

Gene coexpression measures in large heterogeneous samples using count statistics

Y. X. Rachel Wang^a, Michael S. Waterman^{b,1}, and Haiyan Huang^{a,1}

^aDepartment of Statistics, University of California, Berkeley, CA 94720; and ^bProgram in Molecular and Computational Biology, University of Southern California, Los Angeles, CA 90089

Contributed by Michael S. Waterman, September 9, 2014 (sent for review June 27, 2014; reviewed by Wing Hung Wong and Hongyu Zhao)

With the advent of high-throughput technologies making large-scale gene expression data readily available, developing appropriate computational tools to process these data and distill insights into systems biology has been an important part of the “big data” challenge. Gene coexpression is one of the earliest techniques developed that is still widely in use for functional annotation, pathway analysis, and, most importantly, the reconstruction of gene regulatory networks, based on gene expression data. However, most coexpression measures do not specifically account for local features in expression profiles. For example, it is very likely that the patterns of gene association may change or only exist in a subset of the samples, especially when the samples are pooled from a range of experiments. We propose two new gene coexpression statistics based on counting local patterns of gene expression ranks to take into account the potentially diverse nature of gene interactions. In particular, one of our statistics is designed for time-course data with local dependence structures, such as time series coupled over a subregion of the time domain. We provide asymptotic analysis of their distributions and power, and evaluate their performance against a wide range of existing coexpression measures on simulated and real data. Our new statistics are fast to compute, robust against outliers, and show comparable and often better general performance.

local rank patterns | bivariate association | random permutation statistics | Stein’s approximation

A major challenge in systems biology is to understand the intricate interactions and functional relationships between genes and their regulation targets. As advances in high-throughput technologies lead to the generation of enormous amounts of genomic data, the last decade has witnessed a rapidly increasing effort to develop computational tools to reconstruct gene relationships based on a wide range of “omic” data available, in particular transcriptomic or expression data. Coexpression methods, which assess certain types of dependence between the expression profiles of two genes, are one of the earliest tools used for this purpose. The technique has been routinely used for functional gene annotation (1, 2) and more importantly as a measure of edge weights for reconstructing gene networks (3–7).

The problem of finding gene coexpression is closely related to that of detecting bivariate association between two vectors. Since the work by Eisen et al. (8), the Pearson correlation has been adopted as the most widely used coexpression measure (3, 9, 10) for its straightforward conceptual interpretation and computational efficiency. However, it is also known that the Pearson correlation is unsuitable for capturing nonlinear relationships and susceptible to high false discovery rates. Another class of coexpression methods is based on mutual information (MI) (5, 11, 12, 13), which measures general statistical dependence rather than a specific type of bivariate association. The computation of MI involves discretization of the data and tuning parameters, and obtaining P values requires computationally intensive permutation tests. The practical benefits and shortcomings of MI compared with correlation-based methods are still under investigation (11, 12, 14). More comparisons of different coexpression measures and the coexpression networks constructed can be found in refs. 15 and 16.

In the broader statistical literature, other methods available for quantifying bivariate associations include the Renyi correlation (17) measuring the correlation between two variables after suitable transformations; various regression-based techniques (14); and Hoeffding’s D (18), and distance covariance (dCov) (19), for general statistical dependence. These methods are not widely adopted in genomic applications yet. More recently, Reshef et al. (20) proposed the maximal information coefficient (MIC) as an extension of MI, but MIC was shown to have inferior power to dCov (21) and MI (22) in various simulated scenarios.

Most of the methods mentioned so far, perhaps with the exception of MIC, do not specifically target dependence relationships that can be local in nature and often assume the data are random samples from a common distribution in the theoretical analysis. However, real gene interactions may change as the intrinsic cellular state varies or only exist under a specific cellular condition. Furthermore, with data integration now being a routine approach to combat the curse of dimensionality, samples from different experimental conditions or tissue types are likely to prescribe different gene relationships and thus create more complex situations for detecting gene interactions. For instance, a protein that positively regulates expression in one context may act as a repressor in another [e.g., MECP2 (23)], or a gene may participate in either neural development or hematopoiesis depending on tissue type [e.g., EBF1 (24, 25)]. One possible approach to discern local gene interactions is biclustering (26, 27), which simultaneously clusters genes and samples. However, most biclustering techniques are restricted to detecting simple subclasses of linear associations. On the algorithm side, the optimizations of most criteria for measuring the quality of given biclusters can only be achieved locally, and their global behaviors are hard to

Significance

Coexpression analysis is one of the earliest tools for inferring gene associations using expression data but faces new challenges in this “big data” era. In a large heterogeneous dataset, it is likely that gene relationships may change or only exist in a subset of the samples, and they can be nonlinear or nonfunctional. We propose two new robust count statistics to account for local patterns in gene expression profiles. The statistics are generalizable to detect statistical dependence in other application domains. The performance of the statistics is evaluated against a number of popular bivariate dependence measures, showing favorable results. The asymptotic studies of the statistics provide an interesting addition to the combinatorics literature.

Author contributions: M.S.W. and H.H. designed research; Y.X.R.W., M.S.W., and H.H. performed research; Y.X.R.W. analyzed data; and Y.X.R.W., M.S.W., and H.H. wrote the paper.

Reviewers: W.H.W., Stanford University; H.Z., Yale University.

The authors declare no conflict of interest.

Data deposition: The code reported in this paper is available at www.stat.berkeley.edu/~hhuang/coexpression.

See Commentary on page 16236.

¹To whom correspondence may be addressed. Email: hhuang@stat.berkeley.edu or msw@usc.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1417128111/-DCSupplemental.

characterize. Most algorithms also involve a number of tuning parameters with little guidance on how to choose them.

Motivated by these observations, we propose two new coexpression measures based on matching patterns of local expression ranks using count statistics. Our robust statistics specifically take into account the local nature of gene associations while being general enough to detect other common types of dependence relationships. In particular, one of our statistics is designed for time-course data with local dependence structures, such as time series that are coupled over a subregion of the time domain. This is a unique feature compared with other popular coexpression measures. The statistics are fast to compute, and we provide theoretical analysis of their asymptotic properties. We demonstrate their applicability via comparisons to a comprehensive list of existing methods on simulated and real data. Our new methods show better precision, and have the important ability to detect subtle gene relationships that are easily missed by other methods.

Definitions and Asymptotic Properties

For a heterogeneous set of samples with potentially changing gene interactions, we can define a general coexpression measure by aggregating the interactions across all subsamples of size $k \leq n$. For genes x and y with expression levels from n samples $\mathbf{x} = (x_1, \dots, x_n)$ and $\mathbf{y} = (y_1, \dots, y_n)$, we consider

$$W = \sum_{1 \leq i_1 < \dots < i_k \leq n} F(x_{i_1}, \dots, x_{i_k}; y_{i_1}, \dots, y_{i_k}), \quad [1]$$

where $F(\cdot; \cdot)$ is an interaction measure on local expression profiles $(x_{i_1}, \dots, x_{i_k})$ and $(y_{i_1}, \dots, y_{i_k})$ from a subset of k samples. In this paper, we choose $F(\cdot; \cdot)$ to be an indicator function comparing the rank patterns of the subsequences $(x_{i_1}, \dots, x_{i_k})$ and $(y_{i_1}, \dots, y_{i_k})$. Depending on the nature of the expression data studied, we define two corresponding count statistics.

- i) When dealing with time-course data, it is sensible to preserve the order of the samples and consider only interactions within contiguous subsequences. We define W_1 as

$$W_1 = \sum_{i=1}^{n-k+1} \left\{ \mathbb{I}(\phi(x_i, \dots, x_{i+k-1}) = \phi(y_i, \dots, y_{i+k-1})) + \mathbb{I}(\phi(x_i, \dots, x_{i+k-1}) = \phi(-y_i, \dots, -y_{i+k-1})) \right\}, \quad [2]$$

where $\mathbb{I}(\cdot)$ is an indicator function and ϕ is the rank function. That is, ϕ returns the indices of elements in a vector after they have been sorted in an increasing order. W_1 counts the number of contiguous subsequences of length k with matching and reverse rank patterns, indicating positive and negative associations respectively.

- ii) When the order of the samples is not particularly meaningful (e.g., non-time-series data), we consider a more general count that includes all subsequences of length k ,

$$W_2 = \sum_{1 \leq i_1 < \dots < i_k \leq n} \left\{ \mathbb{I}(\phi(x_{i_1}, \dots, x_{i_k}) = \phi(y_{i_1}, \dots, y_{i_k})) + \mathbb{I}(\phi(x_{i_1}, \dots, x_{i_k}) = \phi(-y_{i_1}, \dots, -y_{i_k})) \right\}. \quad [3]$$

It is easy to see that W_2 is equal to the number of increasing (and decreasing) subsequences of length k in a suitably permuted sequence. Suppose σ is a permutation that sorts the elements of \mathbf{y} in an increasing order. Let $\mathbf{z} = \sigma(\mathbf{x})$ be that permutation applied to \mathbf{x} ; W_2 can be rewritten as

$$W_2 = \sum_{1 \leq i_1 < \dots < i_k \leq n} \left\{ \mathbb{I}(z_{i_1} < \dots < z_{i_k}) + \mathbb{I}(z_{i_1} > \dots > z_{i_k}) \right\}. \quad [4]$$

A simple example of the two counts above is given in *SI Appendix*. Both counts are symmetric with respect to \mathbf{x} and \mathbf{y} and efficient

to compute. Counting W_1 has a running time of $O(k(\log k)n)$, while counting W_2 takes $O(kn \log n)$ time using dynamic programming and binary indexed trees. More details on the computation time are given in *SI Appendix, Proofs*.

Asymptotic Distributions. We can derive the asymptotic distributions of W_1 and W_2 for different regimes of k assuming the following: (i) The two sequences \mathbf{x} and \mathbf{y} are independent and have no ties within themselves and (ii) at least one of \mathbf{x} and \mathbf{y} has an exchangeable distribution. Note that the second assumption implies the ranks of the expression vector with an exchangeable distribution is a random permutation of $\{1, 2, \dots, n\}$.

The Stein and Chen–Stein approximations (28, 29) give us the following two asymptotic regimes for W_1 , the proof of which is given in *SI Appendix, Proofs*.

Theorem 1. For $n \rightarrow \infty$, $k \geq 3$ and $k/(\log n)^\alpha \rightarrow 0$ for some $\alpha < 1$,

$$T_1 := \frac{W_1 - \mu_{1,n}}{\sigma_{1,n}} \xrightarrow{D} N(0, 1), \quad [5]$$

where $\mu_{1,n} = 2(n-k+1)/k!$, $\sigma_{1,n}^2 = \text{Var}(W_1)$. For $n \rightarrow \infty$, $\frac{\log n}{k} = O(1)$,

$$d_{TV}(W_1, Z) \rightarrow 0, \quad [6]$$

where $Z \sim \text{Poisson}(\mu_{1,n})$ and d_{TV} is the total variation distance.

When \mathbf{x} and \mathbf{y} satisfy the first assumption and assuming without loss of generality \mathbf{x} satisfies the second assumption, the ranks of \mathbf{z} follow the distribution of a random permutation. While the properties and asymptotic distribution of the longest increasing subsequence (LIS) in a random permutation have been much studied and the statistic itself has been used in a number of applications (30–34), not so much attention has been paid to increasing subsequences of length k . Here we use the results in ref. 35 and the Stein approximation to derive a central limit theorem for W_2 for k growing sufficiently slowly. The proof of the theorem is given in *SI Appendix, Proofs*, and the key lies in obtaining a good upper and lower bound on the variance of W_2 .

Theorem 2. For $n \rightarrow \infty$, $k \geq 3$ and $k/(\log n)^\alpha \rightarrow 0$ for some $\alpha < 1$,

$$T_2 := \frac{W_2 - \mu_{2,n}}{\sigma_{2,n}} \xrightarrow{D} N(0, 1), \quad [7]$$

where $\mu_{2,n} = 2 \binom{n}{k} / k!$ and $\sigma_{2,n}^2 = \text{Var}(W_2)$.

Asymptotic Power. Next we analyze the power of T_1 and T_2 under specific alternative distributions. The first scenario we consider is related to time-course data, where the temporal order of \mathbf{x} and \mathbf{y} are preserved in subsequence analysis.

Theorem 3. Let $\mathbf{x} = (x_1, \dots, x_n)$ and $\mathbf{y} = (y_1, \dots, y_n)$ be two time series with n observations, m of which are perfectly coupled in the sense that $\phi(x_i, \dots, x_{i+m-1}) = \phi(y_i, \dots, y_{i+m-1})$. As $n \rightarrow \infty$, $m \rightarrow \infty$,

- i) T_1 goes to infinity in the following regimes:

- For fixed k , if $m \sim a_1 n$, $a_1 > 2/k!$, then $T_1 = \Omega(\sqrt{n})$.
- For $k \rightarrow \infty$ and $k/(\log n)^\alpha \rightarrow 0$, $\alpha < 1$,
 - if $m \geq a_2 \cdot \frac{n}{k!}$, $a_2 > 2$, then $T_1 = \Omega(\sqrt{n/k!})$;
 - if $m \sim a_3 n$, $a_3 \in (0, 1]$, then $T_1 = \Omega(\sqrt{nk!})$.

- ii) T_2 goes to infinity in the following regimes:

- For fixed k , if $m \sim b_1 n$, $b_1^k > 2/k!$, then $T_2 = \Omega(\sqrt{n})$.
- For $k \rightarrow \infty$ and $k/(\log n)^\alpha \rightarrow 0$, $\alpha < 1$,
 - if $m \geq \frac{en}{k}$, then $T_2 = \Omega(\sqrt{n/k^{3/2}})$;
 - if $m \sim b_2 n$, $b_2 \in (0, 1]$, then $T_2 = \Omega(b_2^k k! \sqrt{n/k^{5/2}})$.

Here $\Omega(\cdot)$ denotes asymptotic lower bound.

Remark 1. In the regimes above, using T_1 and T_2 as statistics both lead to rejection of the null hypothesis with probability 1. We also observe that for both T_1 and T_2 , large k leads to better power in the sense that (i) the statistics have a better

convergence rate when m grows as a fraction of n and (ii) a smaller lower bound on m can be achieved, consequently tolerating more noise in the data, while maintaining the power of the tests going to 1. Comparing T_1 and T_2 , T_1 has better power in the regime of fixed k because T_1 allows for a smaller lower bound on m while maintaining the power going to 1.

The next scenario we consider is when \mathbf{x} and \mathbf{y} follow a perfect functional relationship with d strictly monotonic pieces. This is a reasonable subclass of general functional relationships to study since most smooth functions can be approximated by piecewise strictly monotonic functions. In this case, the order of the data does not have to be preserved, making T_1 a less suitable statistic than T_2 . We only analyze the power of T_2 .

Theorem 4. $\mathbf{y}=f(\mathbf{x})$ for $\mathbf{x}\stackrel{\text{iid}}{\sim}\text{Unif}(0,1)$, f can be decomposed into a fixed number of d strictly monotonic pieces which have lengths ℓ_1, \dots, ℓ_d when projected on to the x axis. As $n \rightarrow \infty$,

- For fixed k , if $d^{k-1} < k!/2$, then $\mathbb{P}(T_2 \geq C\sqrt{n}) \rightarrow 1$;
- For $k \rightarrow \infty$ and $k/(\log n)^\alpha \rightarrow 0$, $\alpha < 1$, then $\mathbb{P}(T_2 \geq C\sqrt{n/k^{5/2}}k!/d^{k-1}) \rightarrow 1$

for some constant $C > 0$.

Remark 2. In the regimes above, the power of the statistic T_2 approaches 1. Larger k and smaller d lead to better convergence rates and thus better power. Having fewer monotonic pieces implies there are more uninterrupted counts in each piece contributing to W_2 .

The proofs of the above theorems are in [SI Appendix](#).

Simulations

To investigate the power of our statistics in more realistic settings, we considered four types of bivariate relationships, all of which are illustrative of gene coexpression relationships likely to exist in an expression dataset. It is essential to include a linear type of relationship since pairwise gene relationships detected by current analyses are still predominantly linear. As an example of nonmonotonic associations, we considered a quadratic relationship. The cross-shaped relationship may occur when two genes switch from activators to repressors across different tissue types or treatment conditions, or simply due to the changes in intrinsic cellular state (36). These relationships have also been used as illustrative scenarios in refs. 20 and 22 in the context of general statistical dependence. An important additional example we considered here pertains to the case of genes with time-course data. We simulated two time series that were coupled over sub-regions of the time domain. The robustness of the statistics was tested against outliers—a ubiquitous feature of biological data. Descriptions of the parameters used for each type of relationship are provided in [SI Appendix, Table S2](#).

Throughout the rest of the paper, the variances of W_1 and W_2 were estimated by Monte Carlo experiments. We compared the power of T_1 and T_2 with seven other popular measures of dependence (the Pearson, Spearman, and Renyi correlations, Hoeffding's D, dCov, MI, and MIC). An additional comparison with LIS-based statistics (34) is provided in [SI Appendix, Fig. S4](#). We chose $k=5$ for T_1 and T_2 guided by the log value of the sample size 220. The results from other values of k are provided in [SI Appendix, Fig. S2](#). We note that the influence of k on the power of T_2 is negligible. While the choice of k has a bigger effect on the power of T_1 due to a smaller number of possible values for the counts, the conclusions we draw from qualitative comparisons with the other measures do not change. More details on the computation of the statistics and their P values can be found in [SI Appendix, Simulations](#).

The power values of various statistics computed under four types of dependence relationships are shown in Fig. 1. Unsurprisingly, the Pearson and Spearman correlations can only detect the linear relationship, with the Pearson correlation being more sensitive to outliers. Across the first three types of dependence, T_2 , Hoeffding's D, MI, dCov, and Renyi are the only statistics maintaining reasonable power throughout. Of these

statistics, Renyi and MI have the best performance on the quadratic relationship, but are underpowered on the linear relationship. For the linear scenario, we also computed a variant of T_1 and T_2 counting only the matching rank patterns (omitting the reverse patterns), which are denoted T_1^+ and T_2^+ in the plot. These unidirectional counts provide a way to significantly improve the power when the monotonicity of the relationship is known. In fact, T_2^+ demonstrates the best power while remaining robust to outliers. On the cross relationship, T_2 has a higher power than all of the other statistics. T_1 does not perform well on the first three types of relationships as it is designed for data with a temporal order.

T_1 and T_2 are the only statistics showing significant power on the time-course data. Without respecting the order of the data points, the scatter plot shows no obvious association pattern, making it difficult for the other measures to detect the dependence structure. T_1 has a slightly better power than T_2 .

We remark here that although other dependence relationships were tested in refs. 20 and 22, most of these are less often observed in real gene coexpression patterns. Such examples include sinusoidal, circular, and checkerboard relationships. For the former two examples, we expect the power of T_2 to be affected by the noise level and the frequency of the sinusoidal wave. As discussed in Theorem 4, the power of T_2 is boosted by having uninterrupted counts from monotonic pieces of the association pattern. Since the checkerboard pattern is not piecewise monotonic, we do not expect T_2 to detect this type of relationships.

In addition, we performed simulations to show the behaviors of the statistics conform to their derived asymptotics. Detailed simulation procedures and results are described in [SI Appendix, Asymptotic Convergence](#).

Real Data Examples

In this section, we evaluate the performance of our new statistics on two gene expression datasets: the classic yeast gene expression dataset (3), and a collection of microarray data for *Arabidopsis* tissues downloaded from the National Center for Biotechnology Information Gene Expression Omnibus.

Yeast Cell Cycle Data. The yeast expression data contain the expression levels of 6,178 genes from four reasonably long time-course experiments with a total of 73 time points. More details on data processing are in [SI Appendix, Yeast Cell Cycle](#). We focused on the coexpression of 133 transcription factors (TFs) ([Dataset S1](#)). Using all of the statistics discussed in simulations, we computed 133×133 coexpression matrices and compared them to a total of 428 curated genetic and physical interactions from BioGrid.

As we expected T_1 to be more suitable for time-course data than T_2 , we examined the interactions identified by T_1 more closely. These interactions reveal the ability of T_1 to capture important bivariate associations missed by the other methods. Fig. 2 shows two pairs of TFs (BAS1 vs. GCN4; MSN2 vs. YAP1) whose coexpression strengths were consistently ranked among the top 10 and top 20 by T_1 with $k=7$ but were assigned very low rankings by all of the other methods. Both pairs correspond to previously reported genetic interactions curated in BioGrid. However, their scatter plots show no obvious trends or dependence patterns, highlighting the importance of preserving the temporal order of the data. More specifically, Gcn4p and Bas1p were shown to be involved in cooperative transcriptional regulation of the ADE3 gene, which encodes an essential regulon enzyme for the biosynthesis of several amino acids (37). MSN2 and YAP1 are both activators required for oxidative stress tolerance, and there is a partial overlap between their H₂O₂-inducible regulons (38). Studies using epistatic miniarray profiles (39, 40) have shown that double mutations in MSN2 and YAP1 lead to severe fitness defect. Two more such examples can be found in [SI Appendix, Fig. S5](#).

[SI Appendix, Table S3](#), shows the number of known interactions between TFs among strongly coexpressed pairs as ranked by various statistics. Overall, T_1 (with various choices of k) and the Pearson correlation have the largest number of overlaps

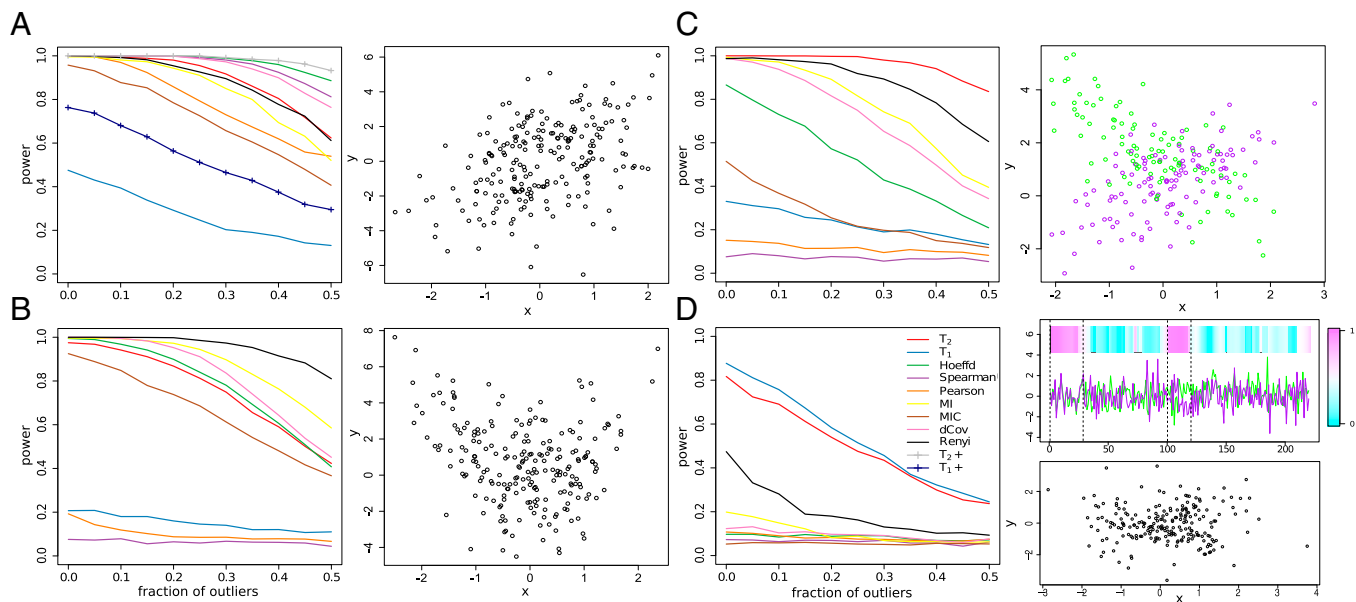


Fig. 1. The power of various statistics rejecting at 5% significance level as level of contamination by outliers increases when the bivariate data follow (A) a linear relationship, (B) a quadratic relationship, (C) a cross-shaped relationship, and (D) two partially coupled time series. The heat map in D shows the absolute values of the Pearson correlations calculated at each time point including its neighboring 15 points.

with the known interactions, with T_1 being the better of the two at most cutoffs. These are followed by T_2 and the Renyi correlation.

Arabidopsis Microarrays. We integrated data from 13 microarray experiments to create a metadata with 220 samples for 22,810 *Arabidopsis* genes. The samples were harvested from shoot

tissues and different regions of root tissues subject to various stress experiments including salt, low pH, and sulfur deficiency treatments. From ref. 41, we downloaded a list of genes involved in the glucosinolates biosynthesis pathway in addition to the 30 pathways in ref. 15 to comprise a total of 510 unique pathway genes (Dataset S2). We computed the pairwise coexpressions between these pathway genes and all of the genes available to

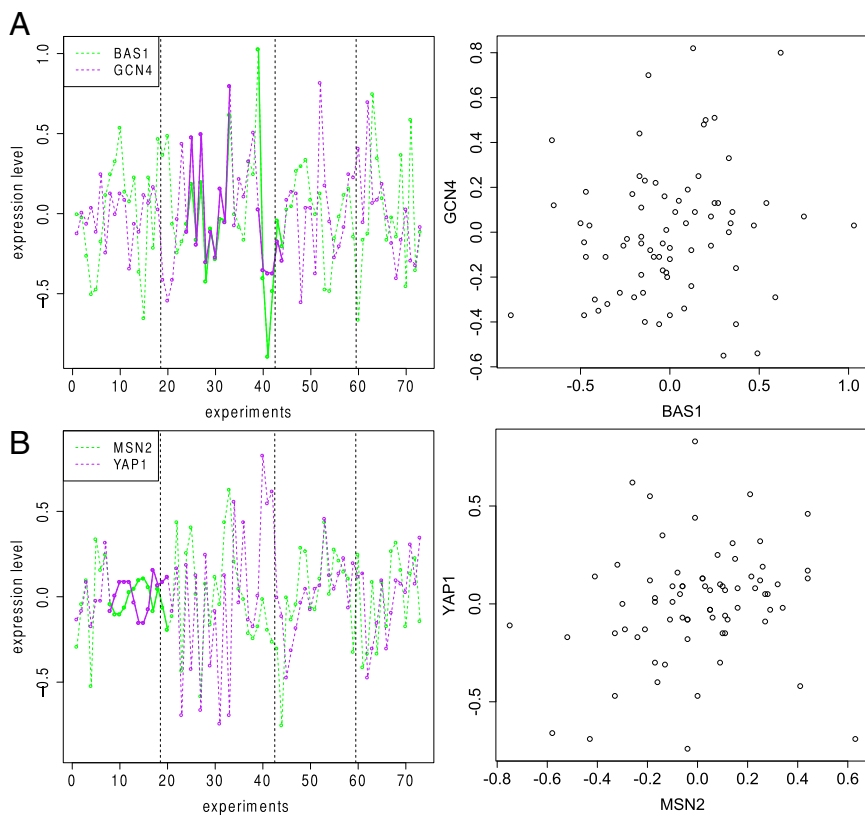


Fig. 2. Expression levels of (A) BAS1 and GCN4 and (B) MSN2 and YAP1 in four time-course experiments (boundaries indicated by the dashed lines). The darker solid lines highlight regions contributing to the counts in T_1 . Both pairs of genes have reported genetic interactions. Their coexpression strengths were consistently ranked among the top 10 and top 20 by T_1 with $k = 7$ but were assigned very low rankings by all of the other methods.

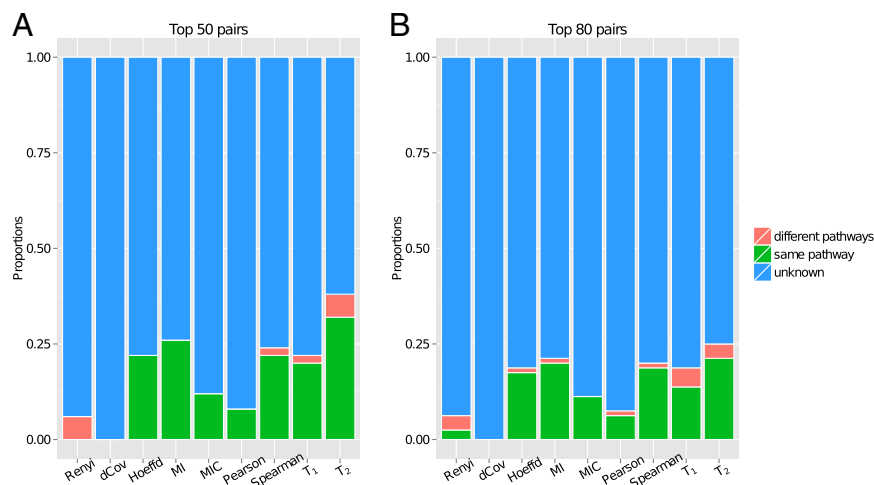


Fig. 3. Number of gene pairs in the same pathway (green), in different pathways (red), and containing a nonpathway gene (blue) among (A) the top 50 pairs and (B) the top 80 pairs as ranked by each method.

test the performance of various measures on distinguishing genes in the same pathway. Our selection of k was guided by the log value of the total sample size, which is ~ 5 . The results presented here were obtained by setting $k=5$ for T_1 and $k=9$ for T_2 . As expected, T_2 is not sensitive to the choice of k , and the results below remain stable for a range of k . More information on data processing can be found in *SI Appendix, Arabidopsis Microarrays*.

Fig. 3 shows the proportions of gene pairs (i) in the same pathway, (ii) in two different pathways, and (iii) containing one nonpathway gene among the top 50 and 80 pairs as ranked by all of the methods. T_2 achieves the best pathway enrichment, followed by MI, the Spearman correlation, Hoeffding's D, and T_1 . As the samples are not composed of long time-course data, it is not surprising that T_1 is a less ideal statistic than T_2 . dCov and Renyi are among the worst performing methods, with almost no pairs in the same pathway, despite their good performance in simulations. Extending the cutoffs to examine more highly ranked pairs, in *SI Appendix, Fig. S6*, the same trend continues for the best four methods until around the top 700 pairs, after which they start to become indistinguishable. dCov remains at the bottom of the list.

Fig. 4 shows two examples where the gene pairs are in the same pathway, but their coexpression values remain significant only under T_2 at 5% level after Bonferroni correction. Some of the sample points are color coded according to their tissue types or treatments to highlight the different patterns of association they exhibit and the lack of a consistent global structure. T_2 is more powerful in this situation due to its definition.

A closer look at the types of relationships detected by T_2 and its closest competitor MI reveals that MI is underpowered on

linear relationships with outliers, an issue also reported by ref. 14. An example is shown in *SI Appendix, Fig. S7*, for two pairs of genes in the same pathway, where the bulk of the samples follow a linear trend but they failed to be identified by MI at an unadjusted significance level of 5%. On the other hand, both pairs were assigned significant P values by T_2 and other statistics, including the Pearson and Spearman correlations.

We also examined the performance of each method in individual pathways. *SI Appendix, Table S4*, shows the methods with the highest counts of same pathway pairs in 20 pathways achieving statistically significant enrichment of pathway genes. T_2 outperforms all of the other methods in 12 pathways out of 20, followed by MI and T_1 , which are the best methods in 4 out of 20 pathways. Note that in the four pathways where MI achieves the highest counts, it is always tied with T_2 , whereas T_2 and T_1 are the unique maxima in six and four pathways, respectively. This implies T_1 and T_2 are potentially more accurate than the other methods in capturing certain coexpression relationships.

Discussion

Statistically, the problem of discovering gene coexpression is to detect bivariate associations between gene expression profiles. In this paper, we propose two new statistics capable of detecting local dependence structures within expression data, motivated by the observation that real gene relationships may have disparaging behaviors in large heterogeneous samples. The statistics are fast to compute, and their asymptotic distributions under the null assumption of independence and exchangeable sample distributions can be derived.

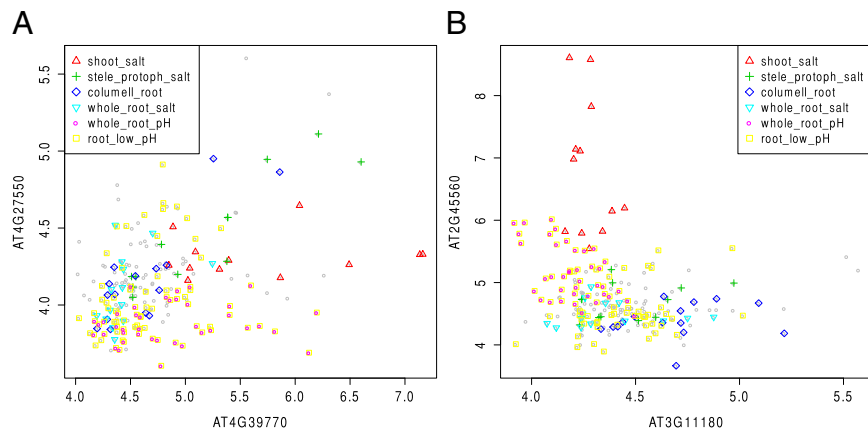


Fig. 4. Expression levels of two gene pairs in the same pathway (A) AT4G27550 and AT4G39770, and (B) AT2G45560 and AT3G11180 with some samples color coded according to their tissue types or treatments: shoot_salt, shoot tissues under salt stress; stele_protoph_salt, stele and protophloem cells under salt stress; columell_root, columella root cap under salt stress, low pH, and sulfur deficiency; whole_root_salt, whole root under salt stress; whole_root_pH, whole root under low pH; root_low_pH, other root cells under low pH.

As demonstrated in both simulation and the yeast cell cycle data, T_1 specializes in detecting local associations in time-course data. In particular, when such associations are not visible within the global association pattern, T_1 offers an attractive alternative to other commonly used coexpression measures. The statistic T_2 , which considers more general local patterns of dependence, is effective on a variety of functional and nonfunctional relationships. However, as T_2 relies on counts from monotonic subpatterns, it is sensitive to noise on high-frequency sinusoidal relationships.

Both statistics involve a tuning parameter k . Asymptotic considerations suggest that values around $\log n$ are reasonable choices since this is within the normal regime of convergence and larger k values are preferable based on the power studies. In simulations, fluctuations of k around $\log n$ have very little effect on the results of T_2 (SI Appendix, Fig. S2). For the *Arabidopsis* data, a range of k can be chosen (5–10) with a small impact on the final results. Due to the more discrete nature of its distribution, T_1 is more sensitive to the choice of k . However, for the yeast cell cycle data, the interacting gene pairs in Fig. 2 received consistent high rankings with $k = 6$ –9. More comparisons of different k values are provided in SI Appendix, Table S3. In practice, choosing k also involves a tradeoff between precision and recall—a common theme of most tuning parameter problems. Larger k would favor higher precision but make the statistics less robust to noise and outliers. More thorough studies investigating how it affects the performance of the statistics in relation to the structure of data would be desirable.

Our definitions and asymptotic analyses of the two unnormalized counts W_1 and W_2 naturally suggest further investigation. Modifying the current definitions to account for ties in the data would be a desirable addition. Extending W_1 to capture temporal dependence patterns with lags would be important for discovering delayed regulations (42). At a more fundamental level, other choices of the interaction measure $F(\cdot; \cdot)$ in Eq. 1 would be interesting to explore. For instance, we can consider relaxing the exact pattern matches to approximate matches, or replacing the indicator function itself with a correlation-based statistic. In terms of asymptotics, it would be of theoretical interest to study the limiting distribution of W_2 for k beyond the log regime. In practice, there often exist inherent dependence structures among the gene samples, especially in time-course data. Thus, removing the exchangeability assumption in the analysis of the null distributions would improve computational accuracy of the P values. Alternatively, it would also be interesting to study the sample dependence directly by reversing the roles of genes and samples and applying a similar technique.

ACKNOWLEDGMENTS. We thank Xianghong Jasmine Zhou for inspiring our study; Fang Fang for her help in initializing this project; and Peter Bickel, Terry Speed, and Larry Goldstein for helpful discussions and feedback. This work was supported by National Institutes of Health Grant U01 HG007031 and National Science Foundation Grant DMS-1160319.

- Zhou X, Kao MC, Wong WH (2002) Transitive functional annotation by shortest-path analysis of gene expression data. *Proc Natl Acad Sci USA* 99(20):12783–12788.
- Fu FF, Xue HW (2010) Coexpression analysis identifies Rice Starch Regulator1, a rice AP2/EREBP family transcription factor, as a novel rice starch biosynthesis regulator. *Plant Physiol* 154(2):927–938.
- Spellman PT, et al. (1998) Comprehensive identification of cell cycle-regulated genes of the yeast *Saccharomyces cerevisiae* by microarray hybridization. *Mol Biol Cell* 9(12):3273–3297.
- Zhang B, Horvath S (2005) A general framework for weighted gene co-expression network analysis. *Stat Appl Genet Mol Biol* 4:e17.
- Basso K, et al. (2005) Reverse engineering of regulatory networks in human B cells. *Nat Genet* 37(4):382–390.
- Yang Y, et al. (2014) Gene co-expression network analysis reveals common system-level properties of prognostic genes across cancer types. *Nat Commun* 5:3231.
- Forrest AR, et al.; FANTOM Consortium and the RIKEN PMI and CLST (DGT) (2014) A promoter-level mammalian expression atlas. *Nature* 507(7493):462–470.
- Eisen MB, Spellman PT, Brown PO, Botstein D (1998) Cluster analysis and display of genome-wide expression patterns. *Proc Natl Acad Sci USA* 95(25):14863–14868.
- Wolfe CJ, Kohane IS, Butte AJ (2005) Systematic survey reveals general applicability of “guilt-by-association” within gene coexpression networks. *BMC Bioinformatics* 6:227.
- Stuart JM, Segal E, Koller D, Kim SK (2003) A gene-coexpression network for global discovery of conserved genetic modules. *Science* 302(5643):249–255.
- Steuer R, Kurths J, Daub CO, Weise J, Selbig J (2002) The mutual information: Detecting and evaluating dependencies between variables. *Bioinformatics* 18(Suppl 2):S231–S240.
- Daub CO, Steuer R, Selbig J, Kloska S (2004) Estimating mutual information using B-spline functions—An improved similarity measure for analysing gene expression data. *BMC Bioinformatics* 5:118.
- Margolin AA, et al. (2006) ARACNE: An algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context. *BMC Bioinformatics* 7(Suppl 1):S7.
- Song L, Langfelder P, Horvath S (2012) Comparison of co-expression measures: Mutual information, correlation, and model based indices. *BMC Bioinformatics* 13(1):328.
- Kumari S, et al. (2012) Evaluation of gene association methods for coexpression network construction and biological knowledge discovery. *PLoS ONE* 7(11):e50411.
- Allen JD, Xie Y, Chen M, Girard L, Xiao G (2012) Comparing statistical methods for constructing large scale gene networks. *PLoS ONE* 7(1):e29348.
- Rényi A (1959) On measures of dependence. *Acta Math Hung* 10(3):441–451.
- Hoeffding W (1948) A non-parametric test of independence. *Ann Math Stat* 19(4):546–557.
- Kosorok MR (2009) Brownian distance covariance. *Ann Appl Stat* 3(4):1266–1269.
- Reshef DN, et al. (2011) Detecting novel associations in large data sets. *Science* 334(6062):1518–1524.
- Simon N, Tibshirani R (2014) Comment on “detecting novel associations in large data sets” by Reshef et al, *Science* Dec 16, 2011. arXiv:1401.7645.
- Kinney JB, Atwal GS (2014) Equitability, mutual information, and the maximal information coefficient. *Proc Natl Acad Sci USA* 111(9):3354–3359.
- Chahrour M, et al. (2008) MeCP2, a key contributor to neurological disease, activates and represses transcription. *Science* 320(5880):1224–1229.
- Milatovich A, Qiu R-G, Grosschedl R, Francke U (1994) Gene for a tissue-specific transcriptional activator (EBF or Olf-1), expressed in early B lymphocytes, adipocytes, and olfactory neurons, is located on human chromosome 5, band q34, and proximal mouse chromosome 11. *Mamm Genome* 5(4):211–215.
- Zhao F, McCarrick-Walmsley R, Åkerblad P, Sigvardsson M, Kadesch T (2003) Inhibition of p300/CBP by early B-cell factor. *Mol Cell Biol* 23(11):3837–3846.
- Cheng Y, Church GM (2000) Biclustering of expression data. *Proc Int Conf Intell Syst Mol Biol* 8:93–103.
- Madeira SC, Oliveira AL (2004) Biclustering algorithms for biological data analysis: A survey. *IEEE/ACM Trans Comput Biol Bioinformatics* 1(1):24–45.
- Stein C (1986) *Approximate Computation of Expectations*. Lecture Notes-Monograph Series, ed Gupta SS (Inst Math Sci, Hayward, CA), Vol 7.
- Chen LHY (1975) Poisson approximation for dependent trials. *Ann Probab* 3(3):534–545.
- Logan BF, Shepp LA (1977) A variational problem for random young tableaux. *Adv Math* 26(2):206–222.
- Baik J, Deift P, Johansson K (1999) On the distribution of the length of the longest increasing subsequence of random permutations. *J Am Math Soc* 12(4):1119–1178.
- Aldous D, Diaconis P (1999) Longest increasing subsequences: From patience sorting to the Baik-Deift-Johansson theorem. *Bull Am Math Soc* 36(4):413–432.
- Arratia R, Barbour AD, Tavaré S (2003) *Logarithmic Combinatorial Structures: A Probabilistic Approach* (Eur Math Soc, Zurich), Vol 1.
- García JE, González-López VA (2014) Independence tests for continuous random variables based on the longest increasing subsequence. *J Multivariate Anal* 127:126–146.
- Pinsky R (2006) Law of large numbers for increasing subsequences of random permutations. *Random Struct Algorithms* 29(3):277–295.
- Li KC (2002) Genome-wide coexpression dynamics: Theory and application. *Proc Natl Acad Sci USA* 99(26):16875–16880.
- Joo YJ, et al. (2009) Cooperative regulation of ADE3 transcription by Gcn4p and Bas1p in *Saccharomyces cerevisiae*. *Eukaryot Cell* 8(8):1268–1277.
- Hasan R, et al. (2002) The control of the yeast H₂O₂ response by the Msn2/4 transcription factors. *Mol Microbiol* 45(1):233–241.
- Zheng J, et al. (2010) Epistatic relationships reveal the functional organization of yeast transcription factors. *Mol Syst Biol* 6(1):420.
- Bandyopadhyay S, et al. (2010) Rewiring of genetic networks in response to DNA damage. *Science* 330(6009):1385–1389.
- Kim K, Jiang K, Teng SL, Feldman LJ, Huang H (2012) Using biologically interrelated experiments to identify pathway genes in *Arabidopsis*. *Bioinformatics* 28(6):815–822.
- Ma P, Castillo-Davis CI, Zhong W, Liu JS (2006) A data-driven clustering method for time course gene expression data. *Nucleic Acids Res* 34(4):1261–1269.

Supporting Information

The supporting information is divided into the following sections:

1. An example of computing W_1 and W_2
2. Additional simulations for asymptotic convergence
3. Simulation details and additional results
4. Additional information and results for the yeast cell cycle data
5. Additional information and results for the *Arabidopsis* data
6. Proofs of the theorems

1 An Example of W_1 and W_2

Suppose $\mathbf{x} = (1, 3, 4, 2, 5)$, $\mathbf{y} = (1, 4, 5, 2, 3)$, and we are interested in computing W_1 and W_2 counts for $k = 3$. For W_1 , there are three possible positions to start a contiguous subsequence of length 3, and only the ones starting at position 1 and 2 have the same ranks in \mathbf{x} and \mathbf{y} . There are no pairs of contiguous subsequences of length 3 with reverse ranks. Hence $W_1 = 2$. To compute W_2 , first sort \mathbf{y} in an ascending order with permutation σ . Applying σ to \mathbf{x} we have $\sigma(\mathbf{x}) = (1, 2, 5, 3, 4)$. The total number of increasing subsequences of length 3 in $\sigma(\mathbf{x})$ is 5, and there are no decreasing subsequences of length 3. Hence $W_2 = 5$.

2 Asymptotic Convergence

Figure S1 shows the convergence of the empirical quantiles of T_1 and T_2 toward the theoretical standard normal quantiles as n increases. Note due to the fact that T_1 can only take $n - k + 2$ possible values, it is easy to produce ties. To examine the asymptotic power of the two statistics under alternative distributions described previously, we generated data that i) were partially coupled time series with the length of dependence $m = n/10$; ii) followed an exact functional relationship with six monotonic pieces, and computed the average power at 5% significance level over 500 iterations. The results for different k and n are shown in Table S1. As predicted by the theoretical analysis, larger k results in better power and T_1 is more powerful than T_2 on the time-course data. In all the cases, as n increases the power tends to 1. The table also displays the average power for the corresponding null distributions of i) and ii) when the two data vectors are independent. Some values are slightly larger than 0.05 due to the heavier tails of the empirical distributions.

3 Simulations

The estimates for Hoeffding's D, dCov, the Renyi correlation, MI and MIC were computed using relevant R packages (`Hmisc` [1], `energy` [2], `acepack` [3], `parmgene` [4] and `minerva` [5]). We used the

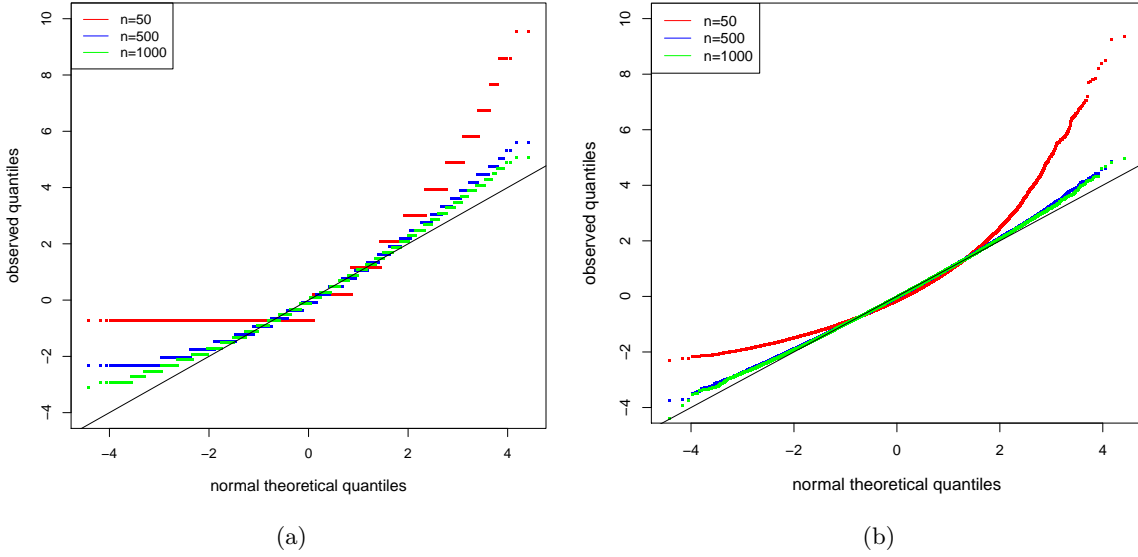


Figure S1: Empirical quantiles for the standardized counts (a) T_1 and (b) T_2 for $n = 50, 500$ and 1000 , $k = 5$, from 10^5 simulated random permutations.

k/n	100	200	300	400	500
3	0.332	0.562	0.690	0.812	0.902
4	0.636	0.976	1	1	1
5	1	1	1	1	1
6	1	1	1	1	1

(a) Power of T_1

k/n	100	200	300	400	500
3	0.302	0.504	0.658	0.796	0.848
4	0.340	0.568	0.726	0.844	0.908
5	0.360	0.650	0.798	0.892	0.952
6	0.392	0.734	0.882	0.924	0.982

(c) Power of T_2

k/n	100	200	300	400	500
3	0.516	0.784	0.884	0.952	0.972
4	0.710	0.952	0.996	1	1
5	0.866	0.992	1	1	1
6	0.946	1	1	1	1

(e) Power of T_2

k/n	100	200	300	400	500
3	0.070	0.050	0.056	0.038	0.036
4	0.060	0.044	0.030	0.054	0.052
5	0.038	0.074	0.064	0.058	0.060
6	0.052	0.042	0.064	0.046	0.076

(b) Power of T_1

k/n	100	200	300	400	500
3	0.048	0.054	0.076	0.052	0.050
4	0.068	0.060	0.044	0.072	0.040
5	0.042	0.056	0.040	0.070	0.038
6	0.050	0.068	0.056	0.046	0.046

(d) Power of T_2

k/n	100	200	300	400	500
3	0.058	0.048	0.062	0.044	0.052
4	0.058	0.062	0.044	0.060	0.052
5	0.078	0.040	0.060	0.062	0.030
6	0.070	0.060	0.064	0.062	0.054

(f) Power of T_2

Table S1: Power at 5% significance level for different choices of k and n when \mathbf{x} and \mathbf{y} are: (a), (c) two independent AR(1) time series (with coefficients 0.1 and -0.2 respectively) but $(x_1, \dots, x_m) = (y_1, \dots, y_m)$ with $m = n/10$; (e) $x_i \stackrel{iid}{\sim} Unif(0, 1)$, $y_i = \cos(6\pi x_i)$. The right panel shows the power under the corresponding null distributions: (b), (d) \mathbf{x} and \mathbf{y} are two independent AR(1) time series (with coefficients 0.1 and -0.2 respectively); (f) \mathbf{x} and \mathbf{y} are iid $Unif(0, 1)$.

standard ACE estimate ([6]) for approximating the Renyi correlation. The computation of some of the measures involve tuning parameters. The MI estimates were computed using the k th nearest neighbor (KNN) algorithm of [7]. A number of bandwidth parameters were tried (6, 10 and 20) and the results corresponding to the one with the best power (20) are shown. The MIC estimates were computed using the R package `minerva` with default parameter settings. For statistics with unknown asymptotic distributions (dCov, ACE and MI), the p-values were calculated by a permutation procedure. For each dataset generated the same statistics were calculated on a null dataset obtained by permuting the orders of y_i . The power was taken to be the fraction of datasets with a statistic value more significant than 95% of the values produced by the permuted datasets. Pre-computed p-values of MIC were downloaded from <http://www.exploredata.net/Downloads/P-Value-Tables>.

Descriptions of the parameters used for the four types of dependence relationships are given in Table S2. 2000 datasets were generated for every scenario with $i \in \{1, \dots, 220\}$, $e_i \stackrel{iid}{\sim} N(0, 1)$ for the first three relationships and $e_i \stackrel{iid}{\sim} N(0, 0.5)$ for the time-course relationship. Outliers were created by randomly choosing a fraction of the data and replacing e_i with η_i .

	x_i	y_i	η_i
Linear	$x_i \stackrel{iid}{\sim} N(0, 1)$	$y_i = x_i + 2e_i$	$\eta_i \stackrel{iid}{\sim} N(0, 5)$
Quadratic	$x_i \stackrel{iid}{\sim} N(0, 1)$	$y_i = x_i^2 + 2e_i$	$\eta_i \stackrel{iid}{\sim} N(0, 5)$
Cross	$x_i \stackrel{iid}{\sim} N(0, 1)$	$y_i = \begin{cases} \frac{1}{2} + x_i + e_i & \text{with probability } \frac{1}{2}, \\ \frac{3}{2} - x_i + e_i & \text{with probability } \frac{1}{2}. \end{cases}$	$\eta_i \stackrel{iid}{\sim} N(0, 3)$
Partially coupled time series	$x_i \sim AR(1)$ with coefficient 0.1	$y_i = \begin{cases} x_i + e_i, & i \in [1, 30] \\ -x_i + e_i, & i \in [101, 120] \\ AR(1) \text{ with coefficient } -0.2, & \\ \text{independent of } x_i, & \text{otherwise.} \end{cases}$	$\eta_i \stackrel{iid}{\sim} N(0, 3)$

Table S2: Parameters for generating the four types of relationships

Power curves for T_1 and T_2 with different choices of k are shown in Figure S2.

We additionally investigated the power loss on a linear relationship with increasing noise level but no outliers. The results are plotted in Figure S3. The linear relationship was generated with $y = x + \beta e$, where $x \stackrel{iid}{\sim} N(0, 1)$, $e \stackrel{iid}{\sim} N(0, 1)$ and $\beta \in \{1, 2, \dots, 10\}$. T_2^+ remains the best performing statistic. As expected, Pearson's correlation shows better performance on data with no outliers and is now ranked the second. T_2 still demonstrates less power than Spearman's correlation, dCov and Hoeffding'D, but remains more powerful than Renyi's correlation and MI.

We provide a power comparison between our statistics and the LIS-based statistics ([8]) computed using their R package `LISest` on simulated data in Figure S4. The four scenarios used the same parameters as described in Table S1, except with $n = 200$ — the largest n allowed by the R package. LIS_Ln represents the LIS; LIS_JLn uses a jackknife version of LIS; and LIS_JLMn uses the longest monotonic subsequence (the maximum of LIS and the longest decreasing subsequence). Overall the power of this class of statistics is not optimal, which can be explained by the relative non-robustness of the length of LIS in the presence of noise and outliers. Furthermore, intuitively LIS-based statistics are better suited to detect global monotonic relationships, which differs from our consideration of potentially changing local dependence patterns. As Figure S4 confirms, their power values are lower for non-monotonic relationships ((b), (c) and (d)).

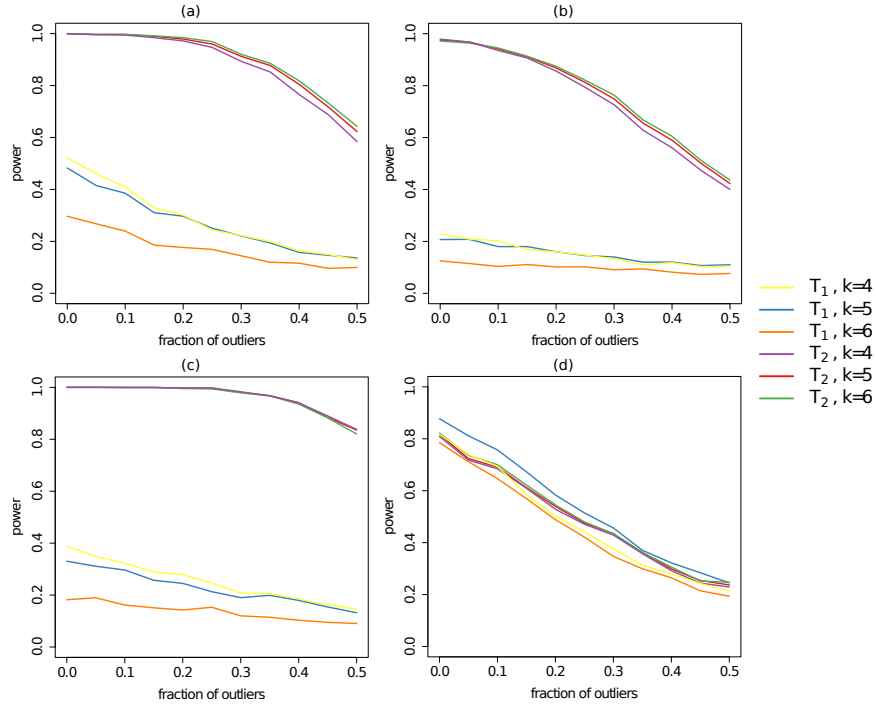


Figure S2: The power of T_1 and T_2 for various k values rejecting at 5% significance level as level of contamination by outliers increases when the bivariate data have (a) a linear relationship; (b) a quadratic relationship; (c) a cross-shaped relationship; (d) are two partially coupled time series.

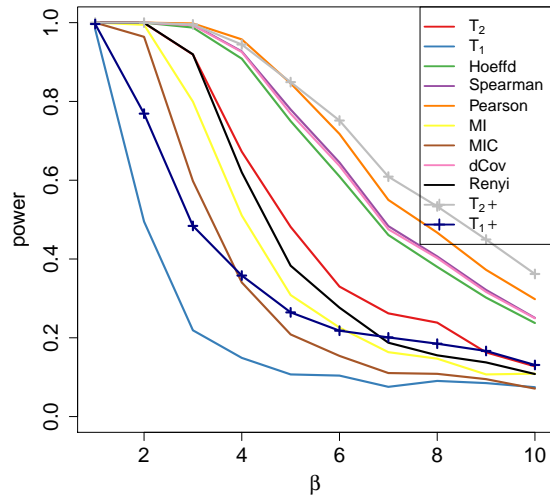


Figure S3: The power of various statistics rejecting at 5% significance level as the noise level increases on a linear relationship.

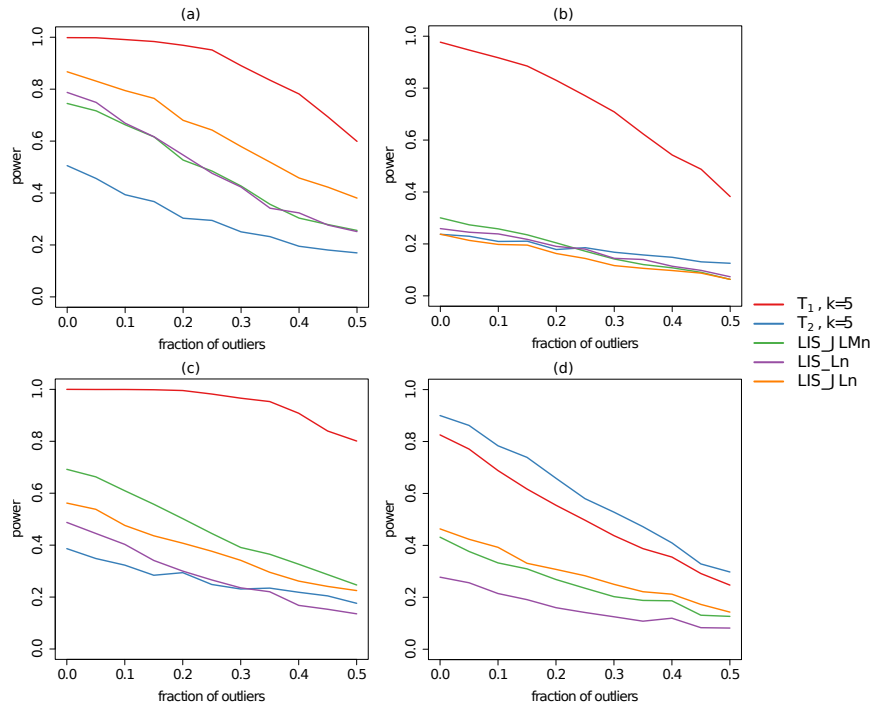


Figure S4: The power of T_1 and T_2 against LIS-based statistics rejecting at 5% significance level as level of contamination by outliers increases when the bivariate data follow (a) a linear relationship; (b) a quadratic relationship; (c) a cross-shaped relationship; (d) two partially coupled time series.

4 Yeast Cell Cycle

The yeast expression data was accessed from <http://genome-www.stanford.edu/celcycle/> and contains the expression levels of 6178 genes from four reasonably long time-course experiments: alpha factor release (18 time points), cdc 15 (24 time points), cdc 28 (17 time points) and elutriation (14 time points). We linearly interpolated some missing data if a point had the two adjacent time points belonging to the same experiment with no missing values. We focused on the coexpression of 133 transcription factors (TFs) with no missing data after interpolation. Since the data has a number of ties, we added small random perturbations for the computation of T_1 and T_2 and took the final results as the maximum counts over 50 iterations.

Figure S5 shows two pairs of TFs (MOT3 and RPN4; PHO2 and SUT1) with genetic interaction identified by T_1 but missed by all the other methods.

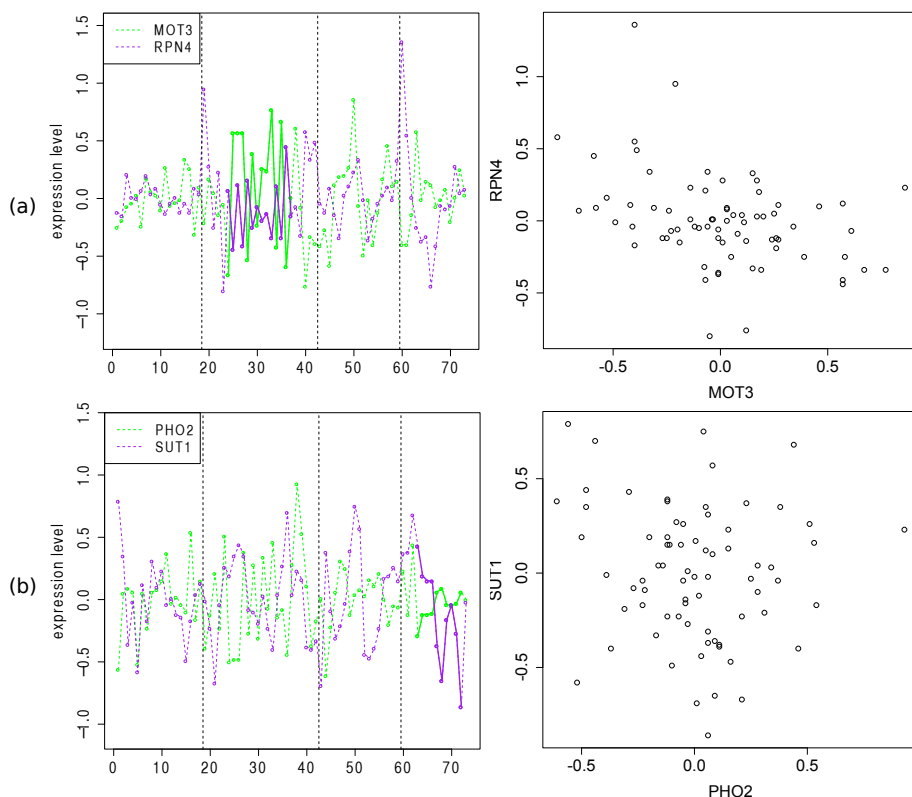


Figure S5: Expression levels of (a) MOT3 and RPN4; (b) PHO2 and SUT1 in four time-course experiments (boundaries indicated by the dashed lines). The solid lines highlight regions contributing to the counts in T_1 . Both have reported genetic interactions ([9, 10]), but received low rankings under methods other than T_1 .

Table S3 shows the number of known TF interactions among strongly coexpressed pairs as ranked by each method. A number of k values were tried for T_1 , while for T_2 only $k = 7$ was shown since the results were quite stable over a range of k values. As T_1 led to many ties, the cutoffs were chosen to include the entire stretches of gene pairs with the same statistic values.

Top rank	$k = 6$				$k = 7$			$k = 8$			$k = 9$		
	4	7	16	31	4	11	22	5	14	44	3	11	37
Pearson	0	2	2	3	0	2	2	1	2	6	0	2	4
Spearman	0	1	2	2	0	1	2	0	1	2	0	1	2
Hoeffding's D	1	1	1	2	1	1	2	1	1	2	0	1	2
MI	0	1	1	1	0	1	1	1	1	1	0	1	1
MIC	1	1	1	1	1	1	1	1	1	2	1	1	2
dCov	1	1	1	2	1	1	1	1	1	2	1	1	2
Renyi	0	0	2	2	0	2	2	0	2	3	0	2	2
T_1	0	1	3	3	1	3	4	1	3	6	0	1	5
T_2	0	1	2	3	0	2	2	0	2	3	0	2	3

Table S3: Number of known interactions in highly ranked coexpression pairs by various statistics. A range of k values were tested for T_1 , and $k = 7$ for T_2 .

5 *Arabidopsis* Microarrays

The original CEL files of the microarrays were downloaded from NCBI GEO (GSE 5623, 7636, 7639, 7641, 7642, 8787 and 30166), then normalized using the robust multi-array analysis (RMA) ([11]) function in the Bioconductor package. After normalization, a small fraction of the data were tied. We added small random perturbations for the computation of T_1 and T_2 and took the final results as the maximum count over 20 iterations. We noted that this had negligible influence on all the final results. Asymptotic p-values were computed for T_1 , T_2 , the Pearson correlation, the Spearman correlation and Hoeffding's D. For dCov, Renyi and MI, null statistic values were calculated by permuting the sample labels of each gene and used as empirical quantiles for determining the significance level of the statistics. Pre-computed p-values for MIC from <http://www.exploredata.net/Downloads/P-Value-Tables> were used.

Figure S7 shows two pairs of genes in the same pathway, where the bulk of the samples follow a linear trend but they failed to be identified by MI at an unadjusted significance level of 5%. On the other hand, both pairs were assigned significant p-values by T_2 and other statistics including the Pearson and Spearman correlations.

For each pathway, we ranked the coexpression between the pathway genes and all the genes available and chose the top L pairs, where L is the number of total gene pairs in this pathway. We then counted the number of gene pairs belonging to this pathway among the chosen pairs, and kept 20 pathways in which at least one method achieved a significant enrichment of pathway genes using Fisher's exact test. Table S4 tallies the methods with the highest counts of same pathway pairs in these 20 pathways.

6 Proofs

6.1 Running time of the algorithms

Lemma 6.1. *Computing W_1 and W_2 takes $O(k(\log k)n)$ and $O(kn \log n)$ time respectively.*

Proof. Computing W_1 involves ranking and comparing the elements of vectors of length k $O(n)$ number of times, thus the running time is $O(k(\log k)n)$.

W_2 counts the total number of subsequences of length k with matching or reverse rank patterns. For any pair of subsequences with matching rank patterns, permuting the two subsequences simultaneously to sort one of them in an increasing order will also sort the other one in an increasing order. Using this observation, let σ be the permutation that sorts \mathbf{y} in an increasing order and $\mathbf{z} = \sigma(\mathbf{x})$ be that

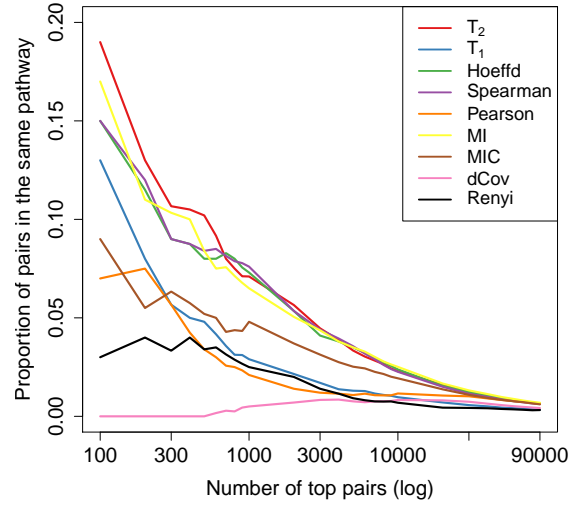


Figure S6: Proportion of gene pairs in the same pathway as a function of the number of highly ranked pairs chosen.

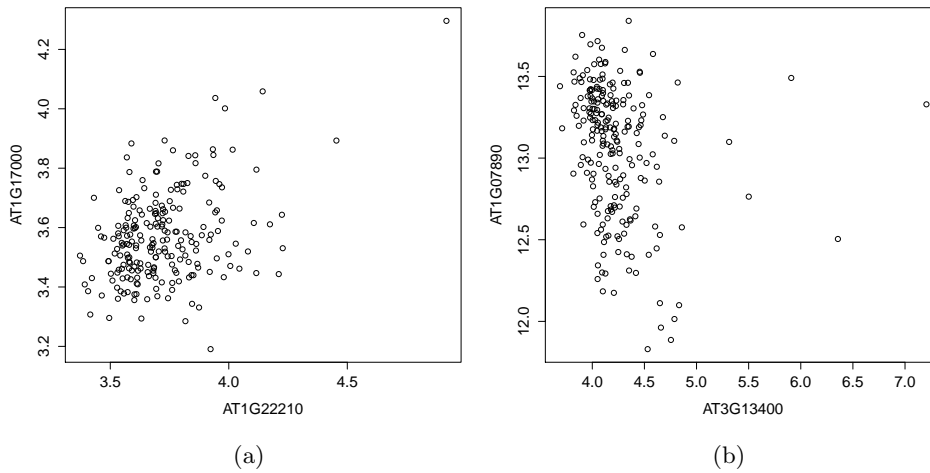


Figure S7: Expression levels of two gene pairs in the same pathway showing a linear relationship with outliers which were not identified as statistically significant by MI.

permutation applied to \mathbf{x} . Then W_2 is the number of increasing (and decreasing) subsequences of length k in \mathbf{z} . To compute W_2 , it suffices to consider counting the increasing subsequences. One obvious solution is dynamic programming. Let $dp[i,1]$ be the number of increasing subsequences of length 1 ending at position i , then the matrix $dp[i,1]$ can be updated as follows.

```

Initialize  $dp[i,1] = 0$ ;  $dp[i,1] = 1$ 
for  $i = 2$  to  $n$ 
  for  $j = 1$  to  $i-1$ 
    if  $z[i] > z[j]$ 

```

	Individual pathways														Tally	
Renyi	✓														1	
dcov															0	
Hoeffd	✓				✓			✓								3
MI	✓						✓		✓					✓		4
MIC															0	
Pearson			✓					✓	✓							3
Spearman	✓						✓			✓						3
T_1				✓		✓							✓		✓	4
T_2	✓	✓		✓		✓	✓		✓	✓		✓	✓	✓	✓	12

Table S4: Methods with the highest counts of pathway genes pairs in 20 pathways with statistically significant enrichment

for $l = 2$ to k
 $dp[i,l] += dp[j,l-1]$

The final answer is obtained by summing $dp[i,k]$ over i . It is easy to see this has a running time of $O(kn^2)$. Note that in the second loop the only entries involved in the update are $z[j]$ whose ranks are smaller than that of $z[i]$. Therefore by first ranking the elements in z , a binary indexed tree structure can be implemented to perform the sum and update efficiently, reducing the running time to $O(kn \log n)$ ([12]). \square

6.2 Asymptotic distributions of W_1

Throughout the sections, C and C_i denote positive constants which may be different at each appearance. Without loss of generality assume \mathbf{x} satisfies the assumption that it has an exchangeable distribution. Then the ranks of any subsequence of \mathbf{x} can be treated as a random permutation. Denote

$$\begin{aligned}
\mathbb{I}_i^+ &= \mathbb{I}(\phi(x_i, \dots, x_{i+k-1}) = \phi(y_i, \dots, y_{i+k-1})), \\
\mathbb{I}_i^- &= \mathbb{I}(\phi(x_i, \dots, x_{i+k-1}) = \phi(-y_i, \dots, -y_{i+k-1})), \\
\mathbb{I}_i &= \mathbb{I}_i^+ + \mathbb{I}_i^-.
\end{aligned} \tag{S1}$$

We have

$$\begin{aligned}
&\mathbb{E}(\mathbb{I}_i^+) \\
&= \sum_{\mathbf{w}} \mathbb{P}(\phi(x_i, \dots, x_{i+k-1}) = \mathbf{w} \mid \phi(y_i, \dots, y_{i+k-1}) = \mathbf{w}) \mathbb{P}(\phi(y_i, \dots, y_{i+k-1}) = \mathbf{w}) \\
&= \frac{1}{k!} \sum_{\mathbf{w}} \mathbb{P}(\phi(y_i, \dots, y_{i+k-1}) = \mathbf{w}) \\
&= \frac{1}{k!}
\end{aligned} \tag{S2}$$

by the independence assumption and the fact that there is only one way to arrange a list of numbers in a given order. Clearly also $\mathbb{E}(\mathbb{I}_i^-) = 1/k!$. In the next lemma, we characterize the behavior of the cross terms $\mathbb{E}(\mathbb{I}_i^+ \mathbb{I}_j^+)$.

Lemma 6.2. 1. When $|j - i| \geq k$, \mathbb{I}_i^+ and \mathbb{I}_j^+ are independent. So are $(\mathbb{I}_i^+, \mathbb{I}_j^-)$ and $(\mathbb{I}_i^-, \mathbb{I}_j^-)$.

2. When $|j - i| = k - l$ with $1 \leq l \leq k - 1$,

$$\frac{1}{(2k - l)!} \leq \mathbb{E}(\mathbb{I}_i^+ \mathbb{I}_j^+) \leq \frac{\binom{2k-2l}{k-l}}{(2k - l)!}. \quad (\text{S3})$$

The same conclusions hold for $(\mathbb{I}_i^-, \mathbb{I}_j^-)$, $1 \leq |j - i| < k$, and $(\mathbb{I}_i^+, \mathbb{I}_j^-)$, $(\mathbb{I}_i^-, \mathbb{I}_j^+)$, $|i - j| = k - 1$.

3. $\mathbb{E}(\mathbb{I}_i^+ \mathbb{I}_j^-) = \mathbb{E}(\mathbb{I}_i^- \mathbb{I}_j^+) = 0$ for $1 \leq |i - j| < k - 1$.

Proof. Note that conditioning on the sequence \mathbf{y} ,

$$\begin{aligned} \mathbb{E}(\mathbb{I}_i^+ \mathbb{I}_j^+) &= \sum_{\mathbf{w}, \mathbf{v}} \mathbb{P}(\phi(x_i, \dots, x_{i+k-1}) = \mathbf{w}, \phi(x_j, \dots, x_{j+k-1}) = \mathbf{v}) \\ &\quad \times \mathbb{P}(\phi(y_i, \dots, y_{i+k-1}) = \mathbf{w}, \phi(y_j, \dots, y_{j+k-1}) = \mathbf{v}). \end{aligned} \quad (\text{S4})$$

For $|j - i| \geq k$, the subsequences (x_i, \dots, x_{i+k-1}) and (x_j, \dots, x_{j+k-1}) do not overlap. Thus their local rank patterns are independent, each having probability $1/k!$ for a given order.

$$\begin{aligned} \mathbb{E}(\mathbb{I}_i^+ \mathbb{I}_j^+) &= \left(\frac{1}{k!}\right)^2 \sum_{\mathbf{w}, \mathbf{v}} P(\phi(y_i, \dots, y_{i+k-1}) = \mathbf{w}, \phi(y_j, \dots, y_{j+k-1}) = \mathbf{v}) \\ &= \left(\frac{1}{k!}\right)^2 = \mathbb{E}(\mathbb{I}_i^+) \mathbb{E}(\mathbb{I}_j^+). \end{aligned} \quad (\text{S5})$$

For $j - i = k - l < k$ (assuming WLOG $j > i$), (x_i, \dots, x_{i+k-1}) and (x_j, \dots, x_{j+k-1}) form a contiguous subsequence $x_i, \dots, x_j, \dots, x_{j+k-1}$. Suppose $\phi(x_i, \dots, x_{j+k-1}) = (u_1, \dots, u_{2k-l})$, then

$$\begin{aligned} \phi(u_1, \dots, u_k) &= (w_1, \dots, w_k), \\ \phi(u_{k-l+1}, \dots, u_{2k-l}) &= (v_1, \dots, v_k), \\ \phi(u_{k-l+1}, \dots, u_k) &= \phi(w_{k-l+1}, \dots, w_k) = \phi(v_1, \dots, v_l) \\ &:= (o_1, \dots, o_l), \quad \text{say.} \end{aligned} \quad (\text{S6})$$

Focusing on the overlapping part (u_{k-l+s}) for $1 \leq s \leq l$, the numbers of elements smaller than u_{k-l+s} in the subsequences (u_1, \dots, u_k) , $(u_{k-l+1}, \dots, u_{2k-l})$ and (u_{k-l+1}, \dots, u_k) are $w_{k-l+s} - 1$, $v_s - 1$, and $o_s - 1$, respectively. Given the overall rank u_{k-l+s} in the sequence $(u_1, \dots, u_{k-l+1}, \dots, u_k, \dots, u_{2k-l})$, we have

$$u_{k-l+s} - 1 = (w_{k-l+s} - 1) + (v_s - 1) - (o_s - 1), \quad (\text{S7})$$

since the elements in the overlapping part are counted twice. In other words, the overlapping part (u_{k-l+s}) for $1 \leq s \leq l$ is fixed, and there are at most $\binom{2k-2l}{k-l}$ ways of arranging the rest $2k - 2l$ numbers. Thus we arrive at the upper bound in (S3). The lower bound is trivial. The same arguments hold for $(\mathbb{I}_i^-, \mathbb{I}_j^-)$, $1 \leq |j - i| < k$, and $(\mathbb{I}_i^+, \mathbb{I}_j^-)$, $(\mathbb{I}_i^-, \mathbb{I}_j^+)$, $|i - j| = k - 1$.

Lastly, for $1 \leq |i - j| < k - 1$, $\mathbb{E}(\mathbb{I}_i^+ \mathbb{I}_j^-) = \mathbb{E}(\mathbb{I}_i^- \mathbb{I}_j^+) = 0$ since no such arrangements of the elements are possible. \square

Let N_i denote the dependency neighborhood of \mathbb{I}_i , the next lemma tries to bound a key quantity in the variance calculation.

Lemma 6.3. *For all $k \geq 3$,*

$$4(n - 2k + 2) \left(\sum_{l=2}^{k-1} \frac{1}{(2k - l)!} + \frac{2}{(2k - 1)!} \right) \leq \sum_{i=1}^{n-k+1} \sum_{j \in N_i \setminus \{i\}} \mathbb{E}(\mathbb{I}_i \mathbb{I}_j) \leq \frac{C(n - k + 1)}{(k + 1)!} \quad (\text{S8})$$

for some $C > 0$.

Proof. First note that

$$\begin{aligned} \sum_{i=1}^{n-k+1} \sum_{j \in N_i \setminus \{i\}} \mathbb{E}(\mathbb{I}_i \mathbb{I}_j) &= 2 \sum_{i=1}^{n-k+1} \sum_{j \in N_i \setminus \{i\}} \mathbb{E}(\mathbb{I}_i^+ \mathbb{I}_j^+) + 2 \sum_{i=1}^{n-k+1} \sum_{|j-i|=k-1} \mathbb{E}(\mathbb{I}_i^+ \mathbb{I}_j^-) \\ &\leq 8(n-k+1) \sum_{l=1}^{k-1} \gamma_l \end{aligned} \quad (\text{S9})$$

by (S3), where

$$\gamma_l = \frac{\binom{2k-2l}{k-l}}{(2k-l)!}. \quad (\text{S10})$$

It remains to bound $\sum_{l=1}^{k-1} \gamma_l$. Taking the ratio of successive terms,

$$\begin{aligned} r_l &= \frac{\gamma_{l+1}}{\gamma_l} = \frac{\binom{2k-2l-2}{k-l-1}}{(2k-l-1)!} \cdot \frac{(2k-l)!}{\binom{2k-2l}{k-l}} \\ &= \frac{(k-l)^2(2k-l)}{(2k-2l)(2k-2l-1)} \\ &= \frac{(k-l)(2k-l)}{2(2k-2l-1)}, \quad l = 1, \dots, k-2. \end{aligned} \quad (\text{S11})$$

For all $k \geq 3$, there exists positive constant C_1 and C_2 (independent of k) such that

$$C_1 k \leq r_l \leq C_2 k, \quad l = 1, \dots, k-2. \quad (\text{S12})$$

Therefore $\sum_{l=1}^{k-1} \gamma_l$ is upper bounded by

$$\begin{aligned} \sum_{l=1}^{k-1} \gamma_l &\leq \gamma_{k-1} \sum_{l=0}^{k-2} \left(\frac{1}{C_1 k} \right)^l \\ &= \gamma_{k-1} \cdot \frac{1 - \left(\frac{1}{C_1 k} \right)^{k-1}}{1 - \frac{1}{C_1 k}} \\ &\leq \frac{C}{(k+1)!} \end{aligned} \quad (\text{S13})$$

for some $C > 0$. Equations (S13) and (S9) give the required upper bound.

For the lower bound, it is easy to see

$$\begin{aligned} \sum_{i=1}^{n-k+1} \sum_{j \in N_i \setminus \{i\}} \mathbb{E}(\mathbb{I}_i \mathbb{I}_j) &= 2 \sum_{i=1}^{n-k+1} \sum_{j \in N_i \setminus \{i\}} \mathbb{E}(\mathbb{I}_i^+ \mathbb{I}_j^+) + 2 \sum_{i=1}^{n-k+1} \sum_{|j-i|=k-1} \mathbb{E}(\mathbb{I}_i^+ \mathbb{I}_j^-) \\ &\geq 2 [2(n-k+1) - 2(k-1)] \left(\sum_{l=1}^{k-1} \frac{1}{(2k-l)!} + \frac{1}{(2k-1)!} \right) \\ &\geq 4(n-2k+2) \left(\sum_{l=2}^{k-1} \frac{1}{(2k-l)!} + \frac{2}{(2k-1)!} \right) \end{aligned} \quad (\text{S14})$$

by the lower bound in (S3). □

With the above bounds we can now prove Theorem 1.

Proof of Theorem 1. In order to use Stein's method for normal approximation, we first give a lower bound of the variance. Note that

$$\begin{aligned}
\sigma_{1,n}^2 &= \sum_{i=1}^{n-k+1} \sum_{j \in N_i} (\mathbb{E}(\mathbb{I}_i \mathbb{I}_j) - (\mathbb{E}\mathbb{I}_i)(\mathbb{E}\mathbb{I}_j)) \\
&= \sum_{i=1}^{n-k+1} \sum_{j \in N_i \setminus \{i\}} \mathbb{E}(\mathbb{I}_i \mathbb{I}_j) + \frac{2(n-k+1)}{k!} - \sum_{i=1}^{n-k+1} \sum_{j \in N_i} \mathbb{E}(\mathbb{I}_i) \mathbb{E}(\mathbb{I}_j) \\
&\geq 4(n-2k+2) \left(\sum_{l=2}^{k-1} \frac{1}{(2k-l)!} + \frac{2}{(2k-1)!} \right) + \frac{2(n-k+1)}{k!} - \frac{4(n-k+1)(2k-1)}{(k!)^2}. \tag{S15}
\end{aligned}$$

by (S8). For k such that $k/n \rightarrow 0$, when n is sufficiently large, $\sigma_{1,n}^2$ is lower bounded by the dominating terms

$$\begin{aligned}
\sigma_{1,n}^2 &\geq C_1 \left(4n \left(\sum_{l=2}^{k-1} \frac{1}{(2k-l)!} + \frac{2}{(2k-1)!} \right) + \frac{2n}{k!} - \frac{4n(2k-1)}{(k!)^2} \right) \\
&= \frac{2C_1 n}{k!} \left(2 \left(\frac{1}{k+1} + \frac{1}{(k+2)(k+1)} + \dots + \frac{2}{(2k-1) \dots (k+1)} \right) + 1 - \frac{2(2k-1)}{k!} \right) \\
&\geq \frac{C_2 n}{k!} \tag{S16}
\end{aligned}$$

for some $C_1, C_2 > 0$ and all $k \geq 3$. One version of Stein's method gives the following error bound for normal approximation ([13]),

$$d_W(T_1, Y) \leq \frac{D^2}{\sigma_{1,n}^3} \sum_{i=1}^{n-k+1} \mathbb{E}|\mathbb{I}_i - 2/k!|^3 + \frac{\sqrt{26}D^{3/2}}{\sqrt{\pi}\sigma_{1,n}^2} \sqrt{\sum_{i=1}^{n-k+1} \mathbb{E}|\mathbb{I}_i - 2/k!|^4} \tag{S17}$$

where d_W is the Wasserstein metric, $Y \sim N(0, 1)$ and $D = \max_i N_i = 2k-1$. This can be further bounded by

$$\begin{aligned}
&C_1 \cdot \frac{D^2 \mu_{1,n}}{\sigma_{1,n}^3} + C_2 \cdot \frac{D^{3/2} \mu_{1,n}^{1/2}}{\sigma_{1,n}^2} \\
&\leq C \cdot \frac{D^2 \mu_{1,n}}{\sigma_{1,n}^3} \\
&\leq C \cdot \frac{k^2 \sqrt{k!}}{\sqrt{n}} \rightarrow 0 \tag{S18}
\end{aligned}$$

using (S16) for $k/(\log n)^\alpha \rightarrow 0$, $\alpha < 1$.

The Chen-Stein method yields the following error bound for Poisson approximation,

$$\begin{aligned}
d_{TV}(W_1, Z) &\leq \min\{1, \mu_{1,n}^{-1}\} \left(\sum_{i=1}^{n-k+1} \sum_{j \in N_i} \mathbb{E}(\mathbb{I}_i) \mathbb{E}(\mathbb{I}_j) + \sum_{i=1}^{n-k+1} \sum_{j \in N_i \setminus \{i\}} \mathbb{E}(\mathbb{I}_i \mathbb{I}_j) \right) \\
&\leq \sum_{i=1}^{n-k+1} \sum_{j \in N_i} \mathbb{E}(\mathbb{I}_i) \mathbb{E}(\mathbb{I}_j) + \sum_{i=1}^{n-k+1} \sum_{j \in N_i \setminus \{i\}} \mathbb{E}(\mathbb{I}_i \mathbb{I}_j) \\
&\leq \frac{4(n-k+1)(2k-1)}{(k!)^2} + \frac{C(n-k+1)}{(k+1)!} \\
&\leq \frac{C(n-k+1)}{(k+1)!}
\end{aligned} \tag{S19}$$

for some $C > 0$ and k sufficiently large. For k growing fast enough such that $\mu_{1,n} = O(1)$, the above bound goes to 0. In particular, using Stirling's approximation one can show in the regime $\log n/k = O(1)$ this condition is satisfied. \square

6.3 Asymptotic distribution of W_2

Assuming \mathbf{x} has an exchangeable distribution, the permuted sequence $\sigma(\mathbf{x})$ also has an exchangeable distribution, and its ranks can be treated as a random permutation. For notational simplicity, take \mathbf{z} as a random permutation of $\{1, \dots, n\}$. For integers $\{i_1, \dots, i_k\}$ satisfying $1 \leq i_1 < \dots < i_k \leq n$, define indicators $\mathbb{I}_{i_1, \dots, i_k}^+(\mathbf{z})$ such that

$$\mathbb{I}_{i_1, \dots, i_k}^+(\mathbf{z}) = \begin{cases} 1 & (i_1, \dots, i_k) \text{ is a subsequence of } \mathbf{z}, \\ 0 & \text{otherwise.} \end{cases} \tag{S20}$$

Similarly define

$$\mathbb{I}_{i_1, \dots, i_k}^-(\mathbf{z}) = \begin{cases} 1 & (i_k, \dots, i_1) \text{ is a subsequence of } \mathbf{z}, \\ 0 & \text{otherwise.} \end{cases} \tag{S21}$$

Then W_2 can be written as the sum of

$$W_2 = \sum_{1 \leq i_1 < \dots < i_k \leq n} \mathbb{I}_{i_1, \dots, i_k}(\mathbf{z}), \tag{S22}$$

where

$$\mathbb{I}_{i_1, \dots, i_k}(\mathbf{z}) = \mathbb{I}_{i_1, \dots, i_k}^+(\mathbf{z}) + \mathbb{I}_{i_1, \dots, i_k}^-(\mathbf{z}). \tag{S23}$$

It is easy to see that if $\{i_1, \dots, i_k\} \cap \{j_1, \dots, j_k\} = \emptyset$, $\mathbb{I}_{i_1, \dots, i_k}(\mathbf{z})$ and $\mathbb{I}_{j_1, \dots, j_k}(\mathbf{z})$ are independent. The variance of W_2 becomes

$$\begin{aligned}
&\text{Var}(W_2) \\
&= \sum_{\{i_1, \dots, i_k\} \cap \{j_1, \dots, j_k\} \neq \emptyset} \{ \mathbb{E}(\mathbb{I}_{i_1, \dots, i_k}(\mathbf{z}) \mathbb{I}_{j_1, \dots, j_k}(\mathbf{z})) - \mathbb{E}(\mathbb{I}_{i_1, \dots, i_k}(\mathbf{z})) \mathbb{E}(\mathbb{I}_{j_1, \dots, j_k}(\mathbf{z})) \} \\
&= 2 \sum_{\{i_1, \dots, i_k\} \cap \{j_1, \dots, j_k\} \neq \emptyset} \mathbb{E}(\mathbb{I}_{i_1, \dots, i_k}^+(\mathbf{z}) \mathbb{I}_{j_1, \dots, j_k}^+(\mathbf{z})) \\
&\quad + 2 \sum_{\{i_1, \dots, i_k\} \cap \{j_1, \dots, j_k\} \neq \emptyset} \mathbb{E}(\mathbb{I}_{i_1, \dots, i_k}^+(\mathbf{z}) \mathbb{I}_{j_1, \dots, j_k}^-(\mathbf{z})) - \frac{4D \binom{n}{k}}{(k!)^2},
\end{aligned} \tag{S24}$$

since $\mathbb{E}(\mathbb{I}^+(z_{i_1}, \dots, z_{i_k})) = 1/k!$. Here D is the size of the dependency neighborhood and equals $\binom{n}{k} - \binom{n-k}{k}$. The sum of the first cross terms can be written as (Proposition 2 in [14])

$$\begin{aligned} & \sum_{\{i_1, \dots, i_k\} \cap \{j_1, \dots, j_k\} \neq \emptyset} \mathbb{E}(\mathbb{I}_{i_1, \dots, i_k}^+(\mathbf{z}) \mathbb{I}_{j_1, \dots, j_k}^+(\mathbf{z})) \\ &= \sum_{j=1}^k \binom{n}{2k-j} \frac{1}{(2k-j)!} A(k-j, j), \end{aligned} \quad (\text{S25})$$

where

$$A(N, j) = \sum_{\substack{\sum_{r=0}^j l_r = N \\ \sum_{r=0}^j m_r = N}} \prod_{r=0}^j \left(\frac{(l_r + m_r)!}{l_r! m_r!} \right)^2. \quad (\text{S26})$$

We will be using the following fact about the constants $A(N, j)$ from Lemma 3 in [14].

Fact 6.4. For sufficiently large k , there exists $C > 0$ such that

$$A(k-1, 1) \geq Ck^{1/2} \binom{2k-2}{k-1}^2. \quad (\text{S27})$$

It is easy to see for all $k \geq 2$, $A(k-1, 1) > \binom{2k-2}{k-1}^2$.

The sum of the second cross terms reduces to

$$\sum_{|\{i_1, \dots, i_k\} \cap \{j_1, \dots, j_k\}|=1} \mathbb{E}(\mathbb{I}_{i_1, \dots, i_k}^+(\mathbf{z}) \mathbb{I}_{j_1, \dots, j_k}^-(\mathbf{z})),$$

since when the size of the intersection is greater than one, it is impossible to find a permutation \mathbf{z} satisfying both conditions specified by the indicators. Using arguments similar to the proof of Proposition 2 in [14], we can show

$$\sum_{|\{i_1, \dots, i_k\} \cap \{j_1, \dots, j_k\}|=1} \mathbb{E}(\mathbb{I}_{i_1, \dots, i_k}^+(\mathbf{z}) \mathbb{I}_{j_1, \dots, j_k}^-(\mathbf{z})) = \binom{n}{2k-1} \frac{1}{(2k-1)!} B(k), \quad (\text{S28})$$

where

$$B(k) = \sum_{\substack{l_0 + l_1 = k-1 \\ m_0 + m_1 = k-1}} \binom{l_0 + m_0}{l_0} \binom{l_1 + m_1}{l_1} \binom{l_0 + m_1}{l_0} \binom{l_1 + m_0}{l_1} \quad (\text{S29})$$

Now we can obtain a lower bound on the variance and use the Stein method to prove Theorem 2.

Proof of Theorem 2. From equations (S24), (S25) and (S28), we have

$$\begin{aligned} \frac{\sigma_{2,n}^2}{\mu_{2,n}^2} &\geq \frac{\binom{n}{2k-1} (k!)^2}{2 \binom{n}{k}^2 (2k-1)!} (A(k-1, 1) + B(k)) - \frac{D}{\binom{n}{k}} \\ &\geq \frac{k^2}{2n} \cdot \left(1 - \frac{k-1}{n-k+1}\right)^{k-1} \binom{2k-1}{k-1}^{-2} (A(k-1, 1) + B(k)) - \frac{D}{\binom{n}{k}}. \end{aligned} \quad (\text{S30})$$

For $k \rightarrow \infty$ and $k = o(n^{1/2})$, it is easy to check $D/\binom{n}{k} = O(k^2/n)$. Applying Fact 6.4,

$$\begin{aligned} \frac{\sigma_{2,n}^2}{\mu_{2,n}^2} &\geq C \cdot \frac{k^{5/2}}{2n} \left(1 - \frac{k-1}{n-k+1}\right)^{k-1} \left[\frac{(2k-2) \cdots k}{(2k-1) \cdots (k+1)} \right]^2 + O(k^2/n) \\ &= C \cdot \frac{k^{5/2}}{2n} (1 + O(k^2/n)) \left(\frac{k}{2k-1}\right)^2 + O(k^2/n) \\ &\geq C \cdot \frac{k^{5/2}}{n} \end{aligned} \tag{S31}$$

for some $C > 0$ and sufficiently large k and n . Applying the bound from the Stein method as in equation (S17), we have

$$\begin{aligned} d_W(T_2, Y) &\leq C_1 \cdot \frac{D^2 \mu_{2,n}}{\sigma_{2,n}^3} + C_2 \cdot \frac{D^{3/2} \mu_{2,n}^{1/2}}{\sigma_{2,n}^2} \\ &\leq C_1 \cdot \frac{k^{1/4} (k!)^2}{n^{1/2}} + C_2 \cdot \frac{k^{1/2} (k!)^{3/2}}{n^{1/2}} \rightarrow 0 \end{aligned} \tag{S32}$$

for $k/(\log n)^\alpha \rightarrow 0$.

For k fixed, $D/\binom{n}{k} \leq k^2/(n-k+1) + o(1/n)$. (S30) becomes

$$\begin{aligned} \frac{\sigma_{2,n}^2}{\mu_{2,n}^2} &\geq \frac{k^2}{2n} (1 + O(1/n)) \binom{2k-1}{k-1}^{-2} (A(k-1, 1) + B(k)) - \frac{k^2}{n-k+1} + o(1/n) \\ &= \left\{ \frac{1}{2} (A(k-1, 1) + B(k)) \binom{2k-1}{k-1}^{-2} - 1 \right\} \frac{k^2}{n} + o(1/n) \\ &:= C(k) \cdot \frac{k^2}{n} + o(1/n), \quad \text{say.} \end{aligned} \tag{S33}$$

When $k = 3$, we can check that $C(3) > 0$ and thus $\sigma_{2,n}^2/\mu_{2,n}^2 \geq C/n$. For other fixed k , the same order lower bound holds. Applying (S17),

$$d_W(T_2, Y) \leq O(n^{-1/2}) \rightarrow 0. \tag{S34}$$

□

6.4 Power analysis

First we prove a lemma upper bounding the variances of T_1 and T_2 .

Lemma 6.5. • $\sigma_{1,n}^2 = O(n)$ for fixed k ; $\sigma_{1,n}^2 = O(n/k!)$ for $k \rightarrow \infty$ and $k/(\log n)^\alpha \rightarrow 0$.

• $\sigma_{2,n}^2 = O(n^{2k-1})$ for fixed k ; $\sigma_{2,n}^2 = O(\mu_{2,n}^2 k^{5/2}/n)$ for $k \rightarrow \infty$ and $k/(\log n)^\alpha \rightarrow 0$.

Proof. By the upper bound in (S8),

$$\begin{aligned} \sigma_{1,n}^2 &= \sum_{i=1}^{n-k+1} \sum_{j \in N_i \setminus \{i\}} \mathbb{E}(\mathbb{I}_i \mathbb{I}_j) + \frac{2(n-k+1)}{k!} - \sum_{i=1}^{n-k+1} \sum_{j \in N_i} \mathbb{E}(\mathbb{I}_i) \mathbb{E}(\mathbb{I}_j) \\ &\leq \frac{C(n-k+1)}{(k+1)!} + \frac{2(n-k+1)}{k!} \\ &= \begin{cases} O(n) & \text{for fixed } k; \\ O(n/k!) & \text{for } k \rightarrow \infty, k/(\log n)^\alpha \rightarrow 0. \end{cases} \end{aligned} \tag{S35}$$

To bound $\sigma_{2,n}^2$, first note that $B(k) \leq A(k-1, 1)$ for all $k \geq 2$. This holds because for every pair of (l_0, l_1) and (m_0, m_1) such that $l_0 + l_1 = k-1$ and $m_0 + m_1 = k-1$, we have

$$\binom{l_0 + m_0}{l_0} \binom{l_1 + m_1}{l_1} + \binom{l_0 + m_1}{l_0} \binom{l_1 + m_0}{l_1} \geq 2 \binom{l_0 + m_0}{l_0} \binom{l_1 + m_1}{l_1} \binom{l_0 + m_1}{l_0} \binom{l_1 + m_0}{l_1}.$$

By equations (S24), (S25) and (S28),

$$\begin{aligned} \sigma_{2,n}^2 &= 2 \sum_{j=1}^k \binom{n}{2k-j} \frac{1}{(2k-j)!} A(k-j, j) + 2 \binom{n}{2k-1} \frac{1}{(2k-1)!} B(k) - \frac{4D \binom{n}{k}}{(k!)^2} \\ &\leq 4 \sum_{j=1}^k \binom{n}{2k-j} \frac{1}{(2k-j)!} A(k-j, j) - \frac{4D \binom{n}{k}}{(k!)^2} \\ &= O\left(\frac{\mu_{2,n}^2 k^{5/2}}{n}\right) \end{aligned} \tag{S36}$$

by Theorem 1 in [14]. The first part of the lemma holds since $\mu_{2,n} = O(n^k)$ for k fixed. \square

Proof of Theorem 3. It is easy to see the count W_1 is bounded below by $m - k + 1$. By the first part of Lemma 6.5,

$$\begin{aligned} T_1 &\geq \frac{m - k + 1 - \mu_{1,n}}{\sigma_{1,n}} \\ &\geq C\sqrt{n} \left(\frac{m}{n} - \frac{2}{k!} \right), \end{aligned} \tag{S37}$$

for some $C > 0$, fixed k and m , n sufficiently large. In this case, m has to grow at the same rate as n , that is $m \sim a_1 n$ and $a_1 > 2/k!$. It follows then $T_1 = \Omega(\sqrt{n})$.

When $k \rightarrow \infty$ and $k/(\log n)^\alpha \rightarrow 0$, for n large enough,

$$\begin{aligned} T_1 &\geq C \sqrt{\frac{n}{k!}} \left(\frac{k!(m-k+1)}{n-k+1} - 2 \right) \\ &\geq C \sqrt{\frac{n}{k!}} \left(\frac{a_2 n - k! \cdot k + k!}{n-k+1} - 2 \right) \\ &= \Omega\left(\sqrt{\frac{n}{k!}}\right) \end{aligned} \tag{S38}$$

for $m \geq a_2 n/k!$, $a_2 > 2$. If m grows at the rate of $a_3 n$, $a_3 \in (0, 1]$,

$$\begin{aligned} T_1 &\geq C\sqrt{nk!} \left(\frac{m-k+1}{n-k+1} - \frac{2}{k!} \right) \\ &\geq C\sqrt{nk!} a_3 = \Omega(\sqrt{nk!}). \end{aligned} \tag{S39}$$

Similarly, the count W_2 is lower bounded by $\binom{m}{k}$, using the second part of Lemma 6.5, for fixed k ,

$$\begin{aligned} T_2 &\geq C \cdot \frac{\binom{m}{k} - 2\binom{n}{k}/k!}{n^{k-1/2}} \\ &= C\sqrt{n} \left(\frac{m \cdots (m-k+1)}{n \cdots (n-k+1)} - \frac{2}{k!} \right) \\ &\geq C\sqrt{n} \left(\left(\frac{m}{n}\right)^k - \frac{2}{k!} \right) \end{aligned} \tag{S40}$$

for sufficiently large m and n . m has to grow at the rate of $b_1 n$ for the lower bound to go to infinity, and $b_1^k > 2/k!$. We have $T_2 = \Omega(\sqrt{n})$.

When $k \rightarrow \infty$ and $k/(\log n)^\alpha \rightarrow 0$, again by Lemma 6.5,

$$\begin{aligned} T_2 &\geq C \sqrt{\frac{n}{k^{5/2}}} \left(k! \frac{m \cdots (m-k+1)}{n \cdots (n-k+1)} - 2 \right) \\ &\geq C \sqrt{\frac{n}{k^{5/2}}} \left(k! \left(\frac{m}{n} \right)^k - 2 \right) \\ &\geq C \sqrt{\frac{n}{k^{3/2}}}, \end{aligned} \tag{S41}$$

for $m \geq en/k$. When $m \sim b_2 n$, $b_2 \in (0, 1]$,

$$\begin{aligned} T_2 &\geq C \sqrt{\frac{n}{k^{5/2}}} \left(k! \left(\frac{m}{n} \right)^k - 2 \right) \\ &\geq C b_2^k k! \sqrt{\frac{n}{k^{5/2}}}. \end{aligned} \tag{S42}$$

□

Proof of Theorem 4. Let n_1, \dots, n_d denote the number of points in (\mathbf{x}, \mathbf{y}) falling on to each monotonic piece, then W_2 is lower bounded by $\sum_{t=1}^d \binom{n_t}{k}$. For fixed d and k ,

$$\frac{\sum_{t=1}^d n_t \cdots (n_t - k + 1)}{n \cdots (n - k + 1)} \xrightarrow{P} \sum_{t=1}^d \ell_t^k \tag{S43}$$

Since by Lemma 6.5,

$$T_2 \geq C \sqrt{n} \left(\frac{\sum_{t=1}^d n_t \cdots (n_t - k + 1)}{n \cdots (n - k + 1)} - \frac{2}{k!} \right) \tag{S44}$$

for some $C > 0$, it follows

$$\mathbb{P} \left(T_2 \geq C \sqrt{n} (d^{-(k-1)} - 2/k!) \right) \rightarrow 1 \tag{S45}$$

using Hölder's inequality and the fact $\sum_{t=1}^d \ell_t = 1$. Thus T_2 is lower bounded by $C \sqrt{n}$ with probability tending to 1 when $d^{k-1} < k!/2$.

When $k \rightarrow \infty$ and $k/(\log n)^\alpha \rightarrow 0$, it is easy to check

$$\frac{n_t \cdots (n_t - k + 1)}{n \cdots (n - k + 1)} \cdot \left(\frac{n}{n_t} \right)^k \xrightarrow{P} 1. \tag{S46}$$

Also,

$$\begin{aligned} &\mathbb{P} \left(\left| \left(\frac{n_t}{n \ell_t} \right)^k - 1 \right| \geq \epsilon \right) \\ &\leq \mathbb{P} \left(\frac{n_t}{n \ell_t} - 1 \geq (1 + \epsilon)^{1/k} - 1 \right) + \mathbb{P} \left(\frac{n_t}{n \ell_t} - 1 \leq (1 - \epsilon)^{1/k} - 1 \right) \\ &\leq \exp(-2n \ell_t^2 ((1 + \epsilon)^{1/k} - 1)^2) + \exp(-2n \ell_t^2 ((1 - \epsilon)^{1/k} - 1)^2) \rightarrow 0 \end{aligned} \tag{S47}$$

by Hoeffding’s inequality. It follows then

$$\sum_{t=1}^d \frac{n_t \cdots (n_t - k + 1)}{n \cdots (n - k + 1)} \cdot \left(\sum_{t=1}^d \ell_t^k \right)^{-1} \xrightarrow{P} 1. \quad (\text{S48})$$

Now noting that

$$T_2 \geq C \sqrt{\frac{n}{k^{5/2}}} \left(k! \frac{\sum_{t=1}^d n_t \cdots (n_t - k + 1)}{n \cdots (n - k + 1)} - 2 \right), \quad (\text{S49})$$

we have

$$\mathbb{P} \left(T_2 \geq C \frac{k!}{d^{k-1}} \sqrt{\frac{n}{k^{5/2}}} \right) \rightarrow 1 \quad (\text{S50})$$

again by Hölder’s inequality. □

References

- [1] Frank E Harrell. Hmisc: Harrell miscellaneous. <http://cran.r-project.org/web/packages/Hmisc/index.html>, 2014.
- [2] Maria L Rizzo and Gabor J Szekely. energy: E-statistics (energy statistics). <http://cran.r-project.org/web/packages/energy/index.html>, 2014.
- [3] Phil Spector, Jerome Friedman, Robert Tibshirani, and Thomas Lumley. acepack: ace() and avas() for selecting regression transformations. <http://cran.r-project.org/web/packages/acepack/index.html>, 2013.
- [4] Gabriele Sales and Chiara Romualdi. parmigene – a parallel R package for mutual information estimation and gene network reconstruction. *Bioinformatics*, 27(13):1876–1877, 2011.
- [5] Davide Albanese, Michele Filosi, Roberto Visintainer, Samantha Riccadonna, Giuseppe Jurman, and Cesare Furlanello. minerva and minepy: a C engine for the MINE suite and its R, Python and MATLAB wrappers. *Bioinformatics*, 29(3):407–408, 2013.
- [6] Leo Breiman and Jerome H Friedman. Estimating optimal transformations for multiple regression and correlation. *J Am Stat Assoc*, 80(391):580–598, 1985.
- [7] Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. Estimating mutual information. *Phys Rev E*, 69(6):066138, 2004.
- [8] Jesús E García and V. A. González-López. Independence tests for continuous random variables based on the longest increasing subsequence. *Journal of Multivariate Analysis*, 127:126–146, 2014.
- [9] Xuewen Pan, Ping Ye, Daniel S Yuan, Xiaoling Wang, Joel S Bader, and Jef D Boeke. A DNA integrity network in the yeast *Saccharomyces cerevisiae*. *Cell*, 124(5):1069–1081, 2006.
- [10] Michael Costanzo, Anastasia Baryshnikova, Jeremy Bellay, Yungil Kim, Eric D Spear, Carolyn S Sevier, Huiming Ding, Judice LY Koh, Kiana Toufighi, and Sara Mostafavi. The genetic landscape of a cell. *Science*, 327(5964):425–431, 2010.

- [11] Rafael A Irizarry, Bridget Hobbs, Francois Collin, Yasmin D Beazer-Barclay, Kristen J Antonellis, Uwe Scherf, and Terence P Speed. Exploration, normalization, and summaries of high density oligonucleotide array probe level data. *Biostatistics*, 4(2):249–264, 2003.
- [12] Peter M Fenwick. A new data structure for cumulative frequency tables. *Software: Practice and Experience*, 24(3):327–336, 1994.
- [13] Nathan Ross et al. Fundamentals of Stein’s method. *Prob. Surv*, 8:210–293, 2011.
- [14] Ross G Pinsky. Law of large numbers for increasing subsequences of random permutations. *Random Struct Algorithms*, 29(3):277–295, 2006.