

NETWORK MODELLING OF TOPOLOGICAL DOMAINS USING HI-C DATA

BY Y. X. RACHEL WANG^{*,1}, PURNAMRITA SARKAR[†], OANA URSU^{‡,2},
ANSHUL KUNDAJE[‡] AND PETER J. BICKEL^{§,3}

University of Sydney^{*}, *University of Texas*[†], *Stanford University*[‡] and *University of California, Berkeley*[§]

Chromosome conformation capture experiments such as Hi-C are used to map the three-dimensional spatial organization of genomes. One specific feature of the 3D organization is known as topologically associating domains (TADs), which are densely interacting, contiguous chromatin regions playing important roles in regulating gene expression. A few algorithms have been proposed to detect TADs. In particular, the structure of Hi-C data naturally inspires application of community detection methods. However, one of the drawbacks of community detection is that most methods take exchangeability of the nodes in the network for granted; whereas the nodes in this case, that is, the positions on the chromosomes, are not exchangeable. We propose a network model for detecting TADs using Hi-C data that takes into account this nonexchangeability. In addition, our model explicitly makes use of cell-type specific CTCF binding sites as biological covariates and can be used to identify conserved TADs across multiple cell types. The model leads to a likelihood objective that can be efficiently optimized via relaxation. We also prove that when suitably initialized, this model finds the underlying TAD structure with high probability. Using simulated data, we show the advantages of our method and the caveats of popular community detection methods, such as spectral clustering, in this application. Applying our method to real Hi-C data, we demonstrate the domains identified have desirable epigenetic features and compare them across different cell types.

1. Introduction. In complex organisms, the genomes are very long polymers divided up into chromosomes and tightly packaged to fit in a minuscule cell nucleus. As a result, the packaging and the three-dimensional (3D) conformation of the chromatin have a fundamental impact on essential cellular processes including cell replication and differentiation. In particular, the 3D structure regulates the transcription of genes at multiple levels (Dekker (2008)). At the chromosome level, open (active) and closed (inactive) compartments alternate along chromosomes

Received October 2017; revised August 2018.

¹Supported in part by the ARC DECRA Fellowship.

²Supported in part by the HHMI International Student Research Fellowship and Stanford Gabilan Fellowship.

³Supported in part by NSF Grant DMS-17130833 and NIH Grant 1U01HG007031-01.

Key words and phrases. Hi-C data, topologically associating domains, network models, community detection.

(Lieberman-Aiden et al. (2009)) to form regions with clusters of active genes and repressed transcriptional activities, the latter typically partitioned to the nuclear periphery (Sexton et al. (2012), Smith et al. (2016)). At a smaller scale, chromatin loops make long-range regulations possible by bringing distant enhancers and repressors close to their target promoters.

Recently, one specific feature of chromatin organization known as topologically associating domains (TADs) has attracted much research attention. TADs are contiguous regions of chromatin with high levels of self-interaction and have been found in different cell types and species (Dixon et al. (2012), Sexton et al. (2012), Hou et al. (2012)). A number of studies have shown TADs contain clusters of genes that are co-regulated (Nora et al. (2012)) and may correlate with domains of histone modifications (Le Dily et al. (2014)), suggesting TADs act as functional units to help gene regulation. Disruptions of domain conformation have been associated with various diseases including cancer and limb malformation (Lupiáñez et al. (2015), Meaburn et al. (2009)).

While it is not possible to completely observe the 3D conformation, in the past decade several chromosome conformation capture technologies have been developed to measure the number of ligation events between spatially close chromatin regions. Hi-C is one of such technologies and provides genome-wide measurements of chromatin interactions using paired-end sequencing (Lieberman-Aiden et al. (2009)). The output can be summarized in a raw contact frequency matrix M , where M_{ij} is the total number of read pairs (which are interacting) falling into bins i and j on the genome. These equal-sized bins partition the genome and range from a few kilobases to megabases depending on the data resolution. Since TADs are regions with high levels of self-interactions, they appear as dense squares on the diagonal of the matrix.

A number of algorithms have been proposed to detect TADs, most of which rely on maximizing the intra-domain contact strength. This includes the earlier methods by Dixon et al. (2012) and Sauria et al. (2014), which summarize the 2D matrix as a 1D statistic to capture the changes in interaction strength at domain boundaries; and methods that directly utilize the 2D structure of the matrix to contrast the TAD squares from the background (Filippova et al. (2014), Lévy-Leduc et al. (2014), Malik and Patro (2015), Weinreb and Raphael (2016), Rao et al. (2014)). All of these methods use an optimization framework and apply standard dynamic programming to obtain the solution. The algorithms typically involve a number of tuning parameters with the number of TADs chosen in heuristic ways. More recently, Cabrerós, Abbe and Tsigos (2016) proposed to view the contact frequency matrix as an weighted undirected adjacency matrix for a network and applied community detection algorithms to fit mixed-membership block models.

Statistical networks provide a natural framework for modelling the 3D structure of chromatin as we can consider it as a spatial interaction network with positions on the genome as nodes. Network models have gained much popularity in numerous fields including social science, genomics and imaging; the availability of Hi-C

data opens new ground for applying network techniques, such as community detection, in order to answer important questions in biology. One of the drawbacks of community detection is that most of the methods take exchangeability of the nodes in the network for granted. However, modelling Hi-C data is a typical situation where the nodes, that is, the positions on the genome, are not exchangeable. In particular, since TADs are contiguous regions, treating TADs as densely connected communities imposes a geometric constraint on the community structure.

In this paper, we propose a network model for detecting TADs that incorporates the linear order of the nodes and preserves the contiguity of the communities found. Our main contributions include: (i) It has been observed empirically TADs are conserved across different cell types, but explicit joint analysis remains incomplete. Our likelihood-based method easily generalizes to allow for *joint inference* with multiple cell types. (ii) It has been postulated that CTCF (an insulator protein) acts as anchors at TAD boundaries (Nora et al. (2012), Sanborn et al. (2015)). Empirically, TAD boundaries correlate with CTCF sites, and modifications of binding motifs can lead to TAD disappearance (Sanborn et al. (2015)). Our model is flexible enough to include the positions of CTCF sites as *biological covariates*. (iii) We account for the existence of nested TADs. (iv) The core of our algorithm is based on linear programming, making it fast and efficient. (v) In addition, we provide theoretical justifications by analyzing the asymptotic performance of the algorithm and using automated model selection for choosing the number of TADs. The latter saves the need for many tuning parameters. Among these, (i) and (ii) are unique features of our method with biological significance.

The rest of the paper is organized as follows. We introduce the model and the estimation algorithm with asymptotic analysis in Section 2. In addition, we describe a post-processing step for testing the enrichment of contact within any TAD found. In Section 3, we first use simulated data to demonstrate the necessity of taking into account the linear ordering of the nodes and compare our method with other TAD detection algorithms. We next present the results of real data analysis for multiple human cell types, individually and jointly, using a publicly available Hi-C dataset (Rao et al. (2014)). We end the paper with a discussion of the advantages of our method and aspects for future work.

2. Methods. In this section, we describe a hierarchical network model for detecting nested TADs in a Hi-C contact frequency matrix using cell-line specific CTCF peaks as covariates. At each level of the hierarchy, we show the parameters can be estimated efficiently via coordinate ascent and provide asymptotic analysis of the algorithm. In addition, the model and algorithm can be adapted to identify TADs conserved across multiple cell lines. As further confirmation that the TADs found by the algorithm indeed correspond to regions of the genome with enriched interactions, we post process the candidate regions by performing a nonparametric test.

2.1. *Model description.* We consider a hierarchical model with a set of maximally nonoverlapping TADs at each level. In this section, we focus on describing the model for the base (outermost) level. The model and parameter estimation for the nested levels are identical and will be mentioned at the end of Section 2.2.

Let M denote a $n \times n$ contact frequency matrix. M is first thresholded at the q th quantile to produce a binary adjacency matrix A . Thresholding has been a common practice in network modeling to handle weighted matrices, despite the information loss it incurs. At canonical sequencing depth, the signal to noise ratio in Hi-C data is typically high and the resolution is relatively low. Thresholding can improve the signal to noise ratio. We examine the effect and sensitivity of the choice of q in Section 3.

As mentioned in the [Introduction](#), experimental evidence suggests TAD boundaries tend to coincide with CTCF binding. This motivates us to incorporate the presence of CTCF into our model. Let $Y \in \{0, 1\}^n$ be a binary vector with ones at positions where CTCF binding occurs. We will treat Y as an available covariate, which can be obtained from ChIP-seq data which is cell-type specific.

Let X denote a $n \times n$ binary matrix such that $X_{ab} = 1$ if (i) $Y_a = 1, Y_b = 1$ and (ii) there is a TAD between position a and b . X_{ab} is always 0 when $Y_a Y_b = 0$. This enforces the model to generate TADs which always have CTCF peaks at their boundaries. Thus $X \in \{0, 1\}^{n \times n}$ denotes a binary latent matrix which encodes the positions of all TADs. Also note that, it is possible to have $Y_a Y_b = 1$, but $X_{ab} = 0$, that is, there was no TAD formed between two CTCF binding sites.

We denote by the parameter vector $\Theta = (\beta, \{\alpha_{ab} : X_{ab} = 1\})$ the probabilities of edges between nodes. If $a \leq i < j \leq b$ for $X_{ab} = 1$, then $P(A_{ij} = 1) = \alpha_{ab}$. (Note that we allow for a different edge probability for each TAD.) Otherwise $P(A_{ij} = 1) = \beta$, which is also referred to as the background probability. The diagonal of A is set to 0. For simplicity we have assumed the connectivity within each TAD and the background is uniform, although the TADs may contain nested sub-TADs and can be heterogeneous. In general the contact frequency decreases as a function of the distance between two loci. For now one can think of the homogeneity assumption as approximating the actual distribution with a piecewise constant function, and we make use of the original weights in the post-processing step (Section 2.4). Finally, our model does not require the exact number of TADs, but only an upper bound on it. We will make this more concrete in Section 2.2.

REMARK 2.1. We demonstrate our model using a concrete example. The corresponding edge probability matrix is shown in Figure 1. In this example $Y_i = 1$ for $i \in \{3, 6, 12, 18\}$. We show positions where $Y_a Y_b = 1$ by red dots at the intersection of the grid lines, where the grid lines show the positions of the CTCF sites. Only X_{ab} at these positions are allowed to be one, since according to our model, TADs can only form between two CTCF sites. In this example, there are two TADs between 3 and 6 and between 12 and 18, that is, $X_{3,6} = 1, X_{12,18} = 1$. The edge probabilities may differ between the two TADs. The model naturally enforces contiguous clusters, and one cannot have a TAD with a hole inside.

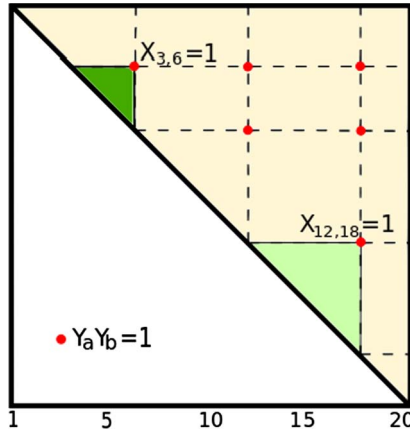


FIG. 1. Example of a probability matrix configuration.

2.2. *Parameter estimation.* Knowing (X, Y, Θ) , the maximization of the log likelihood for A can be written as

$$\begin{aligned}
 & \max_{X \in \{0,1\}^{n \times n}, \Theta} \log p(A; X, Y, \Theta) \\
 &= \frac{1}{2} \sum_{i \neq j} \sum_{a < b} Y_a Y_b X_{ab} \mathbb{I}_{i,j \in [a,b]} \left(A_{ij} \log \frac{\alpha_{ab}}{1 - \alpha_{ab}} + \log(1 - \alpha_{ab}) \right) \\
 (2.1) \quad &+ \frac{1}{2} \sum_{i \neq j} \left(1 - \sum_{a < b} Y_a Y_b X_{ab} \mathbb{I}_{i,j \in [a,b]} \right) \left(A_{ij} \log \frac{\beta}{1 - \beta} + \log(1 - \beta) \right) \\
 & \text{s.t. } \sum_{a < b} Y_a Y_b X_{ab} \leq K \\
 & \text{and } \sum_{c \leq a \leq d} Y_c Y_d X_{cd} \leq 1 \quad \text{for all } a \text{ s.t. } Y_a = 1,
 \end{aligned}$$

where \mathbb{I} is an indicator.

The first constraint upper bounds the total number of TADs at this level, while the second constraint ensures there is at most one TAD covering each position, thus making the TADs nonoverlapping. The likelihood implies it suffices to consider X_{ab} at positions such that both $Y_a = 1$ and $Y_b = 1$, and X is effectively a $m \times m$ matrix, where $m = \sum_a Y_a$. In this way the covariate vector Y helps reduce the search to a smaller grid.

We maximize the likelihood by considering a relaxed objective function and performing coordinate ascent. First note that taking the derivative of $\log p(A; X, Y, \Theta)$ with respect to α_{ab} , the estimate of α_{ab} does not depend on the other param-

eters and is given by

$$(2.2) \quad \hat{\alpha}_{ab} = \frac{\sum_{i,j \in [a,b]} A_{ij}}{(b-a+1)(b-a)}.$$

Therefore it remains to maximize the likelihood with respect to β and X . Since direct maximization of (2.1) over X subject to the constraints involve combinatorial optimization, we propose the following relaxed optimization,

$$\begin{aligned}
 & \max_{\beta, \pi \in [0,1]^{n \times n}} L(A, Y, \beta, \pi) \\
 & := \max_{\beta, \pi \in [0,1]^{n \times n}} \left\{ \frac{1}{2} \sum_{i \neq j} \sum_{a < b} Y_a Y_b \pi_{ab} \mathbb{I}_{i,j \in [a,b]} \right. \\
 & \quad \times \left[A_{ij} \log \frac{\hat{\alpha}_{ab}}{(1 - \hat{\alpha}_{ab})} + \log(1 - \hat{\alpha}_{ab}) \right] \\
 (LP-OPT) \quad & \left. + \frac{1}{2} \sum_{i \neq j} \left(1 - \sum_{a < b} Y_a Y_b \pi_{ab} \mathbb{I}_{i,j \in [a,b]} \right) \right. \\
 & \quad \left. \times \left[A_{ij} \log \frac{\beta}{1 - \beta} + \log(1 - \beta) \right] \right\}, \\
 & \text{s.t. } \sum_{a < b} Y_a Y_b \pi_{ab} \leq K \\
 & \text{and } \sum_{c \leq a \leq d} Y_c Y_d \pi_{cd} \leq 1 \quad \text{for all } a \text{ s.t. } Y_a = 1.
 \end{aligned}$$

The objective and constraints have the same form as (2.1) but with $\pi \in [0, 1]^{n \times n}$ replacing $X \in \{0, 1\}^{n \times n}$. Again since $\pi_{ab} = 0$ if $Y_a Y_b = 0$, the size of π to be estimated is effectively $m \times m$. This relaxed version can be solved via alternating maximization, also denoted by LP-OPT.

1. For each fixed β , (LP-OPT) is linear in π and can be maximized efficiently using linear programming.
2. For each fixed π , the objective is maximized at

$$(2.3) \quad \hat{\beta} = \frac{\sum_{i,j} A_{ij} - \sum_{a,b} \pi_{ab} Y_a Y_b \sum_{i,j \in [a,b]} A_{ij}}{n(n-1) - \sum_{a,b} \pi_{ab} Y_a Y_b (b-a+1)(b-a)}.$$

The above two steps are iterated until convergence in β .

So far we have described the model and parameter estimation for the outermost level of TADs. Within each of these TADs, we can repeat the same algorithm to detect the secondary (nested) level of TADs and continue iterating.

The likelihood approach allows the method to be easily extended to model conserved TADs across multiple cell lines. Assuming the cell lines are independent,

the joint log likelihood can be written as the sum,

$$(2.4) \quad \log p(\{A_\ell\}; X, Y, \{\Theta_\ell\}) = \sum_{\ell} \log p(A_\ell; X, Y, \Theta_\ell),$$

where X represents the latent positions of common TADs, Y is the set of CTCF peaks common to all cell lines; A_ℓ and Θ_ℓ are the adjacency matrix and model parameters specific to cell line ℓ . Similar to the single cell line case, the parameters can be estimated by using a plug-in estimator for each α_ℓ and alternating between maximizing over π and β_ℓ , where π is the relaxed form of X .

2.3. Theoretical guarantees. In this section, we analyze the theoretical properties of the algorithm and discuss the asymptotic performance of the estimates. Given that we have relaxed the original likelihood, it is natural to first check whether the solutions of (2.1) and LP-OPT agree. We have the following lemma stating optimizing the relaxed objective is essentially equivalent to optimizing the original one.

LEMMA 2.2. *For every given β ,*

$$(2.5) \quad \max_{\pi \in \Pi} L(A, Y, \beta, \pi) = \max_{X \in \mathcal{X}} L(A, Y, \beta, X),$$

where Π is the feasible set in LP-OPT and \mathcal{X} is the feasible set in (2.1).

PROOF. Given β , updating π is equivalent to maximizing the function

$$(2.6) \quad \begin{aligned} &L(A, Y, \Theta, \pi) \\ &= \frac{1}{2} \sum_{i \neq j} \sum_{a < b} Y_a Y_b \pi_{ab} 1_{i, j \in [a, b]} \left[A_{ij} \log \frac{\hat{\alpha}_{ab}(1 - \beta)}{(1 - \hat{\alpha}_{ab})\beta} + \log \frac{1 - \hat{\alpha}_{ab}}{1 - \beta} \right] \\ &\quad + \text{constant} \\ &:= l(A; \pi, \beta) + \text{constant}. \end{aligned}$$

Recalling $\hat{\alpha}_{ab}$ is independent of all the parameters, $l(A; \pi, \beta)$ is linear in π . Furthermore, the feasible set for π given in LP-OPT is a convex polyhedron with vertices at X . Since the optimum for a linear function on a convex polyhedron is always attained at the vertices, it follows then maximizing $l(A; \pi, \beta)$ with respect to π is equivalent to maximizing $l(A; X, \beta)$, which is the original objective. \square

The above lemma implies it is valid to analyze the solution of (2.1) even though the algorithm solves a relaxed problem. Furthermore, the optimal π for each run of step 1 in the algorithm belongs to the feasible set \mathcal{X} and defines a set of valid TAD positions (hence no thresholding is needed).

Next we analyze the asymptotics of the alternating optimization algorithm given a reasonable starting value β_0 and the upper bound K for the following setting.

We consider the most general case where each position is allowed a CTCF peak so Y_a will be omitted for the rest of the section. We focus on a single level of the hierarchical model and assume the $n \times n$ adjacency matrix A contains K^* TADs with $\{\alpha_1^*, \dots, \alpha_{K^*}^*\}$ as their connectivity probabilities. Note that to simplify notation, we have changed the subscript for α to a single index. The background has connectivity probability β^* . Let $\{[s_1, t_1], \dots, [s_{K^*}, t_{K^*}]\}$ be the TAD locations with the corresponding sizes $\{n_1^*, \dots, n_{K^*}^*\}$; $t_0 = 0, s_{K^*+1} = n + 1$ for convenience. We consider the case where K^* is fixed, $n_k^*/n \rightarrow p_k > 0$ for all k . In addition the sizes of the inter-TAD regions also follow $(s_{k+1} - t_k - 1)/n \rightarrow q_k$. Denote the number of inter-TAD regions G^* . Define $KL(s||t) = s \log(\frac{s}{t}) + (1 - s) \log(\frac{1-s}{1-t})$.

Assume the given β satisfies the following assumption:

ASSUMPTION 2.3. $\beta^* < \beta < \min_k \alpha_k^*$.

ASSUMPTION 2.4. For large enough n ,

$$\left((s_j - t_i - 1)^2 - \sum_{i < k < j} (n_k^*)^2 \right) KL(\beta^*||\beta) < \sum_{i < k < j} (n_k^*)^2 KL(\alpha_k^*||\beta)$$

for all $j > i + 1$. Note here $(s_j - t_i - 1)$ is the segment between the end of the i th TAD and the beginning of the j th TAD.

Note that when $\beta = \beta^*$, Assumption 2.4 is trivially satisfied.

THEOREM 2.5. Starting with $\beta^{(0)}$ satisfying Assumptions 2.3 and 2.4, for any fixed K and K large enough such that $K \geq K^* + G^*$, the optimal X satisfies

$$(2.7) \quad \exp\left\{ \max_{X \in \mathcal{X}} l(A; X, \beta^{(0)}) \right\} = \exp\{l(A; X_0, \beta^{(0)})\} (1 + o_P(1)),$$

where X_0 is such that $X_{s_k, t_k} = 1$ for all $1 \leq k \leq K^*$ and $X_{t_i+1, s_{i+1}-1} = 1$ for all $0 \leq i \leq K^*$. Furthermore, at the next iteration $\beta^{(1)} = \beta^* + O_P(n^{-1/2})$.

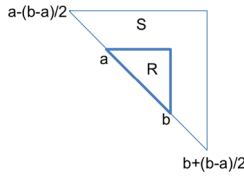
We defer the proof to the [Appendix](#). We have the following remarks.

1. Note that each $X \in \mathcal{X}$ partitions the nodes into $K + 1$ classes, given the partition the distribution of the edges follows a block model and the proofs utilize relevant techniques in this literature.

2. The theorem states that given an appropriate initial $\beta^{(0)}$, the optimal configuration found by the algorithm includes all the TADs as well as the inter-TAD regions. In the next section, we propose a nonparametric test to check enriched interactions within each candidate region called by the algorithm.

3. More importantly, the same optimal X_0 is found for any choice of fixed K , K being large enough. This implies the overfitting problem does not pose a serious concern here since increasing K does not always lead to an increase in the number of candidate TADs. In practice, a reasonable way to choose K is to increase it incrementally until the number of candidate TADs found starts to saturate.

2.4. *Post-processing.* After our algorithm detects the (possibly nested) TAD’s, our goal is to see if these indeed have higher contact frequencies than the surrounding region or the parent TAD. Recall that the contact frequency matrix M has nonnegative weights which are truncated to generate the adjacency matrix of the network. These weights M_{ij} , typically decay as $d = |i - j|$ grows. In order to detect TAD’s with significantly enriched contact frequencies over the surrounding region, we assume the model in equation (2.8). The main idea is that within a TAD, they decay slowly, whereas in the surrounding regions of a TAD they decay faster. Once we have detected the TADs using our linear program, we use these weights to prune weakly connected TADs. Consider the base level; let us assume that we have identified a TAD between positions a and b on the genome. Let the upper triangular region of the this TAD be denoted by R . Now consider the upper triangular region of the square between $a - \frac{a-b}{2}$ and $b + \frac{a-b}{2}$. Denote this by S . We assume the following simple model that dictates how the weights decay within and outside a TAD. Consider two monotonically decaying functions $f, g : \mathbb{N} \rightarrow \mathbb{R}^+ \cup \{0\}$, such that $f(d) > g(d) \forall d \in \mathbb{N}$, that is, $f(d)$ dominates $g(d)$ for any d .



$$(2.8) \quad M_{ij} = \begin{cases} f(|i - j|) + \epsilon_{ij}, & i, j \in R, \\ g(|i - j|) + \epsilon_{ij}, & i, j \in S \setminus R. \end{cases}$$

Here ϵ_{ij} are pairwise independent noise random variables.

Testing. In order to perform a test, for all $d \in \{1, \dots, (b - a)\}$, we calculate

$$\hat{f}(d) = \frac{\sum_{|i-j|=d, i, j \in R} M_{ij}}{b - a + 1 - d}, \quad \hat{g}(d) = \frac{\sum_{|i-j|=d, i, j \in S \setminus R} M_{ij}}{b - a}.$$

Now we take the two sequences \hat{f} and \hat{g} and do a nonparametric rank test (two-sample Wilcoxon test) to determine whether \hat{f} dominates \hat{g} ; if the p -value is smaller than a chosen threshold, we consider the TAD to have significant enrichment over its surrounding neighborhood. Otherwise we discard the TAD. For nested TADs, we are interested in determining whether a TAD found inside a parent TAD (call this T_0) is significant. In such cases, the surrounding region S may go across T_0 . So we simply truncate the outer region so that it does not cross outside T_0 .

3. Results. We first demonstrate key properties of our inference algorithm via simulation experiments, and then provide elaborate real data results.

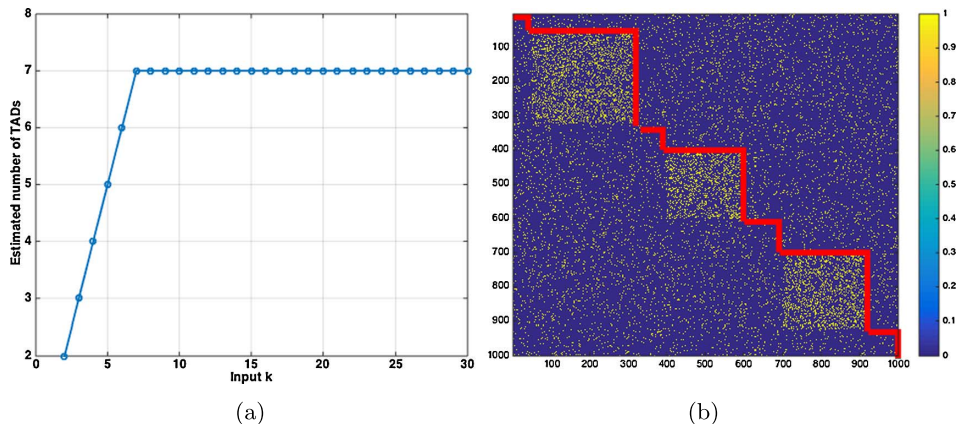


FIG. 2. (a) The y axis shows the estimated number of clusters K , whereas the x axis shows increasing values of K . (b) shows the clustering for input $K = 30$.

3.1. Simulations.

Data simulated under the simple model. First following the basic model described in Section 2 and equation (2.1), we present two sets of experiments to (a) show the robustness of our algorithm LP-OPT to the prespecified number of clusters, and (b) compare with the Spectral Clustering (SC) algorithm. For all the simulations in this setting, all TADs have the same linkage probability α and the background has linkage probability β .

In our first set of experiments (Figure 2(a) and (b)), we show that with somewhat balanced (but not necessarily equal) block-sizes, LP-OPT returns the correct TAD's along with some holes, as shown in Theorem 2.5. Recall that, in our linear program, we use a constraint to specify an upper bound on the number of TADs. This constraint is given by $\sum_{ij} \pi_{ij} \leq K$, where $\sum_{ij} \pi_{ij}$ represents the number of TADs. In Figure 2(a) we plot $\sum_{ij} \pi_{ij}$ after one iteration of the linear program, for the adjacency matrix in Figure 2(b). To be concrete, we set $n = 1000$, $\alpha = 0.2$, $\beta = 0.05$ and three TADs of sizes 270, 200 and 220. We also created CTCF sites at every 10 nodes for this experiment. We see that even though K is increased to 30, the estimated number of clusters levels off at 7, which is precisely three TADs plus four inter-TAD regions, which illustrates our asymptotic result from Theorem 2.5. These TADs detected by LP-OPT are illustrated in Figure 2(b). While one can come up with simple tests to eliminate the “spurious” TADs, we saw that for real data, our post processing step (see Section 2.4) eliminates them effectively for both the base level and nested TADs.

In the second set of simulations (Figure 3(a), (b)), we show that SC often yields clusters with holes, that is, clusters that are not contiguous, whereas we do not. SC is one of the most commonly used algorithms for community detection in networks. It involves performing spectral decomposition on a similarity matrix

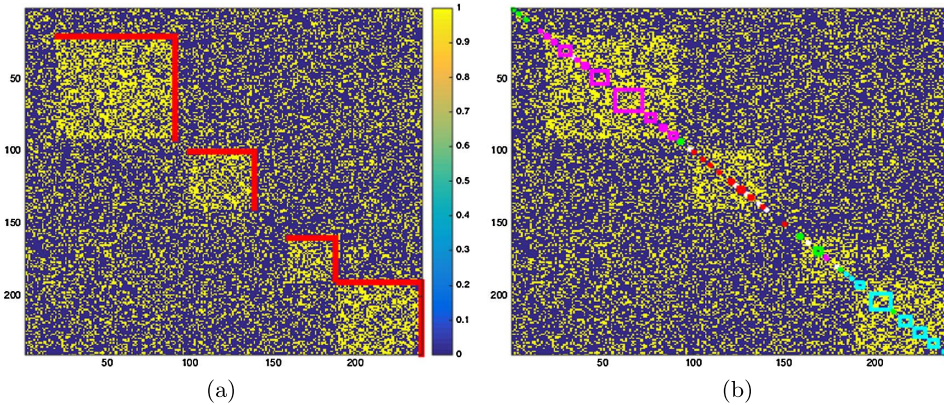


FIG. 3. Clusters identified by (a) LP-OPT and (b) SC. In (a) and (b) different colored squares correspond to different clusters detected by the algorithms. The ideal setting is to see a whole TAD encompassed by one square.

obtained from the data. For networks, one typically uses the normalized adjacency matrix defined as $D^{-1/2}AD^{-1/2}$ where D is the diagonal matrix of degrees, that is, $D = \text{diag}(d_i)$, $d_i = \sum_j A_{ij}$. Now for clustering the nodes into K blocks, one applies k-means clustering to the top K eigenvectors (Rohe, Chatterjee and Yu (2011)). For Figure 3(a) and (b), we set $n = 240$, four TADs with sizes (70, 40, 30, 50). The fifth cluster is the background. We use $\alpha = .5$, $\beta = .25$. In order to have a fair comparison, we do not include CTCF sites for LP-OPT, since SC is not designed to use them either. For both methods, we assume the correct number of blocks is given. In order to be as favorable as possible to SC, we use 4 top eigenvectors, and use k-means with $k = 5$ on these eigenvectors, since the background minus the TADs is one cluster. The results from conventional SC (choosing five top eigenvectors and using k-means with $k = 5$) are worse and hence omitted. For SC the plot reflects the clusterings returned: the colors correspond to different clusters. A square corresponds to a maximal contiguous set of nodes assigned to a cluster. For example, the last TAD (190–240) is assigned to the cyan cluster by SC. However, SC also assigns some nodes from the penultimate TAD (160–190) to this cluster, and moreover the small cyan boxes show that there are many nodes from the last TAD, which are assigned to other clusters, that is, in this setting, SC is unable to create a contiguous cluster will all nodes from one TAD.

We want to point out that while we did the above experiments for Figure 3(a) without the CTCF sites for fairness, including the CTCF sites greatly improves the computational time of LP-OPT. To be concrete, we simulated 10 random networks with the above setting, and obtained the clusterings with and without the CTCF sites. With CTCF sites LP-OPT converges in 0.5 seconds on average, whereas without CTCF sites, the average computation time is 58 seconds.

In Section S1.1 of the Supplementary Material (Wang et al. (2019)) we include additional comparison with SC for varying signal to noise ratio.

Data simulated using real data distribution. We next used a more realistic framework to simulate Hi-C data for chromosome 21 at a resolution of 40 kb, a typical resolution at which Hi-C data are analyzed. TAD positions were generated artificially and contact frequencies were sampled using empirical distributions from a real Hi-C dataset on chromosome 21 provided in Rao et al. (2014). A detailed description of the framework can be found in Section S1.2 of the Supplementary Material. CTCF sites were generated as the union of the true TAD boundaries and randomly sampled positions along the chromosome.

Our procedure led to a 1204×1204 contact frequency matrix, which was processed using a moving window of length 300 with an overlap of 50. The contact frequencies in each 300×300 segment were thresholded at the q th quantile to produce a binary adjacency matrix. Between two adjacent windows, any TADs called by the algorithm falling into overlapping regions are resolved as follows. (i) If the end point of the TAD is the last CTCF site in the first window, it is extended to the first CTCF site in the second window (similarly if the start point of the TAD is the first CTCF site in the second window; (ii) If one TAD is contained in another, the nested one is taken; (iii) If two TADs have a significant overlap (Jaccard index > 0.8 , defined in Section 3.2), they are merged by taking the intersection. A similar procedure is used on the real data (Section 3.2).

Table 1 compares the TADs found by our algorithm with ground truth using normalized mutual information (NMI) for different choices of the threshold q and an increasing number of randomly sampled CTCF positions. Note that the last column corresponds to the case where every position is a CTCF site, since the data generated contains 42 true TADs. In other words, we do not provide the algorithm with partial ground truth. As expected, the performance is better when partial ground truth is supplied but remains overall stable for reasonable choices of q . Figure 4 displays a 24 mb segment of the simulated data with TADs found by our method. LP-OPT was run without additional CTCF information and still achieved high similarity with ground truth.

In comparison, under similar thresholding levels SC achieves a NMI around 0.86–0.89 when the correct cluster number $K = 43$ (the last cluster being the background) is given, and the TADs found contain holes as described above. In

TABLE 1
Normalized mutual information measuring the quality of the TADs found vs. ground truth

# Random CTCF sites	Normalized mutual information				
	50	100	300	600	940
$q = 0.88$	0.90	0.91	0.90	0.89	0.89
$q = 0.9$	0.94	0.92	0.92	0.91	0.92
$q = 0.95$	0.97	0.98	0.95	0.93	0.92
$q = 0.98$	0.98	0.96	0.89	0.87	0.86

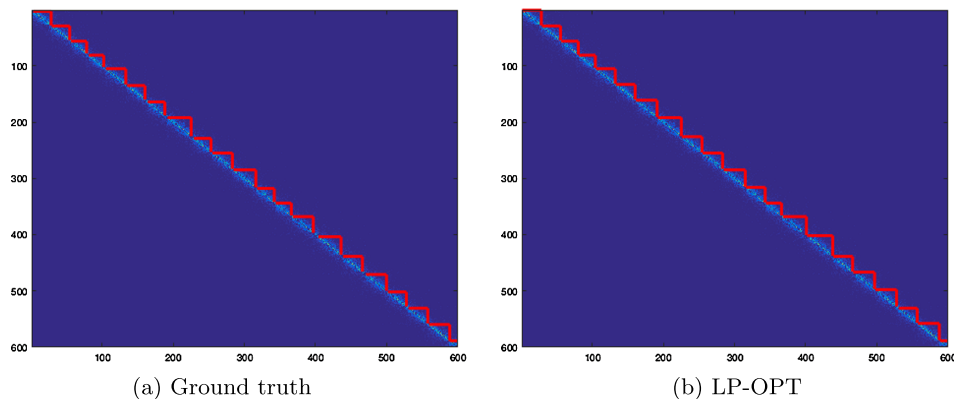


FIG. 4. *chr21:1–24000000*, a 24 mb segment. TADs identified by LP-OPT compared with ground truth. Note that LP-OPT was run without CTCF information thresholded at 95th quantile.

addition, we compare our method with two recently proposed TAD detection algorithms, 3DNetMod (Norton et al. (2018)) and MrTADFinder (Yan, Lou and Gerstein (2017)), which are both based on community detection methods in network analysis. The best NMI achieved by MrTADFinder for a range of tuning parameter values is 0.55. 3DNetMod had difficulty finding TADs on this dataset. These two methods will be included for comparison in the subsequent real data analysis. More details and visual comparisons can be found in Section S1.2 of the Supplementary Material.

3.2. *Real data.* Using the deep-coverage Hi-C data provided in Rao et al. (2014), we ran LP-OPT to identify cell-type specific TADs in five cell types (GM12878, HMEC, HUVEC, K562, NHEK) and common TADs conserved in all of them. We present here a comprehensive analysis of the results from chromosome 21. Similar analysis was also performed on chromosome 1, the results of which can be found in Section S2.2 of the Supplementary Material. Following Rao et al. (2014), the raw contact frequency matrix was normalized using the matrix balancing algorithm in Knight and Ruiz (2013). Using data with 10 kb resolution, the contact frequency matrix of this chromosome has more than 4800 bins. CTCF peaks for each cell type were obtained from the ENCODE pipeline (ENCODE Project Consortium (2012)) and converted into a binary vector of the same resolution as the contact frequency matrix, where each entry represents whether or not the corresponding genome bin contains at least one CTCF peak. This led to around 900 nonzero entries in each cell type. In the combined analysis for common TADs, we took the intersection of the cell-type specific CTCF binary vectors, so an entry is one only when the genome bin contains at least one CTCF peak in all cell types.

We performed TAD calling for three levels, each level with its own quantile thresholding parameter. At the base level, we processed the chromosome using

a moving window of length 300 (3 mb) with an overlap of 50. The contact frequencies in each 300×300 segment was thresholded at 90% quantile ($q_1 = 0.9$) to produce a binary adjacency matrix. Note that by using a moving window, we avoided using one universal threshold for the entire chromosome, which contains active and inactive regions with different chromatin interaction patterns. Any overlaps between two adjacent windows are resolved using the rules described in Section 3.1. The TADs called at the base level were then post-processed using the nonparametric test described in Section 2.4, and only those passing a p -value cutoff (in this case 0.05) were retained for further TAD calling. For the second level, we thresholded the contact frequencies inside the base-level TADs at 50% quantile ($q_2 = 0.5$), followed by running the algorithm and post-processing. The same steps were followed for the third level with $q_3 = 0.5$. For all three levels, the p -value cutoff was chosen to be 0.05. As a side note, correcting for multiple testing at a false discovery rate (FDR) of 0.05 made almost no difference at the base level. However, the same FDR cutoff led to fewer TADs being called at the nested levels. This is unsurprising as the power of the nonparametric test decreases as the number of data points available decreases at the nested levels.

The combined analysis for conserved TADs was performed in the same way, using the algorithm described in Section 2.2. The nonparametric test was run on the called regions for all cell types, and we required all the p -values to be smaller than the cutoff 0.05.

Choice of thresholds. We first checked the robustness of the results using different thresholding levels and biological replicates. Table 2 shows the number of TADs identified under different scenarios and with significant overlap. To compare two TADs S, T from two different sets, we measure the Jaccard index $J(S, T) = \frac{|S \cap T|}{|S \cup T|}$. When the Jaccard index is high enough, there is a one-to-one correspondence between TADs in the two sets. The first two rows in the table show different thresholds at the base level still lead to quite consistent results. Varying q_2, q_3 between 0.4–0.6 does not lead to noticeable changes and the results are hence omitted. Since two biological replicates (primary and replicate) are available for GM12878, we examined the consistency between them and the combined data, and the results are shown in row 3 and 4 of the table. Finally, as the current results were obtained using normalized data, we compared them with the case using the raw contact frequency matrix (row 5). This case still shows a reasonable degree of consistency despite having the lowest amount of overlap among all.

Enrichment of histone marks at boundaries. One of the most commonly used criteria for checking the accuracy of TAD boundaries is to count the number of histone modification peaks nearby (Filippova et al. (2014), Weinreb and Raphael (2016)) and taking higher levels of histone activity as indicators for the start and end points of TADs. The histone data are available in Kellis et al.

TABLE 2
Number of TADs detected under different scenarios and with significant overlap

# TADs		
$q_1 = 0.85$ (GM12878) 85	$q_1 = 0.9$ (GM12878) 81	Jaccard index > 0.7 70
$q_1 = 0.85$ (HMEC) 123	$q_1 = 0.9$ (HMEC) 114	Jaccard index > 0.7 103
Primary (GM12878) 90	Replicate (GM12878) 83	Jaccard index > 0.7 74
Primary (GM12878) 90	Combined (GM12878) 81	Jaccard index > 0.7 80
Normalized (GM12878) 81	Raw (GM12878) 94	Jaccard index > 0.7 61

(2014) and the processed data were downloaded from <https://sites.google.com/site/anshulkundaje/projects/encodehistonemods>. From bin indices, we obtain the coordinates of TAD boundaries by taking the midpoint of every genome bin. Table 3 shows the average number of peaks within 15 kb upstream or downstream from each detected boundary point for various types of histone modification. We compared LP-OPT with 3DNetMod (Norton et al. (2018)), MrTADFinder (Yan, Lou and Gerstein (2017)) and the Arrowhead domains originally reported in Rao et al. (2014). We found that MrTADFinder produced domains quite different from the other three methods (Figure S3 in the Supplementary Material) and the domain boundaries show less enrichment of histone marks. We have thus omitted the method from further comparison. The tuning parameters for 3DNetMod were chosen so that the number of TADs found is roughly comparable to the other two methods. In Table 3, counting the number of times each method achieves the highest enrichment, LP-OPT and 3DNetMod outperform Arrowhead with LP-OPT being slightly better than 3DNetMod. In addition, we note that LP-OPT is significantly faster than 3DNetMod, taking about 10 minutes on chromosome 21 (and 40 minutes on chromosome 1) using one core on a 3.1 GHz Intel Core i5 processor. In comparison, 3DNetMod takes more than 40 minutes on chromosome 21 (and four hours on chromosome 1) requiring four cores on the same processor. The results of 3DNetMod are also quite sensitive to the choice of tuning parameters.

Conserved and cell-type specific TADs. Although commonly used, the metric in Table 3 does not consider epigenetic features inside each TAD, which are particularly important for confirming shared regulatory structures and mechanisms across different cell types. We first examine the histone modification peaks within highly conserved TADs, which are defined as (i) TADs identified in the combined analysis of all cell types (denote this set \mathcal{I}_c), and (ii) if for $S \in \mathcal{I}_c$,

TABLE 3
Average number of histone modification peaks ± 15 kb upstream or downstream from the boundary points

GM12878					
	# domains	H3k9ac	H3k27ac	H3k4me3	Pol II
LP-OPT	81	1.35	1.68	1.16	1.29
Arrowhead	96	1.40	1.60	1.29	1.22
3DNetMod	129	1.36	1.82	1.25	1.06
HUVEC					
	# domains	H3k9ac	H3k27ac	H3k4me1	H3k4me3
LP-OPT	106	1.07	1.16	2.17	0.84
Arrowhead	59	1.02	1.08	2.06	0.83
3DNetMod	125	0.96	1.08	1.82	0.68
HMEC					
	# domains	H3k9ac	H3k27ac	H3k4me1	H3k4me3
LP-OPT	114	1.06	1.49	3.05	0.73
Arrowhead	44	1.02	1.46	3.09	0.81
3DNetMod	122	0.92	1.24	2.67	0.77
NHEK					
	# domains	H3k9ac	H3k27ac	H3k4me1	H3k4me3
LP-OPT	112	1.21	1.32	2.70	0.92
Arrowhead	78	0.99	1.12	2.19	0.68
3DNetMod	136	1.42	1.55	2.91	1.03
K562					
	# domains	H3k9ac	H3k4me1	H3k4me3	Pol II
LP-OPT	91	0.76	2.31	0.98	0.62
Arrowhead	82	0.57	1.82	0.82	0.55
3DNetMod	101	0.79	2.25	0.95	0.71

$\max_{T \in \mathcal{I}_i} J(S, T) > 0.7$ for all i , where \mathcal{I}_i is the set of TADs identified in cell type i . Out of the 50 TADs found in the combined analysis, 29 of them satisfy (ii).

Figure 5 shows the signal tracks for all five cell types inside one of the 29 conserved TADs (chr21:35275000–35725000) for two types of histone modifications (UCSC genome browser). The signal peaks are visibly correlated between cell types. Using ChIP-seq signals from the ENCODE pipeline, the average pairwise correlations between cell types for this TAD were calculated for different histone

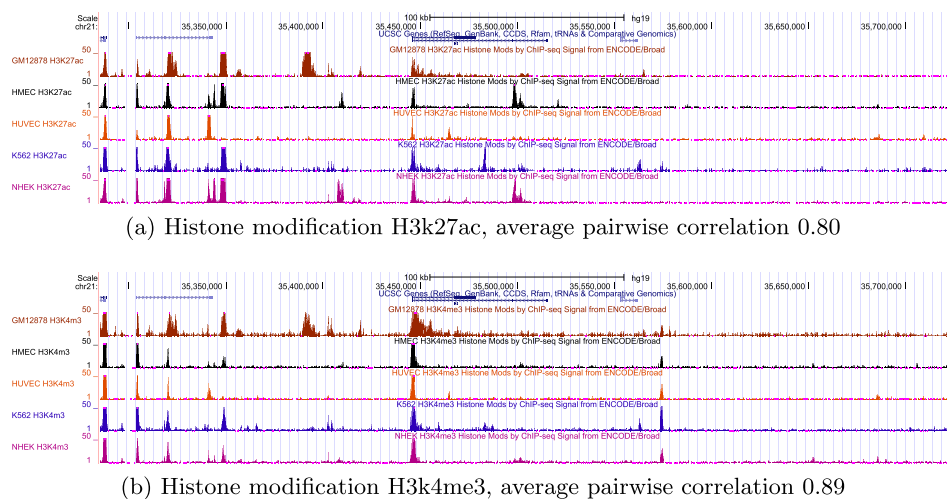


FIG. 5. Histone signal tracks within chr21:35275000–35725000.

modifications. For H3k9ac, H3k27ac, H3k4me1 and H3k4me3, the average correlations are 0.69, 0.80, 0.58, 0.89 respectively. Figure 6 compares the average pairwise correlations inside all 29 conserved TADs with 50 randomly chosen regions of length 290 kb (median length of the conserved TADs) on chromosome 21 for two instances of histone modification. The two-sample Wilcoxon test has p -values 0.05 and 0.006 for H3k27ac and H3k4me1; the results for H3k9ac and H3k4me3 are similar.

Having analyzed TADs with consistent overlaps across all cell types, we now consider TADs which are specific to individual cell types. A TAD is considered specific to that cell type i if (i) $S \in \mathcal{I}_i$; (ii) $\max_{T \in \mathcal{I}_j} J(S, T) < 0.4$ for all $j \neq i$.

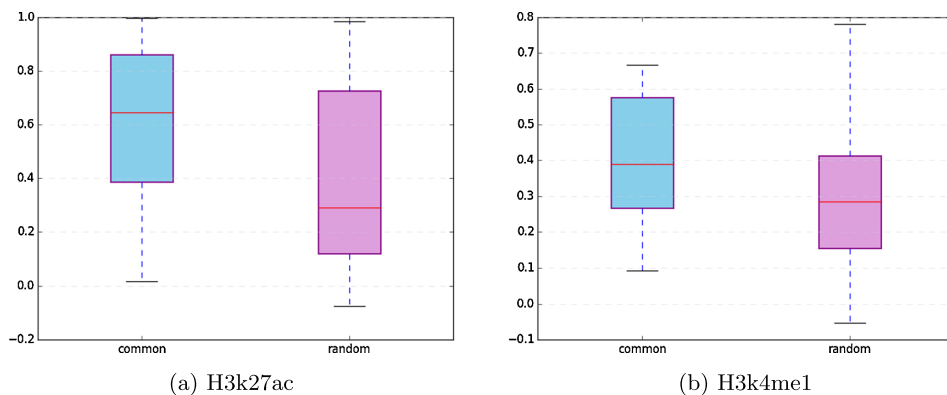


FIG. 6. Comparing conserved TADs with random regions on chr21; pairwise correlations between all cell types for (a) H3k27ac and (b) H3k4me1.

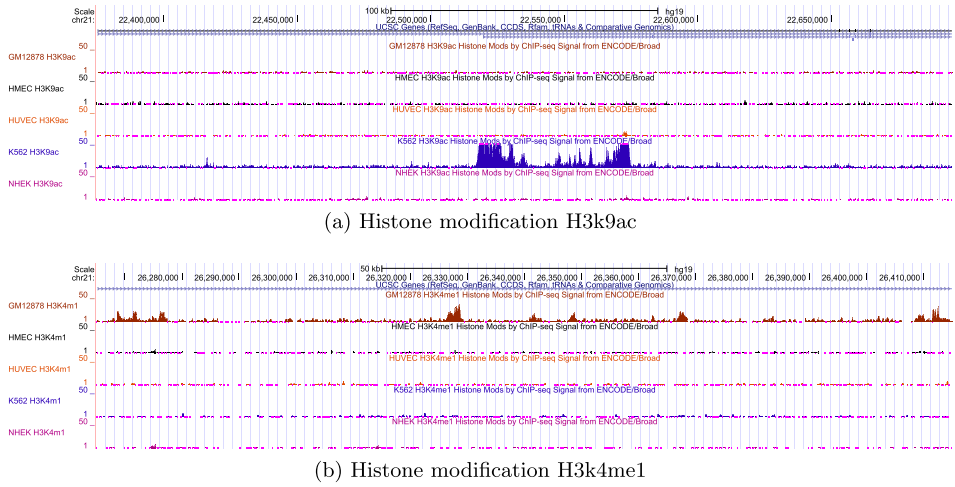


FIG. 7. (a) Signal tracks for H3k9ac within chr21:22375000–22695000, a TAD identified as specific to K562; (b) Signal tracks for H3k4me1 within chr21:26265000–26415000, a TAD identified as specific to GM12878.

This criterion leads to 28 TADs, each specific to one of the cell types. The median length of these TADs is 210 kb, smaller than that of the conserved TADs. As an illustration, Figure 7 shows the histone modification tracks inside two TADs specific to K562 and GM12878 respectively. In these two regions, the histone modifications show a higher level of activity for the two specific cell types. To evaluate whether this is a systematic trend, we next calculated the total signal level for each of the 28 TADs under different types of histone modifications. For each type, we counted the number of TADs which have the highest total signal level in the cell type they are associated with. Comparing to a null distribution under which the cell type with the highest total signal levels is selected randomly, we computed the p -values using a binomial distribution in Table 4. This suggests the cell-type specific TADs tend to be regions with more active histone modifications.

Without CTCF information. As a final remark, we tested whether LP-OPT could reproduce consistent results without CTCF information. Applying the al-

TABLE 4

For each type of histone modification, the number of TADs (out of 28) such that they have the highest total signal level in the cell type they are associated with

	H3k9ac	H3k27ac	H3k4me1	H3k4me3
# TADs with the highest total signal level	13	16	18	10
p -value	1.5×10^{-3}	1.7×10^{-5}	4.2×10^{-7}	3.9×10^{-2}

gorithm to a 5 mb segment of chromosome 21 (chr21:26000000–31000000) using GM12878 data, 15 TADs were identified (vs. 13 TADs with CTCF covariate) at the same p -value cutoffs. 11 pairs of these have a Jaccard index greater than 0.7, suggesting the results are reasonably stable. However, we also note the computational time in this case is significantly longer, as the search space for optimization is considerably larger without incorporating the CTCF sites.

4. Discussion. The 3D structure of chromatin provides key information for understanding the regulatory mechanisms. Recently, technologies such as Hi-C have revealed the existence of an important type of chromatin structure known as TADs, which are regions with enriched contact frequency and have been shown to act as functional units with coordinated regulatory actions inside. In this paper, we propose a statistical network model to identify TADs treating genome segments as nodes and their interactions in 3D as edges. Unlike many traditional networks with exchangeable distributions, our model incorporates the linear ordering of the nodes and guarantees the communities found represent contiguous regions on the genome. Our method also achieves two important biological goals: (i) Considering the empirical observation that TADs boundaries tend to correlate with CTCF binding sites, our method offers the flexibility to include CTCF binding data (or other ChIP-seq data) as biological covariates. (ii) The likelihood-based approach allows for joint inference across multiple cell types. On the theoretical side, we have shown asymptotic convergence of the estimation procedure with appropriate initializations. In practice, we observe the algorithm always converges in a few iterations. Due to the linear nature of the algorithm, our method is computationally efficient; it takes less than 10 minutes to complete the computation on chr21 with CTCF information on a laptop, whereas methods like TADtree (Weinreb and Raphael (2016)) can take up to hours.

Some areas for future work include extending the theoretical analysis to increasing K , and considering modelling higher order interactions between TADs. Our current way of finding conserved and cell-type specific TADs involves computing overlaps between domains and choosing heuristic cutoffs. While we have shown using epigenetic features that the conserved and cell-type specific TADs found have desirable features, it would be more ideal to statistically model the extent of overlaps between different types of TADs.

APPENDIX: PROOFS

Each $X \in \mathcal{X}$ partitions the nodes into $K + 1$ classes, thus we define the corresponding node labels as $Z = (Z_1, \dots, Z_n)$, with $Z_i = k$ if $Z_i \in [s_k, t_k]$, $Z_i = K^* + 1$ if Z_i does not fall inside any TAD. The set of feasible Z is a subset of $\{1, \dots, K^* + 1\}^n$ and can be seen as the latent node labels in a block model. Let \mathcal{X} and \mathcal{Z} be the feasible sets for X and Z respectively. X^* and Z^* are the true latent

positions and the corresponding node labels. Following block model notations, define a $(K^* + 1) \times (K^* + 1)$ matrix H^* , where $H_{k,k} = \alpha_{s_k, t_k}^*$ for $1 \leq k \leq K^*$, and $H_{k,l}^* = \beta^*$ otherwise. For any label Z , let $R(Z, Z^*)$ be the confusion matrix with

$$R_{k,l}(Z, Z^*) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(Z_i = k, Z_i^* = l).$$

Finally set $E = n \times n$ matrix of 1.

With appropriate concentration, it suffices to consider $l(A; \pi, \beta)$ at expectation $\mathbb{E}(A)$. Define

$$(A.1) \quad G(R, \beta) = \sum_{k=1}^{K^*} (RER^T)_{k,k} KL\left(\frac{(RH^*R^T)_{k,k}}{(RER^T)_{k,k}} \parallel \beta\right)$$

for some $Z \in \mathcal{Z}$ and its corresponding R . For simplicity of notation, we assume A has diagonal entries generated in the same way as nondiagonal entries, which does not affect the asymptotic results. We have the following lemma for the maximum of $G(\cdot, \beta)$.

LEMMA A.1. *Suppose β satisfies Assumptions 2.3 and 2.4. Then for all $K \geq K^* + G^*$ and n large enough,*

$$\begin{aligned} \max_{Z \in \mathcal{Z}} G(R(Z, Z^*), \beta) &= \frac{1}{n^2} \sum_{k=1}^{K^*} (n_k^*)^2 KL(\alpha_k^* \parallel \beta) \\ &\quad + \frac{1}{n^2} \sum_{i=0}^{K^*} (s_{i+1} - t_i - 1)^2 KL(\beta^* \parallel \beta). \end{aligned}$$

The maximum is unique at R_0 such that $X_{s_k, t_k} = 1$ for all $1 \leq k \leq K^*$ and $X_{t_i, s_{i+1}} = 1$ for all $0 \leq i \leq K^*$. Furthermore, for any $R_1 \neq R_0$ and n large enough,

$$(A.2) \quad \left. \frac{\partial G((1 - \epsilon)R_0 + \epsilon R_1, \beta)}{\partial \epsilon} \right|_{\epsilon=0^+} \leq -C < 0$$

for some $C > 0$.

PROOF. For each feasible Z , let $\{[l_1, m_1], \dots, [l_K, m_K]\}$ be the corresponding TAD positions defined by Z . For each row of $R(Z)$,

$$(A.3) \quad \begin{aligned} &(RER^T)_{k,k} KL\left(\frac{(RH^*R^T)_{k,k}}{(RER^T)_{k,k}} \parallel \beta\right) \\ &\leq \max \left\{ \left(\frac{(m_k - l_k)^2}{n^2} - \sum_{i=1}^{K^*} R_{ki}^2 \right) KL(\beta^* \parallel \beta), \sum_{i=1}^{K^*} R_{ki}^2 KL(\alpha_k^* \parallel \beta) \right\} \end{aligned}$$

by Assumption 2.3 and the convexity of $K(\cdot\|\beta)$. Also for the k th row of R , define

$$(A.4) \quad \begin{aligned} i_k &= \min\{i : [s_i, t_i] \cap [l_k, m_k] \neq \emptyset\}, \\ j_k &= \max\{i : [s_i, t_i] \cap [l_k, m_k] \neq \emptyset\}. \end{aligned}$$

We first consider the case where the set above is nonempty. For two adjacent rows k and $k + 1$, it suffices to consider the case $j_k = i_{k+1}$. Denote $S_{k,k+1} = \sum_{l=k}^{k+1} (RER^T)_{l,l} KL(\frac{(RH^*R^T)_{l,l}}{(RER^T)_{l,l}}\|\beta)$. By (A.3), $S_{k,k+1}$ is upper bounded by one of the following:

1. $((m_{k+1} - l_k)^2/n^2 - \sum_{q=k}^{k+1} \sum_{l=i_q}^{j_q} R_{ql}^2) KL(\beta^*\|\beta)$.
- 2.

$$\begin{aligned} &\sum_{l=i_k}^{j_k-1} R_{kl}^2 KL(\alpha_l^*\|\beta) + R_{k,j_k}^2 KL(\alpha_{j_k}^*\|\beta) \\ &+ \left((m_{k+1} - l_{k+1})^2/n^2 - \sum_{l=i_{k+1}}^{j_{k+1}} R_{k+1,l}^2 \right) KL(\beta^*\|\beta), \end{aligned}$$

which is itself upper bounded by

$$\begin{aligned} &\sum_{l=i_k}^{j_k-1} R_{kl}^2 KL(\alpha_l^*\|\beta) \\ &+ \max \left\{ \left(\frac{(m_{k+1} - s_{j_k})^2 - (n_{j_k}^*)^2}{n^2} - \sum_{l=i_{k+1}+1}^{j_{k+1}} R_{k+1,l}^2 \right) KL(\beta^*\|\beta), \right. \\ &\left. \left(\frac{n_{j_k}^*}{n} \right)^2 KL(\alpha_{j_k}^*\|\beta) + \left(\frac{(m_{k+1} - t_{j_k})^2}{n^2} - \sum_{l=i_{k+1}+1}^{j_{k+1}} R_{k+1,l}^2 \right) KL(\beta^*\|\beta) \right\}. \end{aligned}$$

- 3.

$$\begin{aligned} &\left((m_k - l_k)^2/n^2 - \sum_{l=i_k}^{j_k} R_{k,l}^2 \right) KL(\beta^*\|\beta) \\ &+ R_{k+1,i_{k+1}}^2 KL(\alpha_{i_{k+1}}^*\|\beta) + \sum_{l=i_{k+1}+1}^{j_{k+1}} R_{k+1,l}^2 KL(\alpha_l^*\|\beta). \end{aligned}$$

Similar to the case above, this is bounded by

$$\sum_{l=i_{k+1}+1}^{j_{k+1}} R_{k+1,l}^2 KL(\alpha_l^*\|\beta) + \max \left\{ \left(\frac{(t_{j_k} - l_k)^2 - (n_{j_k}^*)^2}{n^2} - \sum_{l=i_k}^{j_k-1} R_{kl}^2 \right) KL(\beta^*\|\beta), \right.$$

$$\left. \left(\frac{n_{j_k}^*}{n} \right)^2 KL(\alpha_{j_k}^* \parallel \beta) + \left(\frac{(s_{j_k} - l_k)^2}{n^2} - \sum_{l=i_k}^{j_k-1} R_{k,l}^2 \right) KL(\beta^* \parallel \beta) \right\}.$$

4. $\sum_{q=k}^{k+1} \sum_{l=i_q}^{j_q} R_{ql}^2 KL(\alpha_l^* \parallel \beta).$

If the set in (A.4) is empty,

$$\begin{aligned} (RER^T)_{k,k} KL\left(\frac{(RH^*R^T)_{k,k}}{(RER^T)_{k,k}} \parallel \beta\right) &= \left(\frac{m_k - l_k}{n}\right)^2 KL(\beta^* \parallel \beta) \\ &\leq (s_{l+1} - t_l)^2 KL(\beta^* \parallel \beta) \end{aligned}$$

for some $1 \leq l \leq K^*$.

The above cases show for any $Z \in \mathcal{Z}$, an upper bound for $G(R(Z, Z^*), \beta)$ is of the form

$$\begin{aligned} \sum_{k=1}^L \left(\frac{(s_{j_k} - t_{i_k})^2 - \sum_{i_k < l < j_k} (n_l^*)^2}{n^2} \cdot KL(\beta^* \parallel \beta) \right) \\ + \sum_{l \in \mathcal{I}} \frac{(n_l^*)^2}{n^2} KL(\alpha_l^* \parallel \beta), \end{aligned} \tag{A.5}$$

where \mathcal{I} is an index set such that $\mathcal{I} \cap_{k=1}^L [i_k, j_k] = \emptyset$. By Assumption 2.3, this is bounded by

$$\frac{1}{n^2} \sum_{k=1}^{K^*} (n_k^*)^2 KL(\alpha_k^* \parallel \beta) + \frac{1}{n^2} \sum_{i=0}^{K^*} (s_{i+1} - t_i - 1)^2 KL(\beta^* \parallel \beta)$$

with equality achieved only at R_0 for any $K \geq K^* + G^*$. The second part of the lemma can be checked with differentiation. \square

Let $[l_k, m_k]$ be the k th domain in a configuration Z corresponding to the k th row in R . Next we state a concentration lemma for the averages $\hat{\alpha}_{l_k, m_k}(Z)$. Denote $O_{l_k, m_k}(Z) = (m_k - l_k)^2 \hat{\alpha}_{l_k, m_k}(Z)$ and $\Delta_k(Z) = O_{l_k, m_k}(Z)/n^2 - (RH^*R^T(Z))_{k,k}$.

LEMMA A.2. For $\epsilon \leq 3$,

$$\mathbb{P}\left(\max_{Z \in \mathcal{Z}} \max_{1 \leq k \leq K} |\Delta_k(Z)| \geq \epsilon\right) \leq 2(K)^{n+1} \exp(-C_1(H^*)\epsilon^2 n^2).$$

Let $Z_0 \in \mathcal{Z}$ be a fixed set of labels, then for $\epsilon \leq 3m/n$,

$$\begin{aligned} \mathbb{P}\left(\max_{Z: |Z-Z_0| \leq m} \max_{1 \leq k \leq K} |\Delta_k(Z) - \Delta_k(Z_0)| > \epsilon\right) \\ \leq 2 \binom{n}{m} (K)^{m+1} \exp\left(-C_2(H^*) \frac{n^3 \epsilon^2}{m}\right). \end{aligned} \tag{A.7}$$

$C_1(H^*)$ and $C_2(H^*)$ are constants depending only on H^* .

PROOF. The proof follows from [Bickel and Chen \(2009\)](#) with minor modifications. \square

PROOF OF THEOREM 2.5. Suppose $\beta^{(0)}$ satisfies Assumptions 2.3 and 2.4. We consider the most general setup where every position is a CTCF binding site. The likelihood objective is given by

$$\begin{aligned}
 & l(A; Z, \beta^{(0)}) \\
 \text{(A.8)} \quad &= \frac{1}{2} \sum_{k=1}^K \sum_{\substack{i \neq j, \\ i, j \in [l_k, m_k]}} \left[A_{ij} \log \frac{\hat{\alpha}_{l_k, m_k}(Z)(1 - \beta^{(0)})}{(1 - \hat{\alpha}_{l_k, m_k}(Z))\beta^{(0)}} + \log \frac{1 - \hat{\alpha}_{l_k, m_k}(Z)}{1 - \beta^{(0)}} \right] \\
 &= \frac{1}{2} \sum_{k=1}^K (m_k - l_k)^2 KL(\hat{\alpha}_{l_k, m_k}(Z) \parallel \beta^{(0)}).
 \end{aligned}$$

Let R_0 (and the corresponding X_0, Z_0) be the optimal configuration in Lemma A.1.

We first consider X far away from X_0 . Define

$$I_{\delta_n} = \{X \in \mathcal{X} : G(R(X), \beta^{(0)}) - G(R_0, \beta^{(0)}) < -\delta_n\},$$

where δ_n is a sequence converging to 0 slowly. First by (A.6) in Lemma A.2,

$$\begin{aligned}
 & |l(A; X, \beta) - n^2 G(R(X), \beta)| \\
 \text{(A.9)} \quad &\leq Cn^2 \sum_{k=1}^K \left| \frac{O_{l_k, m_k}(X)}{n^2} - (RH^* R^T(X))_{k, k} \right| \\
 &= o_P(n^{2-\gamma})
 \end{aligned}$$

for some $\gamma < 1/2$. It follows then

$$\begin{aligned}
 & l(A; X, \beta^{(0)}) - l(A; X_0, \beta^{(0)}) \\
 &\leq o_P(n^{2-\gamma}) - n^2 \delta_n
 \end{aligned}$$

and

$$\text{(A.10)} \quad \exp\left\{ \max_{X \in I_{\delta_n}} l(A; X, \beta^{(0)}) - l(A; X_0, \beta^{(0)}) \right\}$$

$$\text{(A.11)} \quad \leq \sum_{X \in I_{\delta_n}} \exp\{l(A; X, \beta^{(0)}) - l(A; X_0, \beta^{(0)})\}$$

$$\text{(A.12)} \quad \leq \exp(o_P(n^{2-\gamma}) - n^2 \delta_n + n \log K) = o_P(1)$$

choosing $\delta_n \rightarrow 0$ slowly enough.

Next consider the case $X \in I_{\delta_n}^c$ and $X \neq X_0$. By (A.7) in Lemma A.2,

$$\begin{aligned}
 & \mathbb{P}\left(\max_{X \neq X_0} \|\Delta(Z) - \Delta(Z_0)\|_\infty > \epsilon |Z - Z_0|/n\right) \\
 (A.13) \quad & \leq \sum_{m=1}^n \mathbb{P}\left(\max_{Z:|Z-Z_0|=m} \|\Delta_k(Z) - \Delta_k(Z_0)\|_\infty > \epsilon \frac{m}{n}\right) \\
 & \leq \sum_{m=1}^n 2n^m K^{m+1} \exp(-Cmn) \rightarrow 0.
 \end{aligned}$$

It follows then if $|Z - Z_0| = m$, $\frac{\|\Delta(Z) - \Delta(Z_0)\|_\infty}{m/n} = o_P(1)$, and $\frac{1}{n^2} \|O(Z) - O(Z_0)\|_\infty \geq \frac{m}{n}(C + o_P(1))$ since $\|RH^*R^T(Z) - RH^*R^T(Z_0)\|_\infty \geq C\frac{m}{n}$. Note that in the set $I_{\delta_n}^c$, $|Z - Z_0| \rightarrow 0$. Then (A.2) implies

$$(A.14) \quad G(R(Z), \beta^{(0)}) - G(R_0(Z_0), \beta^{(0)}) < -C\frac{m}{n}$$

if $|Z - Z_0| = m$. Since $G(R, \beta^{(0)})$ is the population version of $\frac{1}{n^2}l(A; R, \beta^{(0)})$ and $O(Z)/n^2$ approaches $RH^*R^T(Z)$ uniformly in probability, by the continuity of the derivative,

$$(A.15) \quad \frac{1}{n^2}(l(A; R(Z_0), \beta^{(0)}) - l(A; R(Z), \beta^{(0)})) = \Omega_P(m/n)$$

for $|Z - Z_0| = m$. It follows then

$$\begin{aligned}
 & \exp\left\{\max_{X \in I_{\delta_n}^c, X \neq X_0} l(A; X, \beta^{(0)}) - l(A; X_0, \beta^{(0)})\right\} \\
 (A.16) \quad & \leq \sum_{X \in I_{\delta_n}^c, X \neq X_0} \exp\{l(A; X, \beta^{(0)}) - l(A; X_0, \beta^{(0)})\} \\
 & \leq \sum_{m=1}^n \binom{n}{m} (K + 1)^m e^{-\Omega_P(mn)} = o_P(1). \quad \square
 \end{aligned}$$

Acknowledgements. The authors would like to thank Nathan Boley for helpful initial discussions and the reviewers and the AE for their valuable comments which led to an improved version of the paper.

SUPPLEMENTARY MATERIAL

Supplementary information for “Network modelling of topological domains using Hi-C data” (DOI: [10.1214/19-AOAS1244SUPP](https://doi.org/10.1214/19-AOAS1244SUPP); .pdf). The supplementary material includes code for TAD calling and the TAD coordinates called on chr1 and chr21 as the txt files. Additional simulations and real data results can be found in the supplementary file Wang et al. (2019).

REFERENCES

- BICKEL, P. J. and CHEN, A. (2009). A nonparametric view of network models and Newman–Girvan and other modularities. *Proc. Natl. Acad. Sci. USA* **106** 21068–21073.
- CABREROS, I., ABBE, E. and TSIRIGOS, A. (2016). Detecting community structures in hi-c genomic data. In *Information Science and Systems (CISS), 2016 Annual Conference on* 584–589. IEEE Press, New York.
- DEKKER, J. (2008). Gene regulation in the third dimension. *Science* **319** 1793–1794.
- DIXON, J. R. et al. (2012). Topological domains in mammalian genomes identified by analysis of chromatin interactions. *Nature* **485** 376–380.
- ENCODE PROJECT CONSORTIUM (2012). An integrated encyclopedia of DNA elements in the human genome. *Nature* **489** 57–74.
- FILIPPOVA, D., PATRO, R., DUGGAL, G. and KINGSFORD, C. (2014). Identification of alternative topological domains in chromatin. *Algorithms Mol. Biol.* **9** 14.
- HOU, C., LI, L., QIN, Z. S. and CORCES, V. G. (2012). Gene density, transcription, and insulators contribute to the partition of the drosophila genome into physical domains. *Molecular Cell* **48** 471–484.
- KELLIS, M. et al. (2014). Defining functional dna elements in the human genome. *Proc. Natl. Acad. Sci. USA* **111** 6131–6138.
- KNIGHT, P. A. and RUIZ, D. (2013). A fast algorithm for matrix balancing. *IMA J. Numer. Anal.* **33** 1029–1047. [MR3081493](#)
- LE DILY, F. et al. (2014). Distinct structural transitions of chromatin topological domains correlate with coordinated hormone-induced gene regulation. *Genes & Development* **28** 2151–2162.
- LÉVY-LEDUC, C., DELATTRE, M., MARY-HUARD, T. and ROBIN, S. (2014). Two-dimensional segmentation for analyzing hi-c data. *Bioinformatics* **30** i386–i392.
- LIEBERMAN-AIDEN, E. et al. (2009). Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science* **326** 289–293.
- LUPIÁÑEZ, D. G. et al. (2015). Disruptions of topological chromatin domains cause pathogenic rewiring of gene-enhancer interactions. *Cell* **161** 1012–1025.
- MALIK, L. I. and PATRO, R. (2015). Rich chromatin structure prediction from hi-c data. *bioRxiv*, page 032953.
- MEABURN, K. J., GUDLA, P. R., KHAN, S., LOCKETT, S. J. and MISTELI, T. (2009). Disease-specific gene repositioning in breast cancer. *J. Cell Biol.* **187** 801–812.
- NORA, E. P. et al. (2012). Spatial partitioning of the regulatory landscape of the x-inactivation centre. *Nature* **485** 381–385.
- NORTON, H. K., EMERSON, D. J., HUANG, H., KIM, J., TITUS, K. R., GU, S., BASSETT, D. S. and PHILLIPS-CREMIN, J. E. (2018). Detecting hierarchical genome folding with network modularity. *Nat. Methods* **15** 119–122.
- RAO, S. S. et al. (2014). A 3d map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell* **159** 1665–1680.
- ROHE, K., CHATTERJEE, S. and YU, B. (2011). Spectral clustering and the high-dimensional stochastic blockmodel. *Ann. Statist.* **39** 1878–1915. [MR2893856](#)
- SANBORN, A. L., RAO, S. S. P., HUANG, S.-C., DURAND, N. C., HUNTLEY, M. H., JEWETT, A. I., BOCHKOV, I. D., CHINNAPPAN, D., CUTKOSKY, A. et al. (2015). Chromatin extrusion explains key features of loop and domain formation in wild-type and engineered genomes. *Proc. Natl. Acad. Sci. USA* **112** E6456–E6465.
- SAURIA, M. E., PHILLIPS-CREMIN, J. E., CORCES, V. G. and TAYLOR, J. (2014). Hifive: A normalization approach for higher-resolution hic and 5c chromosome conformation data analysis. Available at <https://www.biorxiv.org/content/10.1101/009951v1.full>.
- SEXTON, T. et al. (2012). Three-dimensional folding and functional organization principles of the drosophila genome. *Cell* **148** 458–472.

- SMITH, E. M., LAJOIE, B. R., JAIN, G. and DEKKER, J. (2016). Invariant TAD boundaries constrain cell-type-specific looping interactions between promoters and distal elements around the CFTR locus. *Am. J. Hum. Genet.* **98** 185–201.
- WANG, Y. X. R., SARKAR, P., URSU, O., KUNDAJE, A. and BICKEL, P. J. (2019). Supplement to “Network modelling of topological domains using Hi-C data.” DOI:[10.1214/19-AOAS1244SUPP](https://doi.org/10.1214/19-AOAS1244SUPP).
- WEINREB, C. and RAPHAEL, B. J. (2016). Identification of hierarchical chromatin domains. *Bioinformatics* **32** 1601–1609.
- YAN, K.-K., LOU, S. and GERSTEIN, M. (2017). MrTADFinder: A network modularity based approach to identify topologically associating domains in multiple resolutions. *PLoS Comput. Biol.* **13** e1005647.

Y. X. R. WANG
SCHOOL OF MATHEMATICS AND STATISTICS F07
UNIVERSITY OF SYDNEY
NEW SOUTH WALES 2006
AUSTRALIA
E-MAIL: rachel.wang@sydney.edu.au

O. URSU
DEPARTMENT OF GENETICS
STANFORD UNIVERSITY
300 PASTEUR DR.
LANE L301
STANFORD, CALIFORNIA 94305
USA
E-MAIL: oursu@broadinstitute.org

P. SARKAR
DEPARTMENT OF STATISTICS AND DATA SCIENCES
UNIVERSITY OF TEXAS, AUSTIN
2317 SPEEDWAY STOP D9800
AUSTIN, TEXAS 78712
USA
E-MAIL: purna.sarkar@austin.utexas.edu

A. KUNDAJE
DEPARTMENT OF GENETICS
STANFORD UNIVERSITY
AND
DEPARTMENT OF COMPUTER SCIENCE
STANFORD UNIVERSITY
300 PASTEUR DR.
LANE L301
STANFORD, CALIFORNIA 94305
USA
E-MAIL: akundaje@stanford.edu

P. J. BICKEL
DEPARTMENT OF STATISTICS
UNIVERSITY OF CALIFORNIA, BERKELEY
367 EVANS HALL
BERKELEY, CALIFORNIA 94720
USA
E-MAIL: bickel@stat.berkeley.edu