  are correlated on a set  $R \in \mathcal{R}_k$  for some  $k$ ,  $\tilde{T}^* \rightarrow \infty$ . By definition, there exists a  $\tilde{R} \in \tilde{\mathcal{R}}_k$  such that

$$d(R, \tilde{R}) \leq \frac{1}{4k^2}. \quad (19)$$

Observe that

$$\tilde{T}^* \geq \tilde{T}_k^* \geq (\log \log(n/|\tilde{R}|))^{-1} \left( L_{\tilde{R}} - 2 \log(n/|\tilde{R}|) \right).$$

It now suffices to show that the rightmost hand side is unbounded with probability approaching to 1 when  $k \leq k_*$ . To this end, we first consider the case when  $\tilde{R} \subseteq R$ .

Note that if  $\tilde{R} \subseteq R$ , then (2) holds for any  $i \in \tilde{R}$ . Following an identical argument for consistency of  $T^*$ , it suffices to show that there exists a  $\tilde{\delta}_n > 0$  such that

$$\tilde{\delta}_n \sqrt{\log(n/|\tilde{R}|)} \rightarrow \infty \quad (20)$$

and

$$-|\tilde{R}| \log(1 - \rho^2) \geq (2 + \tilde{\delta}_n) \log \left( \frac{n}{|\tilde{R}|} \right). \quad (21)$$

Observe that (19) implies that

$$|\tilde{R}| \geq \left( 1 - \frac{1}{4k^2} \right) |R|.$$

Thus

$$\begin{aligned} |\tilde{R}| \log \frac{1}{1-\rho^2} &\geq \left(1 - \frac{1}{4k^2}\right) |R| \log \frac{1}{1-\rho^2} \\ &\geq \left(1 - \frac{1}{4k^2}\right) (2 + \delta_n) \log \left(\frac{n}{|\tilde{R}|}\right). \end{aligned}$$

Because

$$\log \left(\frac{n}{|R|}\right) = \log \left(\frac{n}{|\tilde{R}|}\right) + \log \left(\frac{|\tilde{R}|}{|R|}\right) \geq \log \left(\frac{n}{|\tilde{R}|}\right) + \log \left(1 - \frac{1}{4k^2}\right) \geq \log \left(\frac{n}{|\tilde{R}|}\right) - \frac{1}{4k^2},$$

we get

$$\begin{aligned} |\tilde{R}| \log \frac{1}{1-\rho^2} &\geq \left(1 - \frac{1}{4k^2}\right) (2 + \delta_n) \left[ \log \left(\frac{n}{|\tilde{R}|}\right) - \frac{1}{4k^2} \right] \\ &\geq \left(1 - \frac{1}{4k^2}\right)^2 (2 + \delta_n) \log \left(\frac{n}{|\tilde{R}|}\right) \\ &\geq \left(1 - \frac{1}{2k^2}\right) (2 + \delta_n) \log \left(\frac{n}{|\tilde{R}|}\right). \end{aligned}$$

Let

$$\tilde{\delta}_n = \left(1 - \frac{1}{2k^2}\right) \delta_n - \frac{1}{k^2}.$$

Then (21) holds. We now verify (20). Recall that

$$\delta_n^2 (k-1) \log 2 \leq \delta_n^2 \log \left(\frac{n}{|R|}\right) \rightarrow \infty,$$

we get, for sufficiently large  $n$ ,

$$\tilde{\delta}_n \geq \frac{1}{4} \delta_n.$$

This implies that

$$\tilde{\delta}_n^2 \log \left(\frac{n}{|\tilde{R}|}\right) \geq \tilde{\delta}_n^2 \log \left(\frac{n}{|R|}\right) \geq \frac{1}{16} \delta_n^2 \log \left(\frac{n}{|R|}\right) \rightarrow \infty,$$

which completes the proof for the case  $\tilde{R} \subseteq R$ .

Now consider the case when  $\tilde{R} \not\subseteq R$ . By definition,

$$\frac{|\tilde{R} \cap R|}{\sqrt{|R||\tilde{R}|}} \geq 1 - \frac{1}{4k^2}.$$

Because  $\tilde{R} \cap R \subseteq \tilde{R}$ , we get

$$\frac{|\tilde{R}|}{|R|} \geq \left(1 - \frac{1}{4k^2}\right)^2. \quad (22)$$

Thus,

$$|\tilde{R} \cap R| \geq \left(1 - \frac{1}{4k^2}\right)^{3/2} |R| \geq \left(1 - \frac{1}{3k^2}\right) |R|.$$

Similarly, we can derive that

$$|\tilde{R} \cap R| \geq \left(1 - \frac{1}{3k^2}\right) |\tilde{R}|. \quad (23)$$

Following the same treatment as for the previous case, we can derive that

$$\frac{1}{\log \log(n/|\tilde{R} \cap R|)} \left[ L_{\tilde{R} \cap R} - 2 \log \left( \frac{n}{|\tilde{R} \cap R|} \right) \right] \rightarrow_p \infty.$$

Since  $|\tilde{R} \cap R| \leq |\tilde{R}|$ ,

$$\frac{1}{\log \log \frac{n}{|\tilde{R}|}} \left[ L_{\tilde{R} \cap R} - 2 \log \left( \frac{n}{|\tilde{R}|} \right) \right] \geq \frac{1}{\log \log \frac{n}{|\tilde{R} \cap R|}} \left[ L_{\tilde{R} \cap R} - 2 \log \left( \frac{n}{|\tilde{R} \cap R|} \right) \right] \rightarrow \infty.$$

It now suffices to show that

$$|L_{\tilde{R} \cap R} - L_{\tilde{R}}| = O_p \left( \log \log \left( \frac{n}{|\tilde{R}|} \right) \right) \quad (24)$$

In the light of (22),

$$\begin{aligned} \log \log \left( \frac{n}{|\tilde{R}|} \right) &\geq \log \left[ \log \left( \frac{n}{|R|} \right) + 2 \log \left( 1 - \frac{1}{4k^2} \right) \right] \\ &\geq \log \left[ (k-1) \log 2 - \frac{1}{2k^2} \right] = O(\log k). \end{aligned}$$

On the other hand, by Lemma 8,

$$|L_{\tilde{R} \cap R} - L_{\tilde{R}}| \leq C \left[ \frac{1}{3k^2} x + |\tilde{R}|^{-1/2} x^{3/2} \right],$$

with probability at least  $1 - e^{-x}$ . Observe that

$$|\tilde{R}| \geq \left(1 - \frac{1}{4k^2}\right)^2 |R| \geq \left(1 - \frac{1}{2k^2}\right) |R| \geq n2^{-(k+1)}.$$

Equation (24) then follows by taking

$$x = \min \left\{ k^2 \log k, 2^{-k/3} n^{1/3} (\log k)^{2/3} \right\}.$$

The proof is now completed. □

*Proof of Theorem 3.* Our argument is similar to those used earlier by Lepski and Tsybakov (2000) and Walther (2010). We shall outline only the main steps for brevity. Note first that a lower bound for a special case necessarily yields a lower bound for the general case. Thus it suffices to consider the case when  $\mu_1 = \mu_2 = 0$  and  $\sigma_1 = \sigma_2 = 1$ . In the light of (5), for any  $A$ , we can find  $\tilde{\mathcal{R}}(A) \subset \mathcal{R}(A)$  such that  $|\tilde{\mathcal{R}}(A)| = c_3(n/A)$ , and for any  $R_1, R_2 \in \tilde{\mathcal{R}}(A)$ ,  $R_1 \cap R_2 = \emptyset$ . For brevity, we shall assume that  $c_3 = 1$  and for any  $R \in \tilde{\mathcal{R}}(A)$ ,  $|R| = A$ . More general case can be treated in the same fashion albeit the argument becomes considerably more cumbersome.

Denote by  $\mathbb{P}_0$  the joint distribution of  $\{(X_i, Y_i) : i \in \mathbb{I}\}$  under null hypothesis, and by  $\mathbb{P}_R$  the joint distribution under alternative hypothesis where  $X_i$  and  $Y_i$  are correlated on  $R \in \tilde{\mathcal{R}}(A)$  so that (2) holds for  $i \in R$  and (1) holds for  $i \notin R$ . The likelihood ratio between  $\mathbb{P}_0$  and  $\mathbb{P}_R$  can be computed:

$$W_R = \frac{d\mathbb{P}_R}{d\mathbb{P}_0} = \frac{1}{(1 - \rho^2)^{A/2}} \exp \left\{ - \frac{\sum_{i \in R} (\rho^2 X_i^2 - 2\rho X_i Y_i + \rho^2 Y_i^2)}{2(1 - \rho^2)} \right\}$$

To prove the first statement, we first show

$$\mathbb{E}_0(W_R^{1+\delta_n/4}) / (\eta |\tilde{\mathcal{R}}(A)|)^{\delta_n/4} \rightarrow 0 \quad \text{for any } 0 < \eta < 1,$$

where  $\mathbb{E}_0$  stands for expectation taken with respect to  $\mathbb{P}_0$ .

It can be computed that

$$\mathbb{E}_0(W_R^{1+\delta_n/4}) = \frac{1}{(1 - \rho^2 \delta_n^2/16)^{A/2} (1 - \rho^2)^{A\delta_n/8}}.$$

Recall that

$$A \log \frac{1}{1 - \rho^2} \leq (2 - \delta_n) \log \frac{n}{A}.$$

We get

$$\begin{aligned} & -\log \left[ \mathbb{E}_0(W_R^{1+\delta_n/4}) / (\eta |\tilde{\mathcal{R}}(A)|)^{\delta_n/4} \right] \\ & \geq \frac{A\delta_n}{8} \log(1 - \rho^2) + \frac{A}{2} \log(1 - \rho^2 \delta_n^2/16) + \frac{\delta_n}{4} \log \frac{n}{A} + \frac{\delta_n}{4} \log \eta \\ & \geq \frac{\delta_n^2}{8} \log \frac{n}{A} - (1 - \delta_n/2) \left( \log \frac{n}{A} \right) \frac{\log(1 - \rho^2 \delta_n^2/16)}{\log(1 - \rho^2)} + \frac{\delta_n}{4} \log \eta \\ & \geq \frac{\delta_n^2}{8} \log \frac{n}{A} - \frac{\delta_n^2}{16} (1 - \delta_n/2) \left( \log \frac{n}{A} \right) + \frac{\delta_n}{4} \log \eta \\ & = \frac{\delta_n^2}{16} (1 + \delta_n/2) \log \frac{n}{A} + \frac{\delta_n}{4} \log \eta \\ & \geq \frac{\delta_n^2}{16} \log \frac{n}{A} \rightarrow \infty. \end{aligned}$$

Thus,

$$\mathbb{E}_0(W_R^{1+\delta_n/4}) / (\eta |\tilde{\mathcal{R}}(A)|)^{\delta_n/4} \rightarrow 0.$$

Next, we argue that

$$\mathbb{E}_0 \left| |\tilde{\mathcal{R}}(A)|^{-1} \sum_{R \in \tilde{\mathcal{R}}(A)} W_R - 1 \right| \rightarrow 0.$$

To this end, write

$$\bar{W} = |\tilde{\mathcal{R}}(A)|^{-1} \sum_{R \in \tilde{\mathcal{R}}(A)} (W_R - 1),$$

$$\bar{W}_1 = |\tilde{\mathcal{R}}(A)|^{-1} \sum_{R \in \tilde{\mathcal{R}}(A)} (W_R - 1) \mathbf{1}_{(|W_R - 1| > \eta |\tilde{\mathcal{R}}(A)|)},$$

and

$$\bar{W}_2 = |\tilde{\mathcal{R}}(A)|^{-1} \sum_{R \in \tilde{\mathcal{R}}(A)} (W_R - 1) \mathbf{1}_{|W_R - 1| \leq \eta |\tilde{\mathcal{R}}(A)|}.$$

Observe that

$$\mathbb{E}_0 |\bar{W}| \leq \mathbb{E}_0 |\bar{W}_1| + \mathbb{E}_0 |\bar{W}_2| \leq \mathbb{E}_0 |\bar{W}_1| + \eta.$$

On the other hand,

$$\mathbb{E}_0 |\bar{W}_1| \leq \mathbb{E}_0 (W_R \mathbf{1}_{(W_R > \eta |\tilde{\mathcal{R}}(A))}) \leq \mathbb{E}_0 (W_R^{1+\delta_n/4}) / (\eta |\tilde{\mathcal{R}}(A)|)^{\delta_n/4} \rightarrow 0.$$

We can take  $\eta \downarrow 0$  to get

$$\mathbb{E}_0 \left| |\tilde{\mathcal{R}}(A)|^{-1} \sum_{R \in \tilde{\mathcal{R}}(A)} W_R - 1 \right| \rightarrow 0.$$

Finally, let  $\mathbb{P}_1$  be the uniform mixture of  $\mathbb{P}_R$  for  $R \in \tilde{\mathcal{R}}(A)$ , that is,

$$\mathbb{P}_1 = |\tilde{\mathcal{R}}(A)|^{-1} \sum_{R \in \tilde{\mathcal{R}}(A)} \mathbb{P}_R.$$

Then for any test  $\Delta$ ,

$$\begin{aligned} \mathbb{P}_0(\Delta = 1) + \mathbb{P}_1(\Delta = 0) &= \mathbb{E}_0(\Delta) + 1 - \min_{R \in \tilde{\mathcal{R}}(A)} \mathbb{E}_R(\Delta) \\ &\geq \mathbb{E}_0(\Delta) + 1 - |\tilde{\mathcal{R}}(A)| \sum_{R \in \tilde{\mathcal{R}}(A)} \mathbb{E}_R(\Delta) \\ &\geq 1 - \mathbb{E}_0(\Delta (1 - |\tilde{\mathcal{R}}(A)|^{-1} \sum_{R \in \tilde{\mathcal{R}}(A)} W_R)) \\ &\geq 1 - \mathbb{E}_0 \left| 1 - |\tilde{\mathcal{R}}(A)|^{-1} \sum_{R \in \tilde{\mathcal{R}}(A)} W_R \right| \rightarrow 1, \end{aligned}$$

which completes the proof of the first statement.

To show the second statement, we assume the contrary that  $c_n$  is bounded from above. Then  $\{c_n\}$  must have a convergent subsequence. Without loss of generality, assume  $c_n$  itself converges to some  $b \in [0, \infty)$ . Then

$$\log W_R \rightarrow_d N \left( -\frac{b}{2}, b \right),$$

which implies that

$$\limsup \mathbb{P}_R(\Delta = 1) < 1.$$

This contradicts with the fact that the type II error of  $\Delta$  goes to 0 as  $n \rightarrow \infty$ . The second statement is therefore proven.  $\square$

## References

- J. Adler, S. Pagakis, and I. Parmryd. Replicate-based noise corrected correlation for accurate measurements of colocalization. *Journal of Microscopy*, 230(1):121–133, 2008.
- E. Arias-Castro, D. Donoho, and X. Huo. Near-optimal detection of geometric objects by fast multiscale methods. *IEEE Transactions on Information Theory*, 51(7):2402–2425, 2005.
- E. Arias-Castro, E. Candès, and A. Durand. Detection of an anomalous cluster in a network. *The Annals of Statistics*, pages 278–304, 2011.
- S. Bolte and F. P. Cordelières. A guided tour into subcellular colocalization analysis in light microscopy. *Journal of Microscopy*, 224(3):213–232, 2006.
- T.T. Cai and M. Yuan. Rate-optimal detection of very short signal segments. *arXiv preprint arXiv:1407.2812*, 2014.
- H. Chan and G. Walther. Detection with the scan and the average likelihood ratio. *Statistica Sinica*, 23:409–428, 2013.
- J. Chen and A. Gupta. Testing and locating variance change points with application to stock prices. *Journal of the American Statistical Association*, 92(438):739–747, 1997.
- J. Comeau, S. Costantino, and P. Wiseman. A guide to accurate fluorescence microscopy colocalization measurements. *Biophysical Journal*, 91(12):4611–4622, 2006.
- S. Costes, D. Daelemans, E. Cho, Z. Dobbin, G. Pavlakis, and S. Lockett. Automatic and quantitative measurement of protein-protein colocalization in live cells. *Biophysical Journal*, 86(6):3993–4003, 2004.

- S. Dasgupta and A. Gupta. An elementary proof of a theorem of Johnson and Lindenstrauss. *Random structures and algorithms*, 22(1):60–65, 2003.
- A. Desolneux, L. Moisan, and J. Morel. Maximal meaningful events and applications to image analysis. *The Annals of Statistics*, pages 1822–1851, 2003.
- L. Dümbgen and V. Spokoiny. Multiscale testing of qualitative hypotheses. *The Annals of Statistics*, pages 124–152, 2001.
- L. Dümbgen and G. Walther. Multiscale inference about a density. *The Annals of Statistics*, pages 1758–1785, 2008.
- J. Fan. Test of significance based on wavelet thresholding and Neyman’s truncation. *Journal of American Statistical Association*, 91:674–688, 1996.
- J. Fan, C. Zhang, and J. Zhang. Generalized likelihood ratio statistics and wilks phenomenon. *The Annals of Statistics*, 29:153–193, 2001.
- J. Fan, X. Han, and W. Gu. Control of the false discovery rate under arbitrary covariance dependence (with discussions). *Journal of American Statistical Association*, 107:1019–1045, 2012.
- J. Glaz, J. Naus, and S. Wallenstein. *Scan statistics*. Springer, New York, 2001.
- P. Hall and J. Jin. Innovated higher criticism for detecting sparse signals in correlated noise. *The Annals of Statistics*, 38(3):1686–1732, 2010.
- H. Herce, C. Casas-Delucchi, and M. Cardoso. New image colocalization coefficient for fluorescence microscopy to quantify (bio-)molecular interactions. *Journal of Microscopy*, 249(3):184–194, 2013.
- H. Hotelling. New light on the correlation coefficient and its transforms. *Journal of the Royal Statistical Society. Series B (Methodological)*, 15(2):193–232, 1953.
- B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection. *The Annals of Statistics*, pages 1302–1338, 2000.



- L. LeCam and G. Yang. *Asymptotics in Statistics: Some Basic Concepts*. Springer, New York, 2000.
- O. Lepski and A. Tsybakov. Asymptotically exact nonparametric hypothesis testing in sup-norm and at a fixed point. *Probability Theory and Related Fields*, 117(1):17–48, 2000.
- E. Manders, J. Stap, G. Brakenhoff, R. Van Driel, and J. Aten. Dynamics of three-dimensional replication patterns during the s-phase, analysed by double labelling of dna and confocal microscopy. *Journal of Cell Science*, 103(3):857–862, 1992.
- E.M. Manders, F. Verbeek, and J. Aten. Measurement of co-localization of objects in dual-colour confocal images. *Journal of Microscopy*, 169(3):375–382, 1993.
- R.J. Muirhead. *Aspects of Multivariate Statistical Theory*. John Wiley & Sons, New York, NY, 2008.
- M. Pacifico, C. Genovese, I. Verdinelli, and L. Wasserman. False discovery control for random fields. *Journal of the American Statistical Association*, 99(468):1002–1014, 2004.
- L. Robinson, V. de la Pena, and Y. Kushnir. Detecting shifts in correlation and variability with applications to enso-monsoon rainfall relationships. *Theoretical and Applied Climatology*, 94:215–224, 2008.
- S. Rodionov. Sequential method of detecting abrupt changes in the correlation coefficient and its application to bering sea climate. *Climate*, 3:474–491, 2015.
- M. Talagrand. *The Generic Chaining*. Springer-Verlag, New York, NY, 2000.
- G. Vidyamurthy. *Pairs Trading: Quantitative Methods and Analysis*. Wiley, New York, NY, 2004.
- G. Walther. Optimal and fast detection of spatial clusters with scan statistics. *The Annals of Statistics*, 38(2):1010–1033, 2010.
- D. Wieda, W. Krämera, and H. Dehling. Testing for a change in correlation at an unknown point in time using an extended functional delta method. *Econometric Theory*, 28:570–589, 2011.

## Appendix – Auxiliary Results and Proofs

We first state tail bounds for  $t$  and  $F$  distributions necessary for our derivations.

**Lemma 1.** *Let  $X$  be a random variable following a  $t$  distribution with degree of freedom  $n > 1$ . There exists a numerical constant  $0 < c_1 < c_2$  such that*

$$c_1 n^{-1/2} \left(1 + \frac{x^2}{n}\right)^{-\frac{n}{2}} \leq \mathbb{P}(|X| > x) \leq c_2 \left(1 + \frac{x^2}{n}\right)^{-\frac{n}{2}} \quad (25)$$

for any  $x \geq 1$ . In particular,

$$c_1 n^{-1/2} e^{-u/2} \leq \mathbb{P} \left\{ n \log \left(1 + \frac{X^2}{n}\right) \geq u \right\} \leq c_2 e^{-u/2}, \quad (26)$$

for any  $u \geq 1$ .

*Proof of Lemma 1.* Recall that the density of a  $t$  distribution with degree of freedom  $n$  is

$$f(x) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi}\Gamma(\frac{n}{2})} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}} \leq C \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}},$$

for an absolute constant  $C > 0$ . Then, for any  $u > 0$ ,

$$\begin{aligned} \mathbb{P}(X > u) &\leq C \int_u^\infty \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}} dx \\ &\leq C \int_u^\infty \frac{x}{u} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}} dx \\ &= \frac{nC}{2u} \int_u^\infty \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}} d\left(1 + \frac{x^2}{n}\right) \\ &= \frac{nC}{(n-1)u} \left(1 + \frac{x^2}{n}\right)^{-\frac{n-1}{2}} \Big|_u^\infty \\ &\leq 2C \frac{1}{u} \left(1 + \frac{u^2}{n}\right)^{-\frac{n-1}{2}}. \end{aligned}$$

The upper bound in (25) follows immediately by taking  $c = 4\sqrt{2}C$ , by symmetry of  $t$

distribution. On the other hand, observe that

$$f(x) \geq C \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}},$$

for some constant  $C > 0$ . Thus,

$$\begin{aligned} \mathbb{P}(X > u) &\geq C \int_u^\infty \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}} dx \\ &\geq C \int_u^\infty \frac{x}{\sqrt{n}} \left(1 + \frac{x^2}{n}\right)^{-\frac{n}{2}-1} dx \\ &= \frac{\sqrt{n}C}{2} \int_u^\infty \left(1 + \frac{x^2}{n}\right)^{-\frac{n}{2}-1} d\left(1 + \frac{x^2}{n}\right) \\ &= \frac{\sqrt{n}C}{(n-1)} \left(1 + \frac{x^2}{n}\right)^{-\frac{n}{2}} \Big|_u^\infty \\ &\geq \frac{\sqrt{n}C}{(n-1)} \left(1 + \frac{u^2}{n}\right)^{-\frac{n}{2}}. \end{aligned}$$

The lower bound in (25) then follows immediately.

Now, taking

$$x = \sqrt{n(e^{u/n} - 1)}$$

in (25) yields (26). □

**Lemma 2.** *Let  $U_1 \sim \chi_{n_1}^2$  and  $U_2 \sim \chi_{n_2}^2$  be two independent random variables. Then for any  $-1 < x < 1$ ,*

$$\mathbb{P} \left\{ \left| \frac{n_1 + n_2}{n_1} \frac{U_1}{U_1 + U_2} - 1 \right| \geq x \right\} \leq 2 \exp \left( -\frac{n_1 x^2}{12} \right).$$

*Proof of Lemma 2.* As shown by Dasgupta and Gupta (2003), for any  $x > 0$ ,

$$\mathbb{P} \left\{ \frac{n_1 + n_2}{n_1} \frac{U_1}{U_1 + U_2} \leq 1 - x \right\} \leq \exp \left( \frac{n_1}{2} (x + \log(1 - x)) \right),$$

and

$$\mathbb{P} \left\{ \frac{n_1 + n_2}{n_1} \frac{U_1}{U_1 + U_2} \geq 1 + x \right\} \leq \exp \left( \frac{n_1}{2} (-x + \log(1 + x)) \right).$$

The claim then follows from the fact that

$$\log(1+x) \leq x - \frac{x^2}{6}$$

for all  $x$  such that  $|x| < 1$ . □

The following observation on the sample correlation coefficient is useful:

**Lemma 3.** *Assume that  $\{(X_i, Y_i) : i \in R\}$  are iid copies of  $(X, Y) \sim N((\mu_1, \mu_2)^\top, \Sigma)$  where*

$$\Sigma = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}.$$

Then

$$\sum_{i \in R} (Y_i - \bar{Y}_R)^2 (1 - r_R^2) \sim (1 - \rho^2) \chi_{|R|-2}^2.$$

*Proof of Lemma 3.* Consider a linear regression of  $Y$  over  $X$ :

$$Y = \beta_0 + \beta X + \epsilon.$$

Recall that

$$\hat{\beta} = \frac{\sum_{i \in R} (X_i - \bar{X}_R)(Y_i - \bar{Y}_R)}{\sum_{i \in R} (X_i - \bar{X}_R)^2}$$

and  $\hat{\beta}_0 = \bar{Y}_R - \hat{\beta} \bar{X}_R$  are the least squares estimate of  $Y$  over  $X$  where

$$\bar{X}_R = \frac{1}{|R|} \sum_{i \in R} X_i \quad \text{and} \quad \bar{Y}_R = \frac{1}{|R|} \sum_{i \in R} Y_i.$$

The residual sum of squares of the regression can then be written as

$$\sum_{i \in R} (Y_i - \hat{\beta}_0 - \hat{\beta} X_i)^2 = \sum_{i \in R} (Y_i - \bar{Y}_R)^2 (1 - r_R^2)$$

Conditioned on  $X_i$ , the residual sum of squares will follow  $(1 - \rho^2) \chi_{|R|-2}^2$ . Thus the margin distribution of the residual sum of squares is also  $(1 - \rho^2) \chi_{|R|-2}^2$ . □

Next we derive a tail bound for the sample correlation coefficient. For brevity, we work with the case when  $(X, Y)$  are known to be centered so that

$$r_R = \frac{\sum_{i \in R} X_i Y_i}{\sqrt{\sum_{i \in R} X_i^2} \sqrt{\sum_{i \in R} Y_i^2}} \quad (27)$$

where  $(X_i, Y_i)$ s are independent copies of  $(X, Y)$ . Treatment for the more general case is completely analogous, yet this simplification allows us to avoid lengthy discussions about the smaller order effects due to centering by sample means, and repeatedly switching between  $|R| - 1$  or  $|R| - 2$  as the appropriate degrees of freedom.

**Lemma 4.** *Assume that  $\{(X_i, Y_i) : i \in R\}$  are iid copies of  $(X, Y) \sim N(0, I_2)$ . Then for any  $x > 0$ ,*

$$\mathbb{P} \left\{ \left| \sum_{i \in R} X_i Y_i \right| \geq 2\sqrt{x|R|} + 2x \right\} \leq 4e^{-x}.$$

*If in addition,  $0 < x < |R|/16$ , then*

$$\mathbb{P}\{|r_R| \geq x\} \leq 6 \exp(-|R|x^2/64).$$

*Proof of Lemma 4.* Write

$$\sum_{i \in R} X_i Y_i = \frac{1}{2} \sum_{i \in R} \left( \frac{1}{\sqrt{2}}(X_i + Y_i) \right)^2 - \frac{1}{2} \sum_{i \in R} \left( \frac{1}{\sqrt{2}}(X_i - Y_i) \right)^2.$$

Then

$$\begin{aligned} \mathbb{P} \left\{ \left| \sum_{i \in R} X_i Y_i \right| \geq 2\sqrt{u|R|} + 2u \right\} &\leq \mathbb{P} \left\{ \left| \sum_{i \in R} \left( \frac{1}{\sqrt{2}}(X_i + Y_i) \right)^2 - |R| \right| \geq 2\sqrt{u|R|} + 2u \right\} \\ &\quad + \mathbb{P} \left\{ \left| \sum_{i \in R} \left( \frac{1}{\sqrt{2}}(X_i - Y_i) \right)^2 - |R| \right| \geq 2\sqrt{u|R|} + 2u \right\} \\ &\leq 4e^{-u}, \end{aligned}$$

where the second inequality follows from the  $\chi^2$  upper and lower tail bound of Laurent and Massart (2000). Applying the  $\chi^2$  lower tail bound from Laurent and Massart (2000), we can

also derive that

$$\mathbb{P} \left\{ \left| \sum_{i \in R} X_i^2 \right| \leq |R| - 2\sqrt{u|R|} \right\} \leq e^{-u} \quad (28)$$

and

$$\mathbb{P} \left\{ \left| \sum_{i \in R} Y_i^2 \right| \leq |R| - 2\sqrt{u|R|} \right\} \leq e^{-u} \quad (29)$$

Therefore, for any  $u < |R|/16$ ,

$$|r_R| \leq \frac{2\sqrt{u|R|} + 2u}{|R| - 2\sqrt{u|R|}} \leq \frac{2}{|R|} \left( 2\sqrt{u|R|} + 2u \right) \leq 8\sqrt{\frac{u}{|R|}}.$$

with probability at least  $1 - 6e^{-u}$ . The claim follows immediately.  $\square$

We are also interested in the difference in correlation coefficients between two different regions. The following lemma provides a useful probabilistic tool for such purposes.

**Lemma 5.** *Assume that  $\{(X_i, Y_i) : i \in R_1 \cup R_2\}$  are iid copies of  $(X, Y) \sim N(0, I_2)$ , and  $2|R_1 \cap R_2| \geq |R_1 \cup R_2|$ . Then there exist numerical constants  $c_0, c_1, c_2 > 0$  such that for any  $x < c_0|R_1|$ ,*

$$\mathbb{P}(|R_1|r_{R_1}^2 - |R_2|r_{R_2}^2| \geq x) \leq c_1 \exp \left( -c_2 \min \left\{ \left( \frac{|R_1 \cap R_2|}{|R_1 \cup R_2| - |R_1 \cap R_2|} \right)^{1/2} x, |R_1 \cap R_2|^{1/3} x^{2/3} \right\} \right).$$

In particular, if

$$\zeta := \frac{|R_1 \cap R_2|}{\sqrt{|R_1||R_2|}} \geq \frac{1}{4},$$

then there exists a numerical constant  $c_3 > 0$  such that for any  $x < c_0|R_1|$ ,

$$\mathbb{P}(|R_1|r_{R_1}^2 - |R_2|r_{R_2}^2| \geq x) \leq c_1 \exp \left( -c_3 \min \left\{ (1 - \zeta)^{-1/2} x, |R_1 \cap R_2|^{1/3} x^{2/3} \right\} \right).$$

*Proof of Lemma 5.* We first consider the case when  $R_2 \subseteq R_1$ . Recall that

$$r_{R_1} = \frac{\sum_{i \in R_1} X_i Y_i}{\sqrt{\sum_{i \in R_1} X_i^2 \sum_{i \in R_1} Y_i^2}}, \quad \text{and} \quad r_{R_2} = \frac{\sum_{i \in R_1} X_i Y_i}{\sqrt{\sum_{i \in R_2} X_i^2 \sum_{i \in R_2} Y_i^2}}.$$

Therefore,

$$\begin{aligned}
|R_2|r_{R_2}^2 - |R_1|r_{R_1}^2 &= \left( \frac{1}{\sqrt{|R_2|}} \sum_{i \in R_2} X_i Y_i \right)^2 \left( \frac{|R_2|^2}{\sum_{i \in R_2} X_i^2 \sum_{i \in R_2} Y_i^2} - \frac{|R_1|^2}{\sum_{i \in R_1} X_i^2 \sum_{i \in R_1} Y_i^2} \right) \\
&\quad + \frac{|R_1|^2}{\sum_{i \in R_1} X_i^2 \sum_{i \in R_1} Y_i^2} \left[ \frac{1}{|R_2|} \left( \sum_{i \in R_2} X_i Y_i \right)^2 - \frac{1}{|R_1|} \left( \sum_{i \in R_1} X_i Y_i \right)^2 \right]
\end{aligned}$$

We now bound the terms on the right hand side separately.

Observe that

$$\begin{aligned}
&\left| \left( \frac{1}{\sqrt{|R_2|}} \sum_{i \in R_2} X_i Y_i \right)^2 \left( \frac{|R_1|^2}{\sum_{i \in R_1} X_i^2 \sum_{i \in R_1} Y_i^2} - \frac{|R_2|^2}{\sum_{i \in R_2} X_i^2 \sum_{i \in R_2} Y_i^2} \right) \right| \\
&= |R_2|r_{R_2}^2 \left| 1 - \left( \frac{|R_1|}{|R_2|} \frac{\sum_{i \in R_2} X_i^2}{\sum_{i \in R_1} X_i^2} \right)^{-1} \left( \frac{|R_1|}{|R_2|} \frac{\sum_{i \in R_2} Y_i^2}{\sum_{i \in R_1} Y_i^2} \right)^{-1} \right|.
\end{aligned}$$

By Lemma 2,

$$\mathbb{P} \left\{ \left| \frac{|R_1|}{|R_2|} \frac{\sum_{i \in R_2} X_i^2}{\sum_{i \in R_1} X_i^2} - 1 \right| \geq x \right\} \leq 2 \exp \left( -\frac{|R_2|}{12} x^2 \right),$$

and

$$\mathbb{P} \left\{ \left| \frac{|R_1|}{|R_2|} \frac{\sum_{i \in R_2} Y_i^2}{\sum_{i \in R_1} Y_i^2} - 1 \right| \geq x \right\} \leq 2 \exp \left( -\frac{|R_2|}{12} x^2 \right).$$

We get, for any  $x < 1/2$ ,

$$\left| 1 - \left( \frac{|R_1|}{|R_2|} \frac{\sum_{i \in R_2} X_i^2}{\sum_{i \in R_1} X_i^2} \right)^{-1} \left( \frac{|R_1|}{|R_2|} \frac{\sum_{i \in R_2} Y_i^2}{\sum_{i \in R_1} Y_i^2} \right)^{-1} \right| \leq 4x$$

with probability at least  $1 - 4 \exp(-|R_2|x^2/12)$ . On the other hand, by Lemma 4,

$$\mathbb{P}\{|r_{R_2}| \geq x\} \leq 6 \exp(-|R_2|x^2/64).$$

Thus, by taking  $u = 4|R_2|x^3$ ,

$$|R_2|r_{R_2}^2 \left| 1 - \left( \frac{|R_1|}{|R_2|} \frac{\sum_{i \in R_2} X_i^2}{\sum_{i \in R_1} X_i^2} \right)^{-1} \left( \frac{|R_1|}{|R_2|} \frac{\sum_{i \in R_2} Y_i^2}{\sum_{i \in R_1} Y_i^2} \right)^{-1} \right| \leq u$$

with probability at least

$$1 - 10 \exp\left(-\frac{|R_2|^{1/3} u^{2/3}}{64 \cdot 4^{2/3}}\right).$$

Denote by  $\mathcal{E}_1(u)$  the event that the above inequality holds.

To bound the second term, first note that

$$\begin{aligned} & \frac{1}{|R_2|} \left( \sum_{i \in R_2} X_i Y_i \right)^2 - \frac{1}{|R_1|} \left( \sum_{i \in R_1} X_i Y_i \right)^2 \\ &= \left( \frac{1}{|R_2|} - \frac{1}{|R_1|} \right) \left( \sum_{i \in R_1} X_i Y_i \right)^2 - \frac{1}{|R_2|} \left( \sum_{i \in R_1 \setminus R_2} X_i Y_i \right)^2 - \frac{2}{|R_2|} \left( \sum_{i \in R_2} X_i Y_i \right) \left( \sum_{i \in R_1 \setminus R_2} X_i Y_i \right). \end{aligned}$$

By Lemma 4,

$$\frac{|R_1|^2}{\sum_{i \in R_1} X_i^2 \sum_{i \in R_1} Y_i^2} \left( \frac{1}{|R_2|} - \frac{1}{|R_1|} \right) \left( \sum_{i \in R_1} X_i Y_i \right)^2 = \left( \frac{|R_1|}{|R_2|} - 1 \right) |R_1| r_{R_1}^2 \leq \left( \frac{|R_1|}{|R_2|} - 1 \right) |R_1| x^2,$$

with probability at least  $1 - 6 \exp(-|R_1| x^2 / 64)$ . On the other hand, again by Lemma 4,

$$\mathbb{P} \left\{ \left| \sum_{i \in R_1 \setminus R_2} X_i Y_i \right| \geq 2\sqrt{x(|R_1| - |R_2|)} + 2x \right\} \leq 4e^{-x},$$

Recall that, by  $\chi^2$  lower tail bounds from Laurent and Massart (2000), we get

$$\mathbb{P} \left\{ \sum_{i \in R_1} X_i^2 \leq |R_1| - 2\sqrt{x|R_1|} \right\} \leq e^{-x},$$

and

$$\mathbb{P} \left\{ \sum_{i \in R_1} Y_i^2 \leq |R_1| - 2\sqrt{x|R_1|} \right\} \leq e^{-x}.$$

Thus, for any  $x < |R_1|/16$ ,

$$\begin{aligned} \frac{|R_1|^2/|R_2|}{\sum_{i \in R_1} X_i^2 \sum_{i \in R_1} Y_i^2} \left( \sum_{i \in R_1 \setminus R_2} X_i Y_i \right)^2 &\leq \frac{|R_1|^2/|R_2|}{\left( |R_1| - 2\sqrt{x|R_1|} \right)^2} \left( 2\sqrt{x(|R_1| - |R_2|)} + 2x \right)^2 \\ &\leq \frac{16}{|R_2|} \left( \sqrt{x(|R_1| - |R_2|)} + x \right)^2, \end{aligned}$$



with probability at least  $1 - 6e^{-x}$ . In other words,

$$\begin{aligned} \frac{|R_1|^2/|R_2|}{\sum_{i \in R_1} X_i^2 \sum_{i \in R_1} Y_i^2} \left( \sum_{i \in R_1 \setminus R_2} X_i Y_i \right)^2 &\leq \frac{1}{|R_2|} \left( \frac{1}{2} x \sqrt{|R_1|(|R_1| - |R_2|)} + \frac{|R_1|x^2}{16} \right)^2 \\ &\leq \frac{1}{2} \left( \frac{|R_1|}{|R_2|} - 1 \right) |R_1|x^2 + \frac{|R_1|^2 x^4}{128|R_2|}, \end{aligned}$$

with probability at least  $1 - 6 \exp(-|R_1|x^2/64)$  for any  $x < 2$ . Following a similar argument, we can also show that

$$\begin{aligned} &\frac{2|R_1|^2/|R_2|}{\sum_{i \in R_1} X_i^2 \sum_{i \in R_1} Y_i^2} \left( \sum_{i \in R_2} X_i Y_i \right) \left( \sum_{i \in R_1 \setminus R_2} X_i Y_i \right) \\ &\leq \frac{1}{|R_2|} \left( \frac{1}{2} x \sqrt{|R_1||R_2|} + \frac{|R_1|x^2}{16} \right) \left( \frac{1}{2} x \sqrt{|R_1|(|R_1| - |R_2|)} + \frac{|R_1|x^2}{16} \right), \end{aligned}$$

with probability at least  $1 - 10 \exp(-|R_1|x^2/64)$  for any  $x < 2$ . Note  $2|R_2| \geq |R_1|$ . In summary, we get

$$\begin{aligned} &\frac{|R_1|^2}{\sum_{i \in R_1} X_i^2 \sum_{i \in R_1} Y_i^2} \left| \frac{1}{|R_2|} \left( \sum_{i \in R_2} X_i Y_i \right)^2 - \frac{1}{|R_1|} \left( \sum_{i \in R_1} X_i Y_i \right)^2 \right| \\ &\leq 2 \left( \frac{|R_1|}{|R_2|} - 1 \right)^{1/2} |R_1|x^2 + \frac{|R_1|x^3}{16} \end{aligned}$$

with probability at least  $1 - 22 \exp(-|R_1|x^2/64)$  for any  $x < 2$ . Hence, with probability at least

$$1 - 22 \exp \left( -\frac{1}{256} \left( \frac{|R_1|}{|R_2|} - 1 \right)^{-1/2} u \right) - 22 \exp \left( -\frac{1}{16} |R_1|^{1/3} u^{2/3} \right)$$

we have

$$\frac{|R_1|^2}{\sum_{i \in R_1} X_i^2 \sum_{i \in R_1} Y_i^2} \left| \frac{1}{|R_2|} \left( \sum_{i \in R_2} X_i Y_i \right)^2 - \frac{1}{|R_1|} \left( \sum_{i \in R_1} X_i Y_i \right)^2 \right| \leq u.$$

Denote this event by  $\mathcal{E}_2(u)$ .

In summary, for any  $u < |R_1|/256$ ,

$$||R_1|r_{R_1}^2 - |R_2|r_{R_2}^2| \leq 2u$$

with probability at least

$$\begin{aligned} \mathbb{P}\left\{\mathcal{E}_1(u) \cap \mathcal{E}_2(u)\right\} &\geq 1 - 22 \exp\left(-\frac{1}{256} \left(\frac{|R_1|}{|R_2|} - 1\right)^{-1/2} u\right) - 22 \exp\left(-\frac{1}{16} |R_1|^{1/3} u^{2/3}\right) \\ &\quad - 10 \exp\left(-\frac{|R_2|^{1/3} u^{2/3}}{64 \cdot 4^{2/3}}\right) \\ &\geq 1 - 22 \exp\left(-\frac{1}{256} \left(\frac{|R_1|}{|R_2|} - 1\right)^{-1/2} u\right) - 32 \exp\left(-\frac{1}{128} |R_2|^{1/3} u^{2/3}\right). \end{aligned}$$

The statement, when  $R_2 \subseteq R_1$ , then follows.

Now consider the general case when  $R_2 \not\subseteq R_1$ . In this case,

$$\left||R_1|r_{R_1}^2 - |R_2|r_{R_2}^2\right| \leq \left||R_1|r_{R_1}^2 - |R_1 \cap R_2|r_{R_1 \cap R_2}^2\right| + \left||R_2|r_{R_2}^2 - |R_1 \cap R_2|r_{R_1 \cap R_2}^2\right|.$$

We can now appeal to the bounds we derived for nested sets before to get

$$\begin{aligned} \mathbb{P}(\left||R_1|r_{R_1}^2 - |R_1 \cap R_2|r_{R_1 \cap R_2}^2\right| \geq x) &\leq 22 \exp\left(-\frac{1}{256} \left(\frac{|R_1 \cap R_2|}{|R_1| - |R_1 \cap R_2|}\right)^{1/2} u\right) \\ &\quad + 32 \exp\left(-\frac{1}{128} |R_1 \cap R_2|^{1/3} u^{2/3}\right), \end{aligned}$$

and

$$\begin{aligned} \mathbb{P}(\left||R_2|r_{R_2}^2 - |R_1 \cap R_2|r_{R_1 \cap R_2}^2\right| \geq x) &\leq 22 \exp\left(-\frac{1}{256} \left(\frac{|R_1 \cap R_2|}{|R_2| - |R_1 \cap R_2|}\right)^{1/2} u\right) \\ &\quad + 32 \exp\left(-\frac{1}{128} |R_1 \cap R_2|^{1/3} u^{2/3}\right). \end{aligned}$$

The first claim then follows from an application of the union bound.

To show the second statement, assume that  $|R_1| \geq |R_2|$  without loss of generality. Observe that

$$\rho = \frac{|R_1 \cap R_2|}{\sqrt{|R_1||R_2|}} \leq \sqrt{\frac{|R_2|}{|R_1|}},$$

which implies that  $|R_2| \geq \rho^2 |R_1|$ . Therefore,

$$|R_1 \cap R_2| = \rho \sqrt{|R_1||R_2|} \geq \rho^2 |R_1|,$$

and

$$|R_1 \cap R_2| = \rho \sqrt{|R_1||R_2|} \geq \rho |R_2|.$$

Thus,

$$\frac{|R_1 \cap R_2|}{|R_1 \cup R_2| - |R_1 \cap R_2|} \geq \frac{1}{\rho^{-2} + \rho^{-1} - 2} = \frac{1}{1 - \rho} \frac{\rho^2}{2\rho + 1} \geq \frac{1}{48} (1 - \rho)^{-1},$$

where the last inequality follows from the fact that  $1/4 \leq \rho \leq 1$ .  $\square$

We are now in position to derive bounds for the likelihood ratio statistic  $L_R$ . Since we work with centered random variables as stated earlier, it is natural to redefine  $L_R$  as:

$$L_R = -(|R| - 1) \log(1 - r_R^2).$$

where  $r_R$  is given by (27).

**Lemma 6.** *Assume that  $\{(X_i, Y_i) : i \in R\}$  are iid copies of  $(X, Y) \sim N(0, I_2)$  for some  $|R| > 1$ . Then there exists numerical constants  $0 < c_1 < c_2$  such that for any  $x > 1$ ,*

$$c_1 |R|^{-1/2} e^{-x/2} \leq \mathbb{P}(L_R > x) \leq c_2 e^{-x/2}.$$

*Proof of Lemma 6.* Observe that

$$L_R = (|R| - 1) \log \left( 1 + \frac{T_R^2}{|R| - 1} \right)$$

where

$$T_R = r_R \sqrt{\frac{|R| - 1}{1 - r_R^2}}$$

and

$$r_R = \frac{\sum_{i \in R} X_i Y_i}{\sqrt{\sum_{i \in R} X_i^2 \sum_{i \in R} Y_i^2}}$$

It is well known that, under the null hypothesis,

$$T_R \sim t_{|R|-1}.$$

See, e.g., Hotelling (1953). By Lemma 1,

$$c_1 |R|^{-1/2} e^{-x/2} \leq \mathbb{P}(L_R > x) \leq c_2 e^{-x/2},$$

for any  $x > 1$ . □

The following lemma bounds the change in the likelihood ration statistic due to a perturbation of the index set.

**Lemma 7.** *Assume that  $\{(X_i, Y_i) : i \in R_1 \cup R_2\}$  are iid copies of  $(X, Y) \sim N(0, I_2)$ , and  $2|R_1 \cap R_2| \geq |R_1 \cup R_2|$ . Then there exist numerical constants  $c_0, c_1, c_2 > 0$  such that for any  $x < c_0 |R_1|$ ,*

$$\mathbb{P}(|L_{R_1} - L_{R_2}| \geq x) \leq c_1 \exp \left( -c_2 \min \left\{ \left( \frac{|R_1 \cap R_2|}{|R_1 \cup R_2| - |R_1 \cap R_2|} \right)^{1/2} x, |R_1 \cap R_2|^{1/3} x^{2/3} \right\} \right).$$

In particular, if

$$\zeta := \frac{|R_1 \cap R_2|}{\sqrt{|R_1||R_2|}} \geq \frac{1}{4},$$

then there exists a numerical constant  $c_3 > 0$  such that for any  $x < c_0 |R_1|$ ,

$$\mathbb{P}(|L_{R_1} - L_{R_2}| \geq x) \leq c_1 \exp \left( -c_3 \min \left\{ (1 - \zeta)^{-1/2} x, |R_1 \cap R_2|^{1/3} x^{2/3} \right\} \right).$$

*Proof of Lemma 7.* Similar to Lemma 5, it suffices to prove the first statement when  $R_2 \subseteq R_1$ . By the convexity of  $-\log(1 - x)$ , we can ensure

$$L_{R_1} = -|R_1| \log(1 - r_{R_1}^2) \geq -|R_1| \log(1 - r_{R_2}^2) + \frac{|R_1|(r_{R_1}^2 - r_{R_2}^2)}{1 - r_{R_2}^2}$$

and

$$L_{R_2} = -|R_2| \log(1 - r_{R_2}^2) \geq -|R_2| \log(1 - r_{R_1}^2) + \frac{|R_2|(r_{R_2}^2 - r_{R_1}^2)}{1 - r_{R_1}^2}$$

Therefore,

$$|L_{R_1} - L_{R_2}| \leq (|R_2| - |R_1|) \log(1 - \max\{r_{R_1}^2, r_{R_2}^2\}) + \frac{|R_2||r_{R_2}^2 - r_{R_1}^2|}{1 - \max\{r_{R_1}^2, r_{R_2}^2\}}. \quad (30)$$

We now bound the two terms on the right hand side separately.

Denote by  $\mathcal{E}(\alpha)$  the event that

$$\max\{r_{R_1}^2, r_{R_2}^2\} < \alpha.$$

By Lemma 4,

$$\mathbb{P}\{\mathcal{E}(\alpha)\} \geq 1 - 12 \exp(-|R_2|\alpha/64).$$

Note that, for any  $0 < x < \alpha$ ,

$$-\log(1 - x) \leq \frac{x}{1 - \alpha}.$$

We can upper bound the first term on the right hand side of (30) by

$$\frac{1}{1 - \alpha} (|R_1| - |R_2|) \max\{r_{R_1}^2, r_{R_2}^2\}.$$

Therefore,

$$\mathbb{P}\{(|R_2| - |R_1|) \log(1 - \max\{r_{R_1}^2, r_{R_2}^2\}) \geq u\} \leq 12 \exp\left(-\frac{1}{64}|R_2| \min\left\{\alpha, \frac{(1 - \alpha)u}{|R_1| - |R_2|}\right\}\right). \quad (31)$$

The second term of (30) can be upper bounded by

$$\frac{1}{1 - \alpha} (||R_2|r_{R_2}^2 - |R_1|r_{R_1}^2| + (|R_1| - |R_2|)r_{R_1}^2),$$

under the event  $\mathcal{E}(\alpha)$ . By Lemma 5, we get

$$\mathbb{P}\{||R_2|r_{R_2}^2 - |R_1|r_{R_1}^2| \geq x\} \leq c_1 \exp\left(-c_2 \min\left\{\left(\frac{|R_2|}{|R_1| - |R_2|}\right)^{1/2} x, |R_2|^{1/3} x^{2/3}\right\}\right).$$

And by Lemma 4,

$$\mathbb{P}\{(|R_1| - |R_2|)r_{R_1}^2 \geq x\} \leq 6 \exp\left(-\frac{|R_1|x}{64(|R_1| - |R_2|)}\right)$$

Therefore,

$$\mathbb{P} \left\{ \frac{|R_2||r_{R_2}^2 - r_{R_1}^2|}{1 - \max\{r_{R_1}^2, r_{R_2}^2\}} \geq u \right\} \leq 1 - 12 \exp(-|R_2|\alpha/64) - 6 \exp \left( -\frac{(1-\alpha)|R_1|u}{128(|R_1| - |R_2|)} \right) - c_1 \exp \left( -\frac{1-\alpha}{2} c_2 \min \left\{ \left( \frac{|R_2|}{|R_1| - |R_2|} \right)^{1/2} u, |R_2|^{1/3} u^{2/3} \right\} \right).$$

Together with (31), this implies the desired statement for  $R_2 \subseteq R_1$ .  $\square$

A careful inspection of the derivation of Lemma 7 suggests that it can be extended to a more general situation where  $X$  and  $Y$  are correlated for some indices.

**Lemma 8.** *Let  $R_1 \subset R_2$  be two index sets. Assume that  $\{(X_i, Y_i) : i \in R_1\}$  are independent observations so that  $(X_i, Y_i) \sim N(0, I_2)$  for  $i \in R_1$ , and  $X_i, Y_i$  are standard normal random variables with correlation coefficient  $\rho$  for  $i \notin R_1$ . Then there exist numerical constants  $c_0, c_1, c_2 > 0$  such that for any  $x < c_0|R_1|$ ,*

$$\mathbb{P}(|L_{R_1} - L_{R_2}| \geq x) \leq c_1 \exp \left( -c_2 \min \left\{ (1 - \zeta)^{-1/2} x, |R_1|^{1/3} x^{2/3} \right\} \right).$$

provided that  $\zeta := |R_1|/|R_2| \geq 1/4$ .

Finally we derive a perturbation bounds for a polygon which is useful for our discussion in Section 3.1.

**Lemma 9.** *Let  $K_1$  and  $K_2$  be two polygons with vertices  $u_1, u_2, \dots, u_k$  and  $v_1, v_2, \dots, v_k$  respectively. Denote by  $e_j$  the length of the edge between  $u_j$  and  $u_{(j \bmod k)+1}$ . Then*

$$|K_1 \cap K_2^c| \leq r \sum_{j=1}^k (e_j + 2r),$$

where  $r$  is maximum distance between  $u_j$  and  $v_j$ .

*Proof.* Denote by  $Q_i$  the polygon whose first  $i$  vertices are the same with  $K_2$  and whose remaining vertices are the same with  $K_1$ . In particular,  $Q_0 = K_1$  and  $Q_k = K_2$ . It is not hard to see that the  $j$ th edge of  $Q_i$  is no longer than  $e_j + 2r$ . If we compare  $Q_0$  and  $Q_1$ , then only the first vertex might be different, as illustrated in Figure 8.

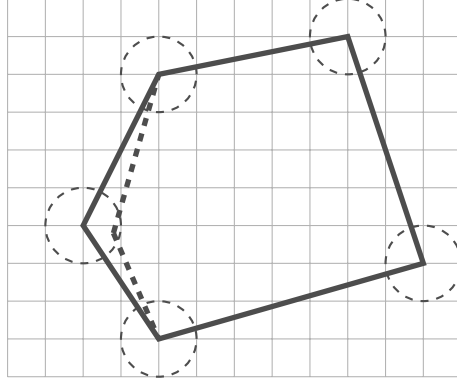


Figure 8: Effect of perturbation of vertices of a polygon.

Because  $Q_0$  and  $Q_1$  are different only in the first vertex, they can only be different in the two edges linked with the first index. It can then be computed that

$$|Q_0 \cap Q_1^c| \leq \frac{1}{2}r(e_1 + e_k + 4r)$$

Similarly,

$$|Q_i \cap Q_{i+1}^c| \leq \frac{1}{2}r(e_i + e_{i+1} + 4r), \quad i = 1, 2, \dots, k-1,$$

It is clear that

$$K_1 \cap K_2^c = Q_0 \cap Q_k^c \subset \bigcup_{i=0}^{k-1} (Q_i \cap Q_{i+1}^c)$$

Therefore,

$$|K_1 \cap K_2^c| \leq r \sum_{i=1}^k (e_i + 2r),$$

which completes the proof. □

*Proof of Proposition 1.* Write

$$\mathcal{C}_{p_1, \dots, p_k} = \{K(\{(a_i, b_i) : 1 \leq i \leq k\}) : 2^{p_i} \leq r_i < 2^{p_i+1}, i = 1, \dots, k\}.$$

It is clear that there exists a constant  $C > 0$  such that

$$|\mathcal{C}_{p_1, \dots, p_k}| \leq Cn2^{2(\sum_{i=1}^k p_i)}.$$

Note that there are constants  $c_1, c_2 > 0$  depending on  $k$  and  $M$  only such that

$$\mathcal{R}_{\text{polygon}}(A; k, M) \subset \{K \in \mathcal{R}_{\text{polygon}}(k, M) : c_1 A^{1/2} \leq r_i \leq c_2 A^{1/2}, i = 1, 2, \dots, k\}.$$

Therefore,

$$\mathcal{R}_{\text{polygon}}(A; k, M) \leq cnA^k$$

which completes the proof because  $A \asymp r_i^2$ . □

*Proof of Proposition 2.* Note that  $\pi_s(K(\{(a_i, b_i) : 1 \leq i \leq k\}))$  is also a polygon. For brevity, we shall hereafter denote it by  $K(\{(\tilde{a}_i, \tilde{b}_i) : 1 \leq i \leq k\})$ . By Lemma 9, we get

$$|K(\{(a_i, b_i) : 1 \leq i \leq k\}) \setminus K(\{(\tilde{a}_i, \tilde{b}_i) : 1 \leq i \leq k\})| \leq C2^s \sum_i r_i \leq Ck2^s r_1.$$

Hence

$$\rho\left(K(\{(a_i, b_i) : 1 \leq i \leq k\}), K(\{(\tilde{a}_i, \tilde{b}_i) : 1 \leq i \leq k\})\right) \geq 1 - \frac{Ck2^s r_1}{\pi r_1^2} \geq 1 - \frac{C2^s}{r_1},$$

which completes the proof. □