

Consistent Testing for Recurrent Genomic Aberrations

Vonn Walter

Lineberger Comprehensive Cancer Center,
University of North Carolina at Chapel Hill, Chapel Hill, NC 27579

Fred A. Wright

Departments of Statistics and Biological Sciences
and the Bioinformatics Research Center,
North Carolina State University, Raleigh, NC 27695

Andrew B. Nobel

Department of Statistics and Operations Research,
Department of Biostatistics, and Lineberger Comprehensive Cancer Center
University of North Carolina at Chapel Hill, Chapel Hill, NC 27579

October 1, 2018

Abstract

Genomic aberrations, such as somatic copy number alterations, are frequently observed in tumor tissue. Recurrent aberrations, occurring in the same region across multiple subjects, are of interest because they may highlight genes associated with tumor development or progression. A number of tools have been proposed to assess the statistical significance of recurrent DNA copy number aberrations, but their statistical properties have not been carefully studied. Cyclic shift testing, a permutation procedure using independent random shifts of genomic marker observations on the genome, has been proposed to identify recurrent aberrations, and is potentially useful for a wider variety of purposes, including identifying regions with methylation aberrations or overrepresented in disease association studies. For data following a countable-state Markov model, we prove the asymptotic validity of cyclic shift p -values under a fixed sample size regime as the number of observed markers tends to infinity. We illustrate cyclic shift testing for a variety of data types, producing biologically relevant findings for three publicly available datasets.

1 Introduction

Many genomic datasets consist of measurements from multiple samples at a common set of genetic markers, with no “phenotype” representing clinical state or experimental condition of the sample. Datasets of this type include genome-wide measurements of DNA copy number or DNA methylation, for which the main goal is to identify aberrant regions on the genome that tend to have extreme measurements in comparison to other regions. Testing for aberrations requires some thought about appropriate test statistics, and constructing a null distribution that appropriately reflects serial correlation structures inherent to genomic data. A meta-analysis across several genome-wide association studies might also be viewed in this framework, in the sense that the testing for association within each study produces a vector of p -values that might be viewed as a vector of “observations.”

The problem of interest to us is the identification of aberrant markers, where multiple samples exhibit a coordinated (unidirectional), departure from the expected state. Aberrant markers

are of particular interest in cancer studies, where tumor suppressors or oncogenes exhibit DNA copy variation or modified methylation levels. Similarly, it may be possible to identify pleiotropic single nucleotide polymorphisms (SNPs) in disease association by identifying genetic markers that repeatedly give rise to small p -values in multiple association studies.

In this paper we provide a rigorous asymptotic analysis of a permutation based testing procedure for identifying aberrant markers in genomic data sets. The procedure, called DiNAMIC, was introduced in Walter et al. (2011), and is described in detail below. In contrast to other procedures which permute all observations, DiNAMIC is based on cyclic shifting of samples. Cyclic shifting eliminates concurrent findings across samples, but retains the adjacency of observations in a sample (with the exception of the first and last entries), thereby largely preserving the correlation structure among markers. Our principal result is that, for a broad family of null data distributions, the sampling distribution of the DiNAMIC procedure is close to the true conditional distribution of the data restricted to its cyclic shifts. As a corollary, we find that the cyclic shift testing provides asymptotically correct Type I error rates.

The outline of the paper is as follows. The next section is devoted to a description of the cyclic shift procedure, a discussion of the underlying testing framework within which our analysis is carried out, and a statement of our principal result. In Section 3 we apply cyclic shift testing to DNA copy number analysis, DNA methylation analysis, and meta-analysis of GWAS data, and show that the results are consistent with the existing biological literature. Because of its broad applicability and solid statistical foundation, we believe that cyclic shift testing is a valuable tool for the identification of aberrant markers in many large scale genomic studies.

2 Asymptotic Consistency of Cyclic Shift Permutation

2.1 Data Matrix

We consider a data set derived from n subjects at m common genomic locations or markers. The data is arranged in an $n \times m$ matrix \mathbf{X} with values in a set $\mathcal{A} \subseteq \mathbb{R}$. Depending on the application, \mathcal{A} may be finite or infinite. The entry x_{ij} of \mathbf{X} contains data from subject i at marker j . Thus the i th row \mathbf{X}_i of \mathbf{X} contains the data from subject i at all markers, and the j th column $\mathbf{X}_{.j}$ of \mathbf{X} contains the data at marker j across subjects. For $1 \leq j \leq m$ let $s_j = s_j(\mathbf{X}_{.j})$ be a local summary statistic for the j th marker. In most applications the simple sum statistic $s_j = \sum_{i=1}^n x_{ij}$ is employed. In order to identify locations with coordinated departures from baseline behavior, we apply a global summary statistic to the local statistics s_1, \dots, s_m . When looking for extreme, positive departures from baseline it is natural to employ the global statistic

$$T(\mathbf{X}) = \max(s_1, \dots, s_m). \quad (2.1)$$

To detect negative departures from baseline, the maximum may be replaced by a minimum. The cyclic shift procedure and the supporting theory in Theorem 1 apply to arbitrary local statistics, as well as a range of global statistics.

2.2 Cyclic Shift Testing

Given a data matrix \mathbf{X} , we are interested in assessing the significance of the observed value $t_0 = T(\mathbf{X})$ of the global statistic. When t_0 is found to be significant, the identity and location of the marker j having the maximum (or minimum) local statistic is of primary biological importance. While in special cases it is possible to compute p -values for t_0 under parametric assumptions, permutation based approaches are often an attractive and more flexible alternative. A permutation based p -value

can be obtained by applying permutations π to the entries of \mathbf{X} , producing the matrices $\pi(\mathbf{X})$, and then comparing t_0 to the resulting values $T(\pi(\mathbf{X}))$ of the global statistic. The maximum global statistic accounts for multiple comparisons across markers, so it is not necessary to apply further multiplicity correction to the permuted values $T(\pi(\mathbf{X}))$.

The performance and suitability of permutation based p -values in the marker identification problem depends critically on the family of allowable permutations π . If π permutes the entries of \mathbf{X} without preserving row or column membership, then the induced null distribution is equivalent to sampling the entries of \mathbf{X} at random without replacement. In this case the induced null distribution does not capture the correlation of measurements within a sample, or systematic differences (e.g. in scale, location, correlation) between samples. In real data, correlations within and systematic differences between samples can be present even in the absence of aberrant markers. As such, p -values obtained under full permutation of \mathbf{X} will be sensitive to secondary features of the data and may yield significant p -values even when no aberrant markers are present. An obvious improvement of full permutation is to separately permute the values in each row (sample) of the data matrix. This approach is used in the GISTIC procedure of Beroukhim et al. (2007). While row-by-row permutation preserves some differences between rows, it eliminates correlations within rows (and correlation differences between rows), so that the induced null distribution is again sensitive to secondary, correlation based features of the data that are not related to the presence of aberrant markers.

The DiNAMIC cyclic shift testing procedure of Walter et al. (2011) addresses the shortcomings of full and row-by-row permutation by further restricting the set of allowable permutations. In the procedure, each row of the data matrix is shifted to the left in a cyclic fashion, as detailed below, so that the first k entries of the vector are placed after the last element; the values of the offsets k are chosen independently from row to row. Cyclic shifting preserves the serial correlation structure with each sample, except at the single break point where the last and first elements of

the unshifted sample are placed next to one another. At the same time, the use of different offsets breaks concurrency among the samples, so that the resulting cyclic null distribution is appropriate for testing the significance of $t_0 = T(\mathbf{X})$.

2.3 Cyclic Shift Testing

Formally, a *cyclic shift* of index $k \in \{0, \dots, m-1\}$ is a map $\sigma_k : \mathcal{A}^m \rightarrow \mathcal{A}^m$ whose action is defined as follows:

$$\sigma_k(x_1, x_2, \dots, x_m) = (x_{k+1}, x_{k+2}, \dots, x_m, x_1, \dots, x_k).$$

Given $\mathbf{k} = (k_1, \dots, k_n)$ with $k_i \in \{0, \dots, m-1\}$, let $\sigma_{\mathbf{k}} = \sigma_{k_1} \otimes \dots \otimes \sigma_{k_n}$ be the map from the set $\mathcal{A}^{n \times m}$ of data matrices to itself defined by applying σ_{k_i} to the i th row of \mathbf{X} , namely,

$$\sigma_{\mathbf{k}}(\mathbf{X}) = (\sigma_{k_1}(\mathbf{X}_{1\cdot}), \dots, \sigma_{k_n}(\mathbf{X}_{n\cdot}))^t$$

The cyclic shift testing procedure of Walter et al. (2011) is as follows.

Cyclic shift procedure to assess the statistical significance of $T(\mathbf{X})$

1. Let $\sigma^1(\cdot), \dots, \sigma^N(\cdot)$ be random cyclic shifts of the form $\sigma_{k_1} \otimes \dots \otimes \sigma_{k_n}$, where k_1, \dots, k_n are independent and each is chosen uniformly from $\{0, \dots, m-1\}$.
2. Compute the values $T(\sigma^1(\mathbf{X})), \dots, T(\sigma^N(\mathbf{X}))$ of the global statistic T at the random cyclic shifts of \mathbf{X} .
3. Define the percentile-based p -value

$$p(T(\mathbf{X})) = \max \left(N^{-1} \sum_{l=1}^N I(T(\sigma^l(\mathbf{X})) \geq T(\mathbf{X})), 1/N \right).$$

Here $I(A)$ is the indicator function of the event A .

2.4 Testing Framework

We wish to assess the performance of the cyclic shift procedure within a formal testing framework. To this end, we regard the observed data matrix \mathbf{X} as an observation from a probability distribution P_m on $\mathcal{A}^{n \times m}$, so that for any (measurable) set $A \subseteq \mathcal{A}^{n \times m}$ the probability $P_m(A) = \mathbb{P}(\mathbf{X} \in A)$. As measurements derived from distinct samples are typically independent, we restrict our attention to the family of measures \mathcal{P} on $\mathcal{A}^{n \times m}$ under which the rows of \mathbf{X} are independent.

Let $\mathcal{P}_0 \subseteq \mathcal{P}$ be the sub-family of \mathcal{P} corresponding to the null hypothesis that \mathbf{X} has no atypical markers, *i.e.*, no markers exhibiting coordinated activity across samples. One may define \mathcal{P}_0 in a variety of ways, but the simplest is to let \mathcal{P}_0 be the set of distributions $P_m \in \mathcal{P}$ such that the rows of \mathbf{X} are stationary and ergodic under P_m ; independence of the rows follows from the definition of \mathcal{P} . Under \mathcal{P}_0 the columns of \mathbf{X} are stationary and ergodic, and the same is true of the local statistics s_j , which are identically distributed and have constant mean and variance. Thus under \mathcal{P}_0 no marker is atypical in a strong distributional sense.

Our principal result shows that the p -value produced by the cyclic shift procedure is approximately consistent for distributions P_m in a subfamily $\mathcal{P}^* \subseteq \mathcal{P}_0$. The family \mathcal{P}^* includes or approximates many distributions of practical interest, including finite order Markov chains with discrete or continuous state spaces. In order to assess the consistency of the cyclic shift p -value we carefully define both the target and the induced distributions of the procedure. As much of what follows concerns probabilities conditional on the observed data matrix, we use \mathbf{X} to denote both the random matrix and its observed realization. Given \mathbf{X} let

$$\mathcal{S}_m(\mathbf{X}) = \{\sigma_{\mathbf{k}}(\mathbf{X}) : \mathbf{k} \in \{0, \dots, m-1\}^n\} \subseteq \mathcal{A}^{n \times m}$$

be the set of all cyclic shifts of \mathbf{X} . Define the *true conditional distribution* to be the conditional

distribution of P_m given $\mathcal{S}_m(\mathbf{X})$, namely

$$P_{\mathbf{X}}(A) = P_m(A | \mathcal{S}_m(\mathbf{X})) \quad A \subseteq \mathcal{A}^{n \times m}.$$

If P_m is discrete with probability mass function $p(\cdot)$ then

$$P_{\mathbf{X}}(A) = \frac{1}{\sum_{\mathbf{Y}' \in \mathcal{S}_m(\mathbf{X})} p(\mathbf{Y}')} \sum_{\mathbf{Y} \in A} p(\mathbf{Y}) \cdot I(\mathbf{Y} \in \mathcal{S}_m(\mathbf{X})).$$

If P_m has probability density function $f(\cdot)$ then $P_{\mathbf{X}}$ may be defined in a similar fashion.

In the cyclic shift procedure, matrices are selected uniformly at random from the set $\mathcal{S}_m(\mathbf{X})$ of cyclic shifts of the observed data matrix \mathbf{X} . The associated *cyclic conditional distribution* has the form

$$Q_{\mathbf{X}}(A) = \sum_{\mathbf{Y} \in A} \frac{1}{|\mathcal{S}_m(\mathbf{X})|} \cdot I(\mathbf{Y} \in A) \quad A \subseteq \mathcal{A}^{n \times m}.$$

Under mild conditions the m^n cyclic shifts of \mathbf{X} are distinct with high probability when m is large (see Lemma 3 in Section 4). In this case, the cyclic conditional distribution may be written as

$$Q_{\mathbf{X}}(A) = \sum_{\mathbf{Y} \in A} \frac{1}{m^n} \cdot I(\mathbf{Y} \in \mathcal{S}_m(\mathbf{X})) = \frac{1}{m^n} \cdot |A \cap \mathcal{S}_m(\mathbf{X})| \quad A \subseteq \mathcal{A}^{n \times m}.$$

The distribution of the cyclic shift p -value is given by

$$p(T(\mathbf{X})) \sim \max(N^{-1} \text{Bin}(N, \alpha), 1/N).$$

Here $\alpha = Q_{\mathbf{X}}(T \geq t_0)$ where t_0 is the observed value of $T(\mathbf{X})$, and $T \geq t_0$ represents the event $\{\mathbf{Y} : T(\mathbf{Y}) \geq t_0\}$. Note that as the number N of cyclic shifts increases, the p -value $p(T(\mathbf{X}))$ will converge in probability to $Q_{\mathbf{X}}(T \geq t_0)$

2.5 Principal Result

Our principal result requires an invariance condition on the global statistic T . Informally, the condition ensures that T does not give special treatment to any column of the data matrix.

Definition: A statistic $T : \mathcal{A}^{n \times m} \rightarrow \mathbb{R}$ is *invariant under constant shifts* if $T(\mathbf{X}) = T(\mathbf{X}')$ whenever \mathbf{X}' is obtained from \mathbf{X} by applying the *same* cyclic shift $\sigma_k(\cdot)$ to each row of \mathbf{X} .

The maximum column sum statistic used in the cyclic shift testing procedure is clearly invariant under constant shifts. More generally, any statistic of the form $T(\mathbf{X}) = g(h(\mathbf{X}_{\cdot 1}), \dots, h(\mathbf{X}_{\cdot m}))$ where $h : \mathcal{A}^n \rightarrow \mathcal{B}$ is an arbitrary local statistic (not necessarily a sum), and $g : \mathcal{B}^m \rightarrow \mathbb{R}$ is invariant under cyclic shifts will be invariant under constant shifts. The following result establishes the asymptotic validity of the cyclic shift procedure in this general setting.

Theorem 1. *Let \mathbf{X} be a random $n \times m$ matrix whose rows \mathbf{X}_i are independent copies of a first-order stationary ergodic Markov chain with countable state space \mathcal{A} and transition probabilities $p(u|v)$.*

Suppose that

$$\max_{u,v \in \mathcal{A}} p(u|v) < 1 \quad \text{and} \quad \frac{p_1(u)p_1(v)}{p_2(u,v)} < \infty \quad \text{for each } u, v \in \mathcal{A} \quad (2.2)$$

where in the second condition we define $0/0$ to be 0. Here $p_1(\cdot)$ and $p_2(\cdot, \cdot)$ denote the one- and two-dimensional marginal distributions of the Markov chain, respectively. For $m \geq 1$ let $T_m : \mathcal{A}^{n \times m} \rightarrow \mathbb{R}$ be a statistic that is invariant under constant shifts. Then

$$\max_{B \subseteq \mathbb{R}} |P_{\mathbf{X}}(T_m \in B) - Q_{\mathbf{X}}(T_m \in B)|$$

tends to zero in probability as m tends to infinity.

The first condition in (2.2) ensures that there are not deterministic transitions between the states of the Markov chain. The second condition can be expressed equivalently as $p_2(u, v) = 0$ implies $p_1(u)p_1(v) = 0$. The proof of Theorem 1 is given in Section 4. As an immediate corollary of

the theorem, we find that

$$\sup_t |P_{\mathbf{X}}(T_m \geq t) - Q_{\mathbf{X}}(T_m \geq t)|$$

tends to zero in P_m -probability as m tends to infinity. Thus, under the conditions of the theorem, when m and N are large, the percentile based p -value $p(T(\mathbf{X}))$ will be close to the true conditional probability $P_{\mathbf{X}}(T_m \geq t_0)$ that $T_m(\mathbf{X})$ exceeds the observed value of T_m . If we define $Q_m(A)$ to be $\mathbb{E}Q_{\mathbf{X}}(A)$, where the expectation is taken under P_m , then conditional convergence also yields the unconditional result

$$\sup_t |P_m(T_m \geq t) - Q_m(T_m \geq t)| \rightarrow 0$$

as m tends to infinity. Thus, under the assumptions of Theorem 1, the percentile based p -value provides asymptotically correct type I error rates.

Theorem 1 can be extended in a number of directions. Under conditions similar to those in (2.2) the theorem extends to matrices \mathbf{X} whose rows are independent copies of a k th order ergodic Markov chain, where $k \geq 2$ is fixed and finite. The theorem can also be extended to settings in which the rows of \mathbf{X} are independent stationary ergodic Markov chains with *different* transition probabilities. In this case we require that the conditions (2.2) hold for each row-chain.

Theorem 1 can also be extended to the setting in which the rows of \mathbf{X} are independent copies of a first-order stationary ergodic Markov chain with a continuous state space and a transition probability density $f(u|v)$. The existence of the transition probability density obviates the need for the first condition in (2.2) and the analysis of Lemmas 2 and 3 in Section 4. The second condition of (2.2) is replaced by the assumption

$$\frac{f_1(X_1) f_1(X_m)}{f_2(X_1, X_m)} = O_P(1), \tag{2.3}$$

where $f_1(\cdot)$ and $f_2(\cdot, \cdot)$ denote the one- and two-dimensional marginal densities of the Markov chain, respectively. Markovity and ergodicity ensure that (X_1, X_m) converges weakly to a pair

(X, X') consisting of independent copies of X_1 , and therefore condition (2.3) holds if the ratio $f_1(u) f_1(v) / f_2(u, v)$ is continuous on \mathbb{R}^2 . Thus Theorem 1 applies, for example, to standard Gaussian AR(1) models. As in the discrete case, one may extend the theorem to settings in which the rows of \mathbf{X} are independent stationary ergodic Markov chains with *different* transition probabilities, provided that (2.3) holds for each row-chain.

2.6 Illustration of Resampling Distributions

Here we present simulation results illustrating the resampling distributions $P_{\mathbf{X}}(A)$ and $Q_{\mathbf{X}}(A)$ defined above. Each simulation was conducted using an $n \times m$ matrix \mathbf{X} with independent, identically distributed rows generated by a stationary first-order r -state Markov chain with a fixed transition matrix M . Figure 1 shows empirical cumulative distribution functions (cdfs) $P_{\mathbf{X}}(T < t)$ and $Q_{\mathbf{X}}(T < t)$ based on simulations conducted with $r = 5$, $n = 4$, and $m = 10$ or 50. Each panel is based on an observed matrix \mathbf{X} produced by the Markov chain, and the results presented here are representative of those obtained from other simulations. Based on Theorem 1, we expect the cdfs to converge as the number of columns m increases. Accordingly, the two curves in each panel of part B of Figure 1 ($m = 50$) exhibit a greater level of concordance than those in part A ($m = 10$). Additional simulation results based on an AR(1) model are presented in Section 7, the Appendix.

3 Application to Genomic Data

In tumor studies DNA copy number values for each subject are measured with respect to a normal reference, typically either a paired normal sample or a pooled reference. In the autosomes the normal DNA copy number is two. Underlying genomic instability in tumor tissue can result in DNA copy number gains and losses, and often these changes lead to increased or decreased expression, respectively, of affected genes (Pinkel and Albertson 2005). Some of these genetic aberrations occur

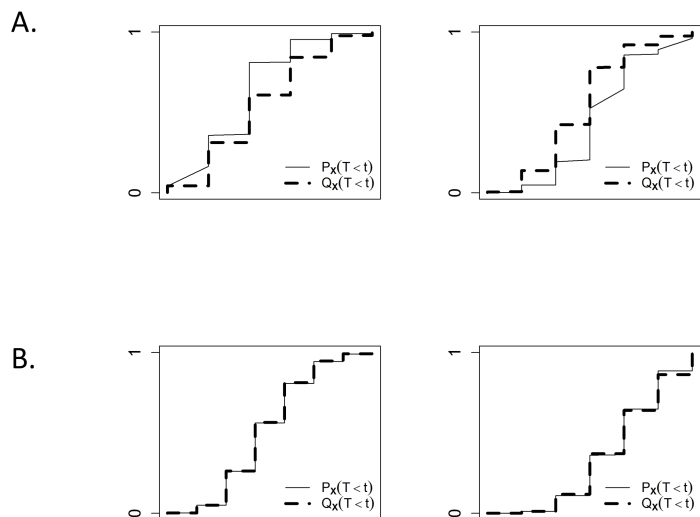


Figure 1: Illustration of Resampling Distributions. Empirical cumulative distribution functions for simulated matrices \mathbf{X} in which independent rows \mathbf{X}_i are generated by a first-order finite-state Markov chain. Each panel corresponds to a simulated $4 \times m$ matrix \mathbf{X} with $m = 10$ (A) or $m = 50$ (B).

at random locations throughout the genome, and these are termed *sporadic*. In contrast, *recurrent* aberrations are found in the same genomic region in multiple subjects. It is believed that recurrent aberrations arise because they lead to changes in gene expression that provide a selective growth advantage. Therefore regions containing recurrent aberrations are of interest because they may harbor genes associated with the tumor phenotype. Distinguishing sporadic and recurrent aberrations is largely a statistical issue, and the cyclic shift procedure was designed to perform this task.

DNA methylation values for a given subject are also measured with respect to a paired or pooled

normal reference. Although DNA methylation values are not constant across the genome, even in normal tissue, at a fixed location they are quite stable in normal samples from a given tissue type. Epigenetic instability can disrupt normal methylation patterns, leading to methylation gains and losses, and these changes can affect gene expression levels (Laird 2003). Regions of the genome that exhibit recurrent hyper- or hypo-methylation in tumor tissue are of interest.

3.1 Peeling

In many applications more than one atypical marker may be present, and as a result multiple columns may produce summary statistics with extreme values. In tumor tissue, for example, underlying genomic instability can result in gains and losses of multiple chromosomal regions; likewise, epigenetic instability can lead to aberrant patterns of DNA methylation throughout the genome. In order to identify multiple atypical markers and assess their statistical significance it is necessary to remove the effect of each discovered marker before initiating a search for the next marker. This task is carried out by a process known as *peeling*. Several peeling procedures have been proposed in the literature, including those employed by GISTIC (Beroukhim et al. 2007) and DiNAMIC (Walter et al. 2011). In the applications here we make use of the procedure described in detail in Walter et al. (2011) .

3.2 DNA Copy Number Data

Walter et al. (2011) used the cyclic shift procedure to analyze the Wilms' tumor data of Natrajan et al. (2006). Here we apply the procedure to the lung adenocarcinoma dataset of Chitale et al. (2009), with $n = 192$ and $m = 40478$. We detected a number of highly significant findings under the null hypothesis that no recurrent copy number gains or losses are present. Table 1 lists the genomic positions of the the three most significant copy number gains and losses, as well as neighboring genes, most of which are known oncogenes and tumor suppressors. Strikingly, Weir et al. (2007)

Table 1: Genomic locations of the three most significant DNA copy number gains (top table) and losses (bottom table) found by applying the cyclic shift procedure to the lung adenocarcinoma dataset of Chitale et al. (2009).

Chromosome	Gain Locus (bp)	Gene
5p15	967984	<i>TERT</i>
1q21	149346163	<i>ARNT</i>
8q24	128816933	<i>MYC</i>
Chromosome	Loss Locus (bp)	Gene
8p23	2795183	<i>CSMD1</i>
13q11	19254995	<i>PSPC1</i>
9p21	21958070	<i>CDKN2A</i>

detected highly significant gains of the oncogenes *TERT*, *ARNT*, and *MYC* in their comprehensive investigation of the disease, each of which appears in Table 1. The loss results for chromosomes 8 and 9 in Table 1 are also highly concordant with previous findings of Weir et al. (2007), and Wistuba et al. (1999). Weir et al. (2007) detected chromosomal loss in a broad region of 13q that contains the locus in Table 1, but it is not clear if the target of this region is the known tumor-suppressor *RB1* or some other gene.

3.3 DNA Methylation Data

Using unsupervised clustering techniques, Fackler et al. (2011) found an association between methylation patterns and estrogen-receptor status in a cohort of breast cancer tumors. This cohort consisted of 20 tumor/normal pairs, and we used differences in methylation signal between tumor and normal tissue as the observations. We applied the cyclic shift procedure to the resulting differences to detect loci that exhibited recurrent hyper- or hypomethylation in tumors. As shown in Table 2, the most significant hypermethylation sites occur in *ABCA3*, *GALR1*, and *NID2*, and these genes have previously been found to be highly methylated in lung adenocarcinoma, head and neck squamous cell carcinoma, and bladder cancer, respectively, by Selamat et al. (2012), Misawa et al. (2008),

Table 2: Genomic locations of the three most significant hypermethylation (top table) and hypomethylation (bottom table) sites found by applying the cyclic shift procedure to the breast cancer dataset of Fackler et al. (2011).

Chromosome	Gain Locus (bp)	Gene
16p13	2331829	<i>ABCA3</i>
18q23	73091357	<i>GALR1</i>
14q22	51605897	<i>NID2</i>
Chromosome	Gain Locus (bp)	Gene
20q13	62266251	<i>MYT1</i>
3q24	144378065	<i>SLC9A9</i>
1q21	150565702	<i>MCL1</i>

and Renard et al. (2009). Hypomethylation of the transcription factor *MYT1* on chromosome 20 was detected; this is notable because Viré et al. (2006) found that *MYT1* could be activated via decreased methylation.

3.4 Meta-Analysis of Genomewide Association Studies

Genome-wide association studies (GWAS) are used to identify genetic markers, typically *single nucleotide polymorphisms* (SNPs), that are associated with a disease of interest. When conducting a GWAS involving a common disease and alleles with small to moderate effect sizes, large numbers of cases and controls are required to have adequate power to detect disease SNPs (Pfeiffer et al. 2009).

The Wellcome Trust Case Control Consortium (WTCCC 2007) performed a genome-wide association study of seven common familial diseases - bipolar disorder (BD), coronary artery disease (CAD), Crohn’s disease (CD), hypertension (HT), rheumatoid arthritis (RA), type I diabetes (T1D), and type II diabetes (T2D) - based on an analysis of 2000 separate cases for each disease and a set of 3000 controls. We applied the inverse of the standard normal cumulative distribution function to the Cochran-Armitage trend test p -values from the WTCCC study, a transformation that produces z -scores whose values are similar those exhibited by a stationary process. We then analyzed the

matrix \mathbf{X} whose entries are negative thresholded z-scores arranged in rows corresponding to the seven disease phenotypes. As seen in Figure 2, a number of regional markers on chromosome 6 produce extremely large column sums. These markers lie in the major histocompatibility complex (MHC), which is noteworthy because MHC class II genes have been shown to be associated with autoimmune disorders, including RA and T1D (Fernando et al. 2008). When applied to \mathbf{X} , cyclic shift testing identified several highly significant apparently pleiotropic SNPs in the MHC region that produced large entries in the rows corresponding to both RA and T1D, including rs9270986, which is upstream of the RA and T1D susceptibility gene *HLA-DRB1*.

The WTCCC dataset serves as a proof of principle for cyclic shift applied to GWAS studies, although the use of a common set of controls may create modest additional correlation not fully captured in the cyclic shifts. We note that the cyclic shift procedure applied to GWAS is sensitive only to small p -values that occur in multiple studies. Thus the procedure is qualitatively different from typical meta-analyses, such as Zeggini et al. (2008), which can be sensitive to large observed effects from a single study.

4 Proof of Theorem 1

Let \mathbf{X} be a random $n \times m$ matrix whose rows $\mathbf{X}_1, \dots, \mathbf{X}_n$ are independent realizations of a first-order stationary ergodic Markov chain with countable state space \mathcal{A} . Denote the distribution of \mathbf{X} in $\mathcal{A}^{n \times m}$ by P_m . Let $p_1(\cdot)$ and $p(\cdot|\cdot)$ denote, respectively, the stationary distribution and the one-step transition probability of the Markov chain defining the rows of \mathbf{X} . Let $p_l(\cdot)$ denote the joint probability mass function of l contiguous variables in the chain. Thus the vectors \mathbf{X}_i have common probability mass function

$$p_m(x_0, \dots, x_{m-1}) = p_1(x_0) \prod_{j=1}^{m-1} p(x_j|x_{j-1}). \quad (4.1)$$

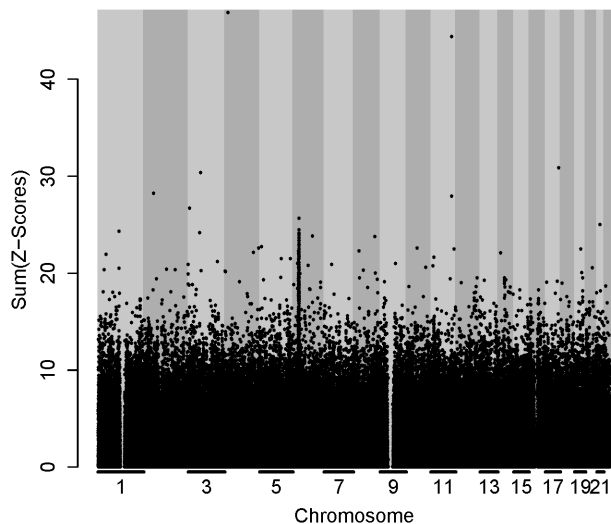


Figure 2: Cyclic Shift Testing Identifies Pleiotropic Single Nucleotide Polymorphisms. Marker-specific summary statistics were obtained from the Wellcome Trust Case Control Consortium study and plotted genome-wide. Numerous regional markers in the multihistocompatibility complex region of chromosome 6 exhibit large summary statistics, including several markers that were highly significant under cyclic shift testing and were associated with multiple disease phenotypes.

In what follows we assume that (2.2) holds. The ergodicity assumption on the Markov chain ensures that the joint probability mass function of (X_0, X_{m-1}) converges to the joint probability mass function of the pair (X, X') where $X, X' \in \mathcal{A}$ are independent with the same distribution as X_0 . It follows that

$$\frac{p_1(X_{i,0})p_1(X_{i,m-1})}{p_2(X_{i,0}, X_{i,m-1})} = O_P(1) \quad 1 \leq i \leq n. \quad (4.2)$$

In other words, for each row i the ratio in (4.2) is stochastically bounded under P_m as m tends to infinity.

Suppose for the moment that m is fixed. For any integer r , define the cyclic shift $\sigma_r : \mathcal{A}^m \rightarrow \mathcal{A}^m$ on sequences of length m by

$$\sigma_r(x_0, x_1, \dots, x_{m-1}) = (x_{[r]}, x_{[r+1]}, \dots, x_{[r+(m-1)]})$$

where $[k] = k \bmod m$. We index vectors as $(x_0, x_1, \dots, x_{m-1})$ rather than (x_1, x_2, \dots, x_m) , as was done in the body of the present manuscript, because this allows us to write the subscripts of the shifted vector in terms of $[k]$, substantially reducing notation. For each $\mathbf{r} = (r_1, \dots, r_n) \in \mathbb{Z}^n$ define $\sigma_{\mathbf{r}}(\mathbf{X})$ to be the $n \times m$ matrix with rows $\sigma_{r_1}(\mathbf{X}_{1\cdot}), \dots, \sigma_{r_n}(\mathbf{X}_{n\cdot})$. If $\mathbf{r}, \mathbf{s} \in \mathbb{Z}^n$, then it is easy to verify that

$$(\sigma_{\mathbf{r}} \circ \sigma_{\mathbf{s}})(\mathbf{X}) = (\sigma_{\mathbf{s}} \circ \sigma_{\mathbf{r}})(\mathbf{X}) = (\sigma_{\mathbf{r}+\mathbf{s}})(\mathbf{X}).$$

Let $\mathcal{S}_m(\mathbf{X}) = \{\sigma_{\mathbf{r}}(\mathbf{X}) : \mathbf{r} \in \mathbb{Z}^n\}$ be the set of cyclic shifts of \mathbf{X} .

Let $P_{\mathbf{X}}$ and $Q_{\mathbf{X}}$ be the true conditional and cyclic conditional distributions given $\mathcal{S}_m(\mathbf{X})$, defined by $P_{\mathbf{X}}(A) = P_m(A | \mathcal{S}_m(\mathbf{X}))$ and $Q_{\mathbf{X}}(A) = m^{-n} |A \cap \mathcal{S}_m(\mathbf{X})|$, respectively. In order to compare the distributions $P_{\mathbf{X}}$ and $Q_{\mathbf{X}}$ we introduce two closely related distributions, $P_{\mathbf{X}}^o$ and $Q_{\mathbf{X}}^o$, that are more amenable to analysis. Let $P_{\mathbf{X}}^o$ be a (random) measure on $\mathcal{A}^{n \times m}$ defined by

$$P_{\mathbf{X}}^o(A) = \sum_{\mathbf{r} \in [m]^n} \eta(\sigma_{\mathbf{r}}(\mathbf{X})) I(\sigma_{\mathbf{r}}(\mathbf{X}) \in A),$$

where $[m] = \{0, 1, \dots, m-1\}$ and

$$\eta(\sigma_{\mathbf{r}}(\mathbf{X})) = \prod_{i=1}^n \left[\frac{p_m(\sigma_{r_i}(\mathbf{X}_{i\cdot}))}{\sum_{s \in [m]} p_m(\sigma_s(\mathbf{X}_{i\cdot}))} \right].$$

Let $Q_{\mathbf{X}}^o$ be a (random) measure on $\mathcal{A}^{n \times m}$ defined by

$$Q_{\mathbf{X}}^o(A) = \sum_{\mathbf{r} \in [m]^n} \frac{1}{m^n} I(\sigma_{\mathbf{r}}(\mathbf{X}) \in A).$$

One may readily verify that $P_{\mathbf{X}}^o(\mathcal{A}^{m \times n}) = Q_{\mathbf{X}}^o(\mathcal{A}^{m \times n}) = 1$, so both $P_{\mathbf{X}}^o$ and $Q_{\mathbf{X}}^o$ are valid probability measures on $\mathcal{A}^{m \times n}$.

We will say that the set of cyclic shifts $\mathcal{S}_m(\mathbf{X})$ is *full* if its cardinality is equal to m^n , or equivalently, if all cyclic shifts of \mathbf{X} are distinct.

Lemma 1. *If $\mathcal{S}_m(\mathbf{X})$ is full, then (a) $P_{\mathbf{X}}^o = P_{\mathbf{X}}$ and (b) $Q_{\mathbf{X}}^o = Q_{\mathbf{X}}$.*

Proof. (a) For any $A \subseteq \mathcal{A}^{n \times m}$ we may write $P_{\mathbf{X}}(A)$ as

$$\frac{P_m(A \cap \mathcal{S}_m(\mathbf{X}))}{P_m(\mathcal{S}_m(\mathbf{X}))} = \sum_{\mathbf{r} \in [m]^n} \frac{P_m(\sigma_{\mathbf{r}}(\mathbf{X}))}{P_m(\mathcal{S}_m(\mathbf{X}))} I(\sigma_{\mathbf{r}}(\mathbf{X}) \in A). \quad (4.3)$$

Since $\mathcal{S}_m(\mathbf{X})$ is full,

$$P_m(\mathcal{S}_m(\mathbf{X})) = P_m\left(\bigcup_{\mathbf{r} \in [m]^n} \sigma_{\mathbf{r}}(\mathbf{X})\right) = \sum_{\mathbf{r} \in [m]^n} P_m(\sigma_{\mathbf{r}}(\mathbf{X})).$$

The independence of the rows of \mathbf{X} allows us to write the last expression as $\sum_{\mathbf{r} \in [m]^n} \prod_{i=1}^n p_m(\sigma_{r_i}(\mathbf{X}_{i \cdot}))$,

but this may be rewritten as $\prod_{i=1}^n \sum_{s \in [m]} p_m(\sigma_s(\mathbf{X}_{i \cdot}))$. Therefore (4.3) is equivalent to

$$\sum_{\mathbf{r} \in [m]^n} \left[\prod_{i=1}^n \frac{p_m(\sigma_{r_i}(\mathbf{X}_{i \cdot}))}{\sum_{s \in [m]} p_m(\sigma_s(\mathbf{X}_{i \cdot}))} \right] I(\sigma_{\mathbf{r}}(\mathbf{X}) \in A) = P_{\mathbf{X}}^o(A).$$

(b) There are m^n elements in $\mathcal{S}_m(\mathbf{X})$ when \mathbf{X} is full, so for any $A \in \mathcal{A}$

$$Q_{\mathbf{X}}(A) = m^{-n} |A \cap \mathcal{S}_m(\mathbf{X})| = \sum_{\mathbf{r} \in [m]^n} \frac{1}{m^n} I(\sigma_{\mathbf{r}}(\mathbf{X}) \in A) = Q_{\mathbf{X}}^o(A).$$

□

Lemma 2. *Let $\mathbf{x} = (x_0, x_1, \dots, x_{m-1}) \in \mathcal{A}^m$ be a sequence of length m . Let k be the least positive integer such that $\sigma_k(\mathbf{x}) = \mathbf{x}$. If $k < m$, then k divides m , and \mathbf{x} is equal to the repeated concatenation of a fixed block of length k .*

Proof. Suppose to the contrary that $1 \leq k < m$ does not divide m . Then we may write $m = kq + r$, where $1 \leq r < k$. Now $\sigma_k(\mathbf{x}) = \mathbf{x}$ implies that $\sigma_{-k}(\mathbf{x}) = \mathbf{x}$, and it follows that

$$\sigma_r(\mathbf{x}) = \sigma_{m-kq}(\mathbf{x}) = \sigma_m \circ \sigma_{-kq}(\mathbf{x}) = \mathbf{x}.$$

As this contradicts the minimality of k , we conclude that k divides m . The second conclusion follows in a straightforward way from the first. □

Corollary 1. *If $\mathbf{x} \in \mathcal{A}^m$ is such that $\sigma_k(\mathbf{x}) = \mathbf{x}$ for some $1 \leq k < m$, then \mathbf{x} contains two disjoint, equal blocks of length at least $m/3$.*

Lemma 3. *If (2.2) holds, then $P_m(\mathcal{S}_m(\mathbf{X}) \text{ is full})$ converges to 1 as m tends to infinity.*

Proof. We begin by noting that $\mathcal{S}_m(\mathbf{X})$ is full if $\mathcal{S}_m(\mathbf{X}_i)$ is full for $i = 1, \dots, n$. Because the rows of \mathbf{X} are independent, it therefore suffices to prove the result in the case $n = 1$. Thus we write $\mathbf{X} = (X_0, \dots, X_{m-1})$. If $\mathcal{S}_m(\mathbf{X})$ is not full, then Corollary 1 implies that there exist integers $l, r \geq m/3$ such that $X_j = X_{r+j}$ for $j = 0, \dots, l-1$. An easy calculation using the Markov property shows that, for fixed r , the P_m -probability of this event is at most ρ^{l-1} , where $\rho < 1$ is the maximum appearing in (2.2). Thus the probability that $\mathcal{S}_m(\mathbf{X})$ is not full is at most $m\rho^{m/3-1}$, which tends to zero as m tends to infinity. □

Definition: A set $A \subset \mathcal{A}^{n \times m}$ is *invariant under constant shifts* if $\sigma_{\mathbf{r}}(A) = A$ whenever $\mathbf{r} = (r, \dots, r)$ is a constant index sequence. Let \mathbb{A}_m be the family of all sets $A \subset \mathcal{A}^{n \times m}$ that are invariant under constant shifts.

Theorem 2. *Suppose that (2.2) holds and that the stationary Markov chain described by (4.1) is ergodic. Then*

$$\max_{A \in \mathbb{A}_m} |P_{\mathbf{X}}^o(A) - Q_{\mathbf{X}}^o(A)| \rightarrow 0$$

in probability as m tends to infinity.

Proof. Fix $m \geq 1$ and $A \in \mathbb{A}_m$. For $k \in \mathbb{Z}$ let $\mathbf{k}^* = (k, k, \dots, k) \in \mathbb{Z}^n$ be the constant sequence each of whose coordinates is equal to k . It follows from the invariance of A and the basic properties of cyclic shifts that for each $k \in \mathbb{Z}$,

$$\begin{aligned} P_{\mathbf{X}}^o(A) &= \sum_{\mathbf{r} \in [m]^n} \eta(\sigma_{\mathbf{r}}(\mathbf{X})) I(\sigma_{\mathbf{r}}(\mathbf{X}) \in A) \\ &= \sum_{\mathbf{r} \in [m]^n} \eta(\sigma_{\mathbf{r}+\mathbf{k}^*}(\mathbf{X})) I(\sigma_{\mathbf{r}+\mathbf{k}^*}(\mathbf{X}) \in A) \\ &= \sum_{\mathbf{r} \in [m]^n} \eta(\sigma_{\mathbf{r}+\mathbf{k}^*}(\mathbf{X})) I(\sigma_{\mathbf{r}}(\mathbf{X}) \in A). \end{aligned}$$

Thus we may express $P_{\mathbf{X}}^o(A)$ in the form of an average over k :

$$P_{\mathbf{X}}^o(A) = \sum_{\mathbf{r} \in [m]^n} \left[\frac{1}{m} \sum_{k \in [m]} \eta(\sigma_{\mathbf{r}+\mathbf{k}^*}(\mathbf{X})) \right] I(\sigma_{\mathbf{r}}(\mathbf{X}) \in A).$$

Combining this last expression with the definition of $Q_{\mathbf{X}}^{\circ}$ yields the bound

$$\begin{aligned}
|P_{\mathbf{X}}^{\circ}(A) - Q_{\mathbf{X}}^{\circ}(A)| &\leq \sum_{\mathbf{r} \in [m]^n} \left| \frac{1}{m} \sum_{k \in [m]} \eta(\sigma_{\mathbf{r}+\mathbf{k}^*}(\mathbf{X})) - \frac{1}{m^n} \right| I(\sigma_{\mathbf{r}}(\mathbf{X}) \in A) \\
&\leq \sum_{\mathbf{r} \in [m]^n} \left| \frac{1}{m} \sum_{k \in [m]} \eta(\sigma_{\mathbf{r}+\mathbf{k}^*}(\mathbf{X})) - \frac{1}{m^n} \right| \\
&= \frac{1}{m^n} \sum_{\mathbf{r} \in [m]^n} \left| \frac{1}{m} \sum_{k \in [m]} m^n \eta(\sigma_{\mathbf{r}+\mathbf{k}^*}(\mathbf{X})) - 1 \right|. \tag{4.4}
\end{aligned}$$

We now turn our attention to the quantity $m^n \eta(\sigma_{\mathbf{r}+\mathbf{k}^*}(\mathbf{X}))$ appearing in (4.4). Let $\mathbf{x} = (x_0, \dots, x_{m-1})$ be a fixed m -vector with entries in \mathcal{A} , and let $t \in \mathbb{Z}$. By expanding the joint probability $p_m(\cdot)$ as a product of one-step conditional probabilities and canceling common terms, a straightforward calculation shows that for all integers t

$$m \cdot \frac{p_m(\sigma_t(\mathbf{x}))}{\sum_{s \in [m]} p_m(\sigma_s(\mathbf{x}))} = \rho_t(\mathbf{x}) \gamma_m^{-1}(\mathbf{x}) \tag{4.5}$$

where

$$\rho_t(\mathbf{x}) = \frac{p_1(x_{[t]}) p_1(x_{[t-1]})}{p_2(x_{[t]}, x_{[t-1]})} \quad \text{and} \quad \gamma_m(\mathbf{x}) = \frac{1}{m} \sum_{j=0}^{m-1} \rho_j(\mathbf{x}).$$

(Recall that $[t] = t \bmod m$.) It follows from the definition of $\eta(\mathbf{X})$ and equation (4.5) that

$$\begin{aligned}
m^n \eta(\sigma_{\mathbf{r}+\mathbf{k}^*}(\mathbf{X})) &= \prod_{i=1}^n \left[\frac{m p_m(\sigma_{r_i+k}(\mathbf{X}_{i \cdot}))}{\sum_{s \in [m]} p_m(\sigma_s(\mathbf{X}_{i \cdot}))} \right] \\
&= \prod_{i=1}^n \rho_{r_i+k}(\mathbf{X}_{i \cdot}) \gamma_m^{-1}(\mathbf{X}_{i \cdot}) = \Gamma_m^{-1}(\mathbf{X}) \prod_{i=1}^n \rho_{r_i+k}(\mathbf{X}_{i \cdot})
\end{aligned}$$

where $\Gamma_m(\mathbf{X}) = \prod_{i=1}^n \gamma_m(\mathbf{X}_{i \cdot})$.

The assumptions of the theorem ensure that the random variables $X_{i,0}, \dots, X_{i,m-1}$ in the i th

row of \mathbf{X} are the initial terms of a stationary ergodic process, and therefore the same is true of the non-negative random variables $\rho_1(\mathbf{X}_i), \dots, \rho_{m-1}(\mathbf{X}_i)$. Note that the random variable $\rho_0(\mathbf{X}_i)$ cannot be included in this sequence because it involves the non-adjacent variables $X_{i,0}$ and $X_{i,m-1}$. It is easy to see that

$$\mathbb{E}\rho_1(\mathbf{X}_i) = \mathbb{E}\left(\frac{p_1(X_{i,1})p_1(X_{i,0})}{p_2(X_{i,0}, X_{i,1})}\right) = \sum_{u,v \in \mathcal{A}} \frac{p(u)p(v)}{p(u,v)} p(u,v) = 1.$$

From the ergodic theorem and the fact that $\rho_0(\mathbf{X}_i)$ is stochastically bounded (see (4.2)), it follows that

$$\gamma_m(\mathbf{X}_i) = \mathbb{E}\rho_1(\mathbf{X}_i) + o_P(1) = 1 + o_P(1), \quad (4.6)$$

and therefore $\Gamma_m(\mathbf{X})$ and $\Gamma_m^{-1}(\mathbf{X})$ are equal to $1 + o_P(1)$ as well. (Here and in what follows the stochastic order symbols $o_P(1)$ and $O_P(1)$ refer to the underlying measure P_m with m tending to infinity). For $\mathbf{r} \in [m]^n$ let $V_0(\mathbf{r}) = \{k \in [m] : r_i + k \equiv 0 \pmod{m} \text{ for some } 1 \leq i \leq n\}$ and let $V_1(\mathbf{r}) = [m] \setminus V_0(\mathbf{r})$. Note that $|V_0(\mathbf{r})| \leq n$ for each $\mathbf{r} \in [m]^n$. Combining the relation (4.6) with inequality (4.4) and equation (4.5), we conclude that

$$\begin{aligned} & |P_X^o(A) - Q_X^o(A)| \\ & \leq \Gamma_m^{-1}(\mathbf{X}) \cdot \frac{1}{m^n} \sum_{\mathbf{r} \in [m]^n} \left| \frac{1}{m} \sum_{k \in [m]} \prod_{i=1}^n \rho_{r_i+k}(\mathbf{X}_i) - 1 \right| + |\Gamma_m^{-1}(\mathbf{X}) - 1| \\ & = O_P(1) \cdot \frac{1}{m^n} \sum_{\mathbf{r} \in [m]^n} \left| \frac{1}{m} \sum_{k \in [m]} \prod_{i=1}^n \rho_{r_i+k}(\mathbf{X}_i) - 1 \right| + o_P(1) \\ & = O_P(1) \cdot \frac{1}{m^n} \sum_{\mathbf{r} \in [m]^n} \left| \frac{1}{m} \sum_{k \in V_1(\mathbf{r})} \prod_{i=1}^n \rho_{r_i+k}(\mathbf{X}_i) - 1 \right| + O_P(1) \Delta_m + o_P(1) \quad (4.7) \end{aligned}$$

where in the last line

$$\Delta_m := \frac{1}{m^n} \sum_{\mathbf{r} \in [m]^n} \frac{1}{m} \sum_{k \in V_0(\mathbf{r})} \prod_{i=1}^n \rho_{r_i+k}(\mathbf{X}_i).$$

As the upper bound in (4.7) is independent of our choice of $A \in \mathbb{A}_m$, it is enough to show that the first two terms in (4.7) are $o_P(1)$. Concerning the first term, by Markov's inequality it suffices to show that

$$\max_{\mathbf{r} \in [m]^n} \mathbb{E} \left| \frac{1}{m} \sum_{k \in V_1(\mathbf{r})} \prod_{i=1}^n \rho_{r_i+k}(\mathbf{X}_i) - 1 \right| \rightarrow 0 \text{ as } m \rightarrow \infty.$$

This follows from Corollary 2 below. As for the second term, note that

$$\begin{aligned} \Delta_m &\leq \prod_{i=1}^n (\rho_0(\mathbf{X}_i) \vee 1) \cdot \left[\frac{1}{m^n} \sum_{\mathbf{r} \in [m]^n} \frac{1}{m} \sum_{k \in V_0(\mathbf{r})} \prod_{i:r_i+k \neq 0} \rho_{r_i+k}(\mathbf{X}_i) \right] \\ &= O_P(1) \cdot \left[\frac{1}{m^n} \sum_{\mathbf{r} \in [m]^n} \frac{1}{m} \sum_{k \in V_0(\mathbf{r})} \prod_{i:r_i+k \neq 0} \rho_{r_i+k}(\mathbf{X}_i) \right] \end{aligned}$$

The term in brackets is non-negative and has expectation at most n/m . Thus $\Delta_m = o_P(1)$ and the result follows. \square

Let $\{U_1(k) : k \geq 0\}, \dots, \{U_n(k) : k \geq 0\}$ be independent, real-valued stationary ergodic processes defined on the same underlying probability space. Suppose that $\mathbb{E}|U_i(0)|$ is bounded for $i = 1, \dots, n$, and define $\mu = \prod_{i=1}^n \mathbb{E}(U_i(0))$. Let $\mathbf{r} = (r_1, \dots, r_n)$ denote a vector with non-negative integer-valued components. For $k, m \geq 1$ define random variables

$$V_m(k : \mathbf{r}) = \prod_{i=1}^n U_i((r_i + k) \bmod m) - \mu.$$

The independence of the processes $\{U_i(\cdot)\}$ ensures that $\mathbb{E}(V_m(k : \mathbf{r})) = 0$.

Lemma 4. *Under the assumptions above, $\max_{\mathbf{r} \in \mathbb{Z}^n} \mathbb{E} \left| m^{-1} \sum_{k=0}^{m-1} V_m(k : \mathbf{r}) \right|$ converges to zero as m*

tends to infinity.

Proof. Standard arguments show that the joint process $\{(U_1(k), \dots, U_n(k)) : k \geq 0\}$ is stationary and ergodic, and therefore the same is true for the process $\{\prod_{i=1}^n U_i(k) : k \geq 0\}$ of products. The L_1 ergodic theorem implies that

$$\Delta(l) = \mathbb{E} \left| \frac{1}{l} \sum_{k=0}^{l-1} \left(\prod_{i=1}^n U_i(k) - \mu \right) \right| \rightarrow 0 \text{ as } l \rightarrow \infty. \quad (4.8)$$

Note also that

$$\Delta(l) \leq \mathbb{E} \left| \prod_{i=1}^n U_i(0) - \mu \right| \leq 2 \prod_{i=1}^n \mathbb{E}|U_i(0)| \triangleq \Delta_0 \quad (4.9)$$

which is bounded by assumption.

Fix $m \geq 1$ and $\mathbf{r} = (r_1, \dots, r_n)$ with $r_i \geq 0$. Because the indices of $U_i(\cdot)$ in $V_m(k : \mathbf{r})$ are assessed modulo m , we may assume without loss of generality that $0 \leq r_1, \dots, r_n \leq m - 1$. Let $0 \leq r(1) < r(2) < \dots < r(n') \leq m - 1$ be the distinct order statistics of r_1, \dots, r_n , and note that $n' \leq n$. Define $r(0) = 0$, $r(n' + 1) = m$, and the differences $m_j = r(j + 1) - r(j)$ for $j = 0, \dots, n'$. Consider the decomposition

$$\sum_{k=0}^{m-1} V_m(k : \mathbf{r}) = \sum_{j=0}^{n'} W_j \text{ where } W_j = \sum_{k=r(j)}^{r(j+1)-1} V_m(k : \mathbf{r}). \quad (4.10)$$

The key feature of the sum W_j is this: for $r(j) \leq k \leq r(j + 1) - 1$ there are no “breaks” in the indexing of the terms $U_i((r_i + k) \bmod m)$ in $V_m(k : \mathbf{r})$ arising from the modular sum. In particular, there exist integers $\tilde{r}_1, \dots, \tilde{r}_n$ such that $(r_i + k) \bmod m = \tilde{r}_i + k$ for each $i = 1, \dots, n$, and each k in the sum defining W_j . As a result, the stationarity and independence of the individual processes

$\{U_i(\cdot)\}$ ensures that W_j is equal in distribution to the random variable

$$\tilde{W}_j = \sum_{k=0}^{m_j-1} \left(\prod_{i=1}^n U_i(k) - \mu \right).$$

We now turn our attention to the expectation in the statement of the lemma. It follows immediately from the decomposition (4.10) that

$$\frac{1}{m} \sum_{k=0}^{m-1} V_m(k : \mathbf{r}) = \sum_{j=0}^{n'} \frac{m_j}{m} \frac{1}{m_j} W_j,$$

which yields the elementary bound

$$\left| \frac{1}{m} \sum_{k=0}^{m-1} V_m(k : \mathbf{r}) \right| \leq \sum_{j=0}^{n'} \frac{m_j}{m} \left| \frac{1}{m_j} W_j \right|.$$

Taking expectations of both sides in the last display yields the inequality

$$\mathbb{E} \left| \frac{1}{m} \sum_{k=0}^{m-1} V_m(k : \mathbf{r}) \right| \leq \sum_{j=0}^{n'} \frac{m_j}{m} \mathbb{E} \left| \frac{1}{m_j} W_j \right| = \sum_{j=0}^{n'} \frac{m_j}{m} \mathbb{E} \left| \frac{1}{m_j} \tilde{W}_j \right| = \sum_{j=0}^{n'} \frac{m_j}{m} \Delta(m_j),$$

where the first equality follows from the distributional equivalence of W_j and \tilde{W}_j . In particular, for each integer $l \geq 1$ we have

$$\mathbb{E} \left| \frac{1}{m} \sum_{k=0}^{m-1} V_m(k : \mathbf{r}) \right| \leq \sum_{j:m_j \leq l} \frac{m_j}{m} \Delta(m_j) + \sum_{j:m_j > l} \frac{m_j}{m} \Delta(m_j) \leq \frac{nl\Delta_0}{m} + \sup_{l' > l} \Delta(l').$$

It follows from (4.8) and (4.9) that the final term in the last display tends to zero with m if $l = l(m)$ is any sequence such that l tends to infinity and $1/m$ converges to 0. Moreover, the final term does not depend on the vector r . This completes the proof of the lemma. \square

An elementary argument using Lemma 4 establishes the following corollary.

Corollary 2. *Under the assumptions of Lemma 4, $\max_{r \in \mathbb{N}^n} \mathbb{E} |m^{-1} \sum_{k'} V_m(k : r)|$ converges to zero as m tends to infinity, where for each r the sum is restricted to those $k' \in [m]$ such that $r_i + k' \not\equiv 0 \pmod{m}$.*

Proof of Theorem 1: Theorem 1 follows from Theorem 2 and the fact that for each $B \subseteq \mathbb{R}$ the event $\{\mathbf{Y} : T_m(\mathbf{Y}) \in B\} \in \mathbb{A}_m$ as T_m is invariant under constant shifts.

5 Discussion

High resolution genomic data is routinely used by biomedical investigators to search for recurrent genomic aberrations that are associated with disease. Cyclic shift testing provides a simple, permutation based approach to identify aberrant markers in a variety of settings. Here we establish finite sample, large marker asymptotics for the consistency of p -values produced by cyclic shift testing. The results apply to a broad family of Markov based null distributions. To our knowledge, this is the first theoretical justification of a testing procedure of this kind. Although cyclic shift testing was developed for DNA copy number analysis, we demonstrate its utility for DNA methylation and meta-analysis of genome wide association studies.

6 Acknowledgements

This research was supported by the National Institutes of Health (T32 CA106209 for VW), the Environmental Protection Agency (RD835166 for FAW), the National Institutes of Health/National Institutes of Mental Health (1R01MH090936-01 for FAW), and the National Science Foundation (DMS-0907177 and DMS-1310002 for ABN).

References

- [1] Beroukhi, R., Getz, G., Nghlemphu, L., Barretina, J., Hsueh, T., Linhart, D., Vivanco, I., Lee, J.C., Huang, J.H., Alexander, S., et al. (2007). “Assessing the significance of chromosomal aberrations in cancer: Methodology and application to glioma,” *Proc. Nat. Acad. Sci.* **104** 20007–20012.
- [2] Chitale, D., Gong, Y., Taylor, B.S., Broderick, S., Brennan, C., Somwar, R., Golas, B., Wang, L., Motoi, N., Szoke, J., et al. (2009). “An integrated genomic analysis of lung cancer reveals loss of *DUSP4* in *EGFR*-mutant tumors,” *Oncogene* **28** 2773–2783.
- [3] Fackler, M.J., Umbricht, C.B., Williams, D., Argani, P., Cruz, L-A., Merino, V.F., Teo, W.W., Zhang, Z., Huang, P., Visanathan, K., et al. (2011). “Genome-wide methylation analysis identifies genes specific to breast cancer hormone receptor status and risk of recurrence,” *Cancer Res.* **71** 6195–6207.
- [4] Fernando, M.M.A., Stevens, C.R., Walsh, E.C., De Jager, P.L., Goyette, P., Plenge, R.M., Vyse, T.J., Rioux, J.D. (2008). “Defining the role of the MHC in autoimmunity: a review and pooled analysis.” *PLoS Genet.* **4**(4): e1000024.
- [5] Laird, P.W., (2003). “The power and the promise of DNA methylation markers,” *Nature Rev.* **3** 253–266.
- [6] Misawa, K., Ueda, Y., Kanazawa, T., Misawa, Y., Jang, I., Brenner, J.C., Ogawa, T., Takebayashi, S., Krenman, R.A. Herman, J.G., et al. (2008). “Epigenetic inactivation of galanin receptor 1 in head and neck cancer,” *Clin. Cancer Res.* **14** 7604–7613.
- [7] Natrajan, R., Williams, R.D., Hing, S.N., Mackay, A., Reis-Filho, J.S., Fenwick, K., Irvani, M., Valgeirsson, H., Grigoriadis, A., Langford, C.F., et al. (2006). “Array CGH profiling of favourable

- histology Wilms tumours reveals novel gains and losses associated with relapse,” *J. Path.* **210** 49 – 58.
- [8] Pfeiffer, R.M., Gail, M.H., and Pee, D., (2009). “On combining data from genome-wide association studies to discover disease-associated SNPs,” *Stat. Sci.* **24**(4) 547–560.
- [9] Pinkel, D. and Albertson, D.G., (2005). “Array comparative genomic hybridization and its applications in cancer,” *Nature Genet.* **37** S11 – S17.
- [10] Renard, I., Joniau, S., van Cleynenbreugel, B., Collette, C., Naome, C., Vlassenbroeck, I., Nicolas, H., de Leval, J., Straub, J., Van Criekinge, W., et al. (2010). “Identification and validation of the methylated *TWIST1* and *NID2* genes through real-time methylation-specific polymerase chain reaction assays for the noninvasive detection of primary bladder cancer in urine samples,” *Eur. Urology* **58** 96–104.
- [11] Selamat, S.A., Chung, B.S., Girard, L., Zhang, W., Zhang, Y., Campan, M., Siegmund, K.D., Koss, M.N., Hagan, J.A., Lam, W.L., et al. (2012). “Genome-scale analysis of DNA methylation in lung adenocarcinoma and integration with mRNA expression,” *Genome Res.* doi:10.1101/gr.132662.111.
- [12] Thompson, J.R., Attia, J., and Minelli, C., (2011). “The meta-analysis of genome-wide association studies,” *Briefings Bioinf.* **12**(3) 259–269
- [13] Vire, E., Brenner, C., Deplus, R., Blanchon, L., Fraga, M., Didelot, C., Morey, L., Van Eynde, A., Bernard, D., Vanderwinder, J-M., et al. (2006). “The polycomb group protein EZH2 directly controls DNA methylation,” *Nature* **439**(16) 871–874.
- [14] Walter, V., Nobel, A.B., and Wright, F.A. (2011). “DiNAMIC: a method to identify recurrent DNA copy number aberrations in tumors,” *Bioinformatics* **27**(5) 678–685.

- [15] Weir, B.A., Woo, M.S., Getz, G., Perner, S., Ding, L., Beroukhi, R., Lin, W.M., Province, M.A., Kraja, A., Johnson, L.A., et al. (2007). “Characterizing the cancer genome in lung adenocarcinoma,” *Nature* **450** 893–901.
- [16] The Wellcome Trust Case Control Consortium, (2007). “Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls,” *Nature* **447**(7) 661–678.
- [17] Wistuba, I.I., Behrens, C., Virmani, A.K., Milchgrub, S., Syed, S., Lam, S., Mackay, B., Minna, J.D., and Gazdar, A.F. (1999). “Allelic losses at chromosome 8p21 - 23 are early and frequent events in the pathogenesis of lung cancer,” *Cancer Research* **59** 1973–1979.
- [18] Zeggini, E., Scott, L.J., Saxena, R., and Voight, B.F., for the Diabetes Genetics Replication and Meta-analysis (DIAGRAM) Consortium (2008). “Meta-analysis of genome-wide association data and large-scale replication identifies additional susceptibility loci for type 2 diabetes,” *Nature Genet.* **40**(5) 638–645.

7 Appendix

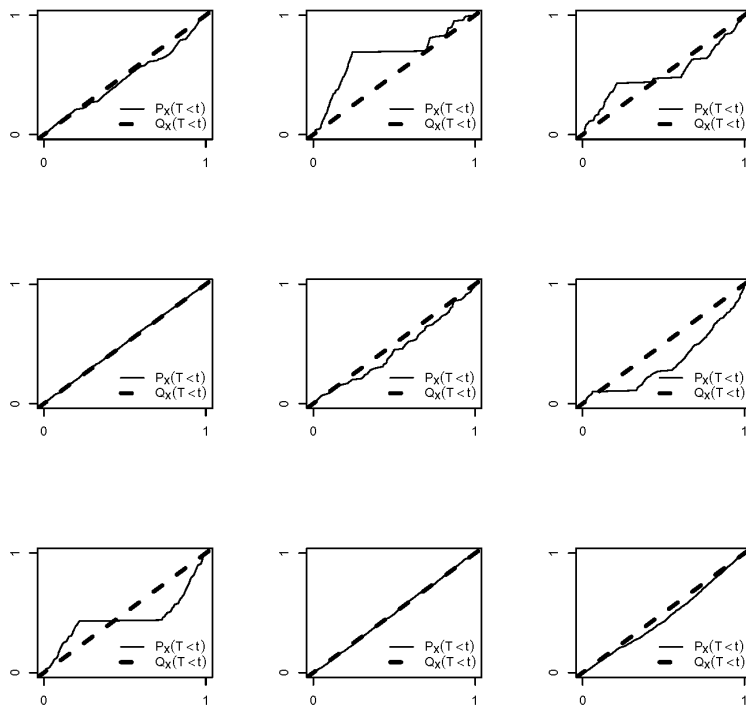


Figure A.1: Illustration of Resampling Distributions. Empirical cumulative distribution functions for simulated matrices \mathbf{X} in which independent rows \mathbf{X}_i are generated by a Gaussian AR(1) process with mean 0, standard deviation 1, and correlation .9. Each panel corresponds to a simulated 2×100 matrix \mathbf{X} .

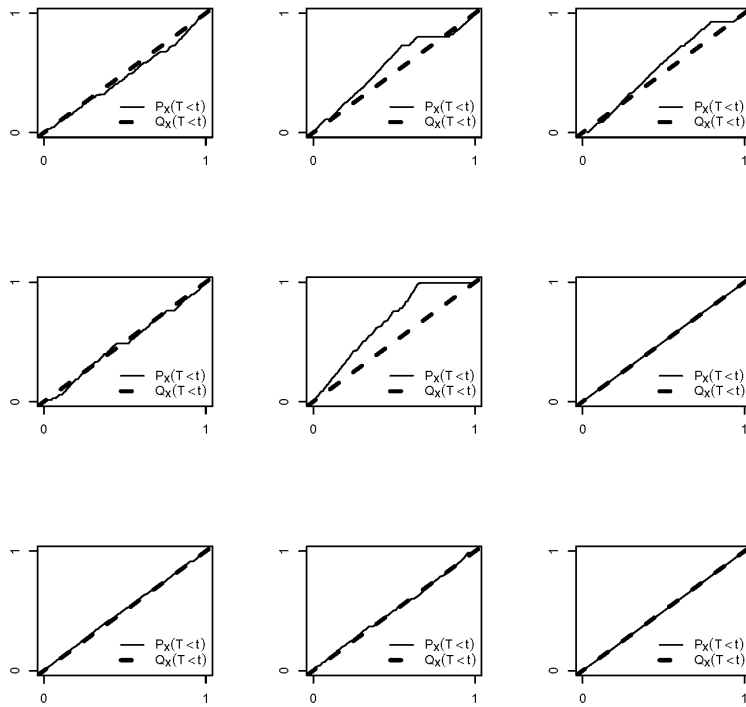


Figure A.2: Illustration of Resampling Distributions. Empirical cumulative distribution functions for simulated matrices \mathbf{X} in which independent rows \mathbf{X}_i are generated by a Gaussian AR(1) process with mean 0, standard deviation 1, and correlation .9. Each panel corresponds to a simulated 2×1000 matrix \mathbf{X} .

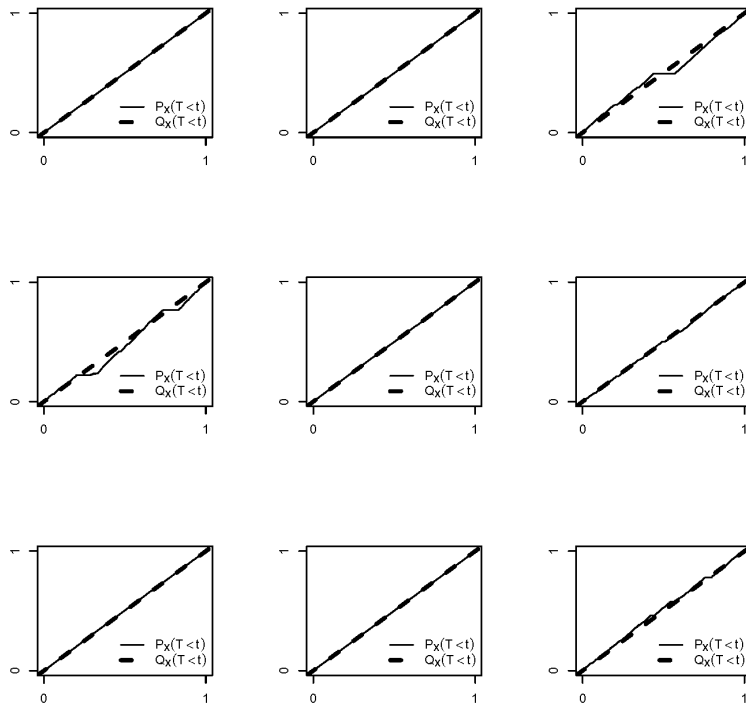


Figure A.3: Illustration of Resampling Distributions. Empirical cumulative distribution functions for simulated matrices \mathbf{X} in which independent rows \mathbf{X}_i are generated by a Gaussian AR(1) process with mean 0, standard deviation 1, and correlation .9. Each panel corresponds to a simulated 2×10000 matrix \mathbf{X} .

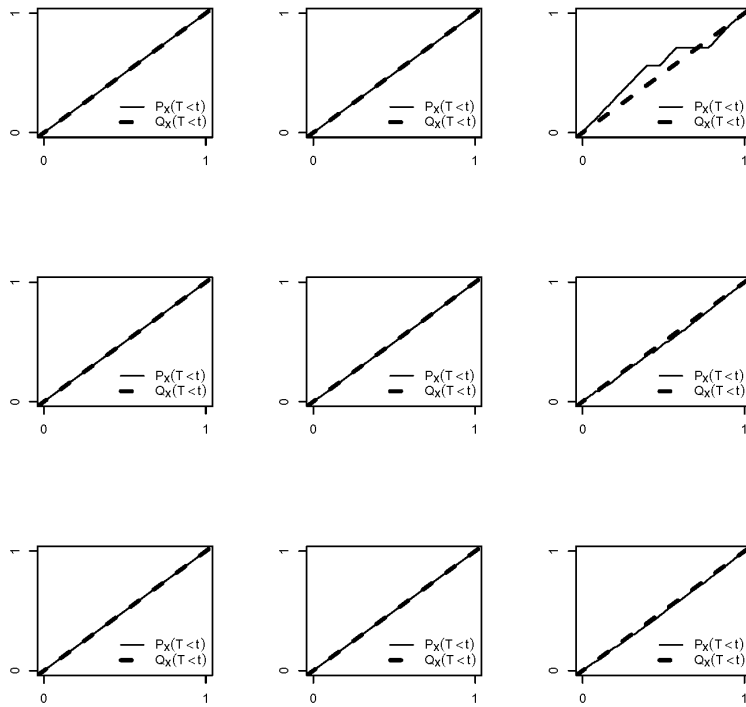


Figure A.4: Illustration of Resampling Distributions. Empirical cumulative distribution functions for simulated matrices \mathbf{X} in which independent rows \mathbf{X}_i are generated by a Gaussian AR(1) process with mean 0, standard deviation 1, and correlation .9. Each panel corresponds to a simulated 2×100000 matrix \mathbf{X} .

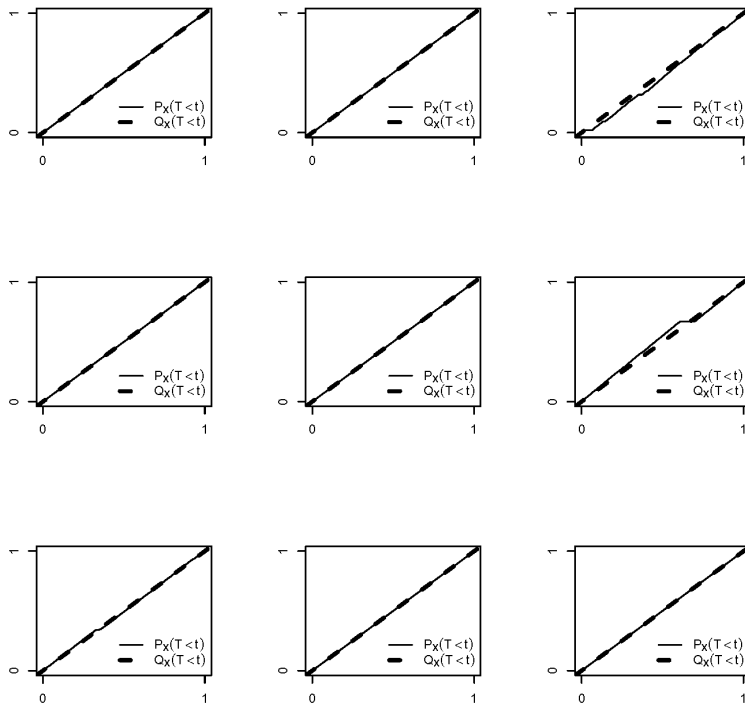


Figure A.5: Illustration of Resampling Distributions. Empirical cumulative distribution functions for simulated matrices \mathbf{X} in which independent rows \mathbf{X}_i are generated by a Gaussian AR(1) process with mean 0, standard deviation 1, and correlation .9. Each panel corresponds to a simulated 2×500000 matrix \mathbf{X} .