

NORMAL APPROXIMATION FOR FUNCTIONS OF HIDDEN MARKOV MODELS

CHRISTIAN HOUDRÉ,* *Georgia Institute of Technology*
GEORGE KERCHEV,** *Université du Luxembourg*

Abstract

The generalized perturbative approach is an all-purpose variant of Stein's method used to obtain rates of normal approximation. Originally developed for functions of independent random variables, this method is here extended to functions of the realization of a hidden Markov model. In this dependent setting, rates of convergence are provided in some applications, leading, in each instance, to an extra log-factor vis-à-vis the rate in the independent case.

Keywords: Stein's method; Markov chains; generalized perturbative approach; normal approximation; stochastic geometry

2020 Mathematics Subject Classification: Primary 60F05
Secondary 60K35; 60D05

1. Introduction

Let $X = (X_1, \dots, X_n)$ be a random vector with coordinates in a Polish space E , and let $f : E^n \rightarrow \mathbb{R}$ be a measurable function such that $f(X)$, for n large, is square-integrable. For a large class of such functions f it is expected that as n grows without bound, $f(X)$ behaves like a normal random variable. To quantify such estimates one is interested in bounding the distance between $f(X)$ and the normal random variable $\mathcal{N} \sim N(m_f, \sigma_f^2)$ where $m_f = \mathbb{E}[f(X)]$ and $\sigma_f^2 = \text{Var}(f(X))$. Two such distances of interest are the Kolmogorov distance

$$d_K(f(X), \mathcal{N}) := \sup_{t \in \mathbb{R}} |\mathbb{P}(f(X) \leq t) - \mathbb{P}(\mathcal{N} \leq t)|$$

and the Wasserstein distance

$$d_W(f(X), \mathcal{N}) := \sup_h |\mathbb{E}[h(f(X))] - \mathbb{E}[h(\mathcal{N})]|,$$

where this last supremum is taken over real-valued Lipschitz functions h such that $|h(x) - h(y)| \leq |x - y|$ for all $x, y \in \mathbb{R}$.

For the case where the components of X are independent random variables, upper bounds on $d_W(f(X), \mathcal{N})$ were first obtained in [2], and these were extended to $d_K(f(X), \mathcal{N})$ in [14]. Both results rely on a class of difference operators that will be described in Section 2.

Received 9 October 2020; revision received 15 July 2021.

* Postal address: School of Mathematics, Georgia Institute of Technology, Atlanta, Georgia 30332-0160, USA. Email address: houdre@math.gatech.edu

** Postal address: Université du Luxembourg, Unité de Recherche en Mathématiques, Maison du Nombre, 6 Avenue de la Fonte, L-4364 Esch-sur-Alzette, Grand Duché du Luxembourg. Email address: gkerchev@gmail.com

© The Author(s), 2022. Published by Cambridge University Press on behalf of Applied Probability Trust.

Very few results address the (weakly) dependent case, and in the present work we provide estimates on $d_K(f(X), \mathcal{N})$ and $d_W(f(X), \mathcal{N})$ when X is generated by a hidden Markov model. Such models are of interest because of their many applications in fields such as computational biology and speech recognition; see e.g. [7]. Recall that a hidden Markov model (Z, X) consists of a Markov chain $Z = (Z_1, \dots, Z_n)$ which emits the observed variables $X = (X_1, \dots, X_n)$. The possible states in Z are each associated with a distribution on the values of X . In other words the observation X is a mixture model where the choice of the mixture component for each observation depends on the component of the previous observation. The mixture components are given by the sequence Z . Note also that given Z , X is a Markov chain.

The content of the paper is as follows. Section 2 contains a short overview of results on normal approximation in the independent setting and introduces a simple transformation involving independent and identically distributed (i.i.d.) random variables that allows us to adapt these estimates to the hidden Markov model. Moreover, further quantitative bounds are provided for the special case when f is a Lipschitz function with respect to the Hamming distance. Applications to variants of the ones analyzed in [2] and [14] are developed in Sections 3 and 4, leading to an extra log-factor in the various rates obtained there. The majority of the more technical computations are carried out and presented in Section 5.

2. Main results

For a comprehensive review of Stein’s method we refer the reader to [4]. The exchangeable pairs approach was outlined in [2]. We now recall below a few of its main points.

Let $W := f(X)$. Originally in [2], and then in [14], various bounds on the distance between W and the normal distribution were obtained through a variant of Stein’s method. As is well known, Stein’s method is a way to obtain normal approximations based on the observation that the standard normal distribution \mathcal{N} is the only centered and unit-variance distribution that satisfies

$$\mathbb{E}[g'(\mathcal{N})] = \mathbb{E}[\mathcal{N}g(\mathcal{N})]$$

for all absolutely continuous g with almost-everywhere (a.e.) derivative g' such that $\mathbb{E}|g'(\mathcal{N})| < \infty$ (see [4]), and for the random variable W , $|\mathbb{E}[Wg(W) - g'(W)]|$ can be thought of as a distance measuring the proximity of W to \mathcal{N} . In particular, for the Kolmogorov distance, the solutions g_t to the differential equation

$$\mathbb{P}(W \leq t) - \mathbb{P}(\mathcal{N} \leq t) = g'_t(W) - Wg_t(W)$$

are absolutely continuous with a.e. derivative such that $\mathbb{E}|g'_t(\mathcal{N})| < \infty$ (see [4]). Then,

$$d_K(W, \mathcal{N}) = \sup_{t \in \mathbb{R}} |\mathbb{E}[g'_t(W) - Wg_t(W)]|. \tag{1}$$

Further properties of the solutions g_t allow for upper bounds on $\mathbb{E}[g'_t(W) - Wg_t(W)]$ using difference operators associated with W that were introduced in [2] (see [14]). This is called the *generalized perturbative approach* in [3], and we describe it next. First, we recall the perturbations used to bound the right-hand side of (1) in [2, 14]. Let $X' = (X'_1, \dots, X'_n)$ be an independent copy of X , and let $W' = f(X')$. Then (W, W') is an exchangeable pair, since it has the same joint distribution as (W', W) . A perturbation $W^A = f^A(X) := f(X^A)$ of W is defined through the change X^A of X as follows:

$$X_i^A = \begin{cases} X'_i & \text{if } i \in A, \\ X_i & \text{if } i \notin A, \end{cases}$$

for any $A \subseteq [n] := \{1, \dots, n\}$, including $A = \emptyset$. With these definitions, still following [2], difference operators are defined for any $\emptyset \subseteq A \subseteq [n]$ and $i \notin A$, as follows:

$$\Delta_i f^A = f(X^A) - f(X^{A \cup \{i\}}).$$

Moreover, set

$$T_A(f) := \sum_{j \notin A} \Delta_j f(X) \Delta_j f(X^A),$$

$$T'_A(f) := \sum_{j \notin A} \Delta_j f(X) |\Delta_j f(X^A)|,$$

and for $k_{n,A} = 1/\binom{n}{|A|}(n - |A|)$, set

$$T_n(f) := \sum_{\emptyset \subseteq A \subsetneq [n]} k_{n,A} T_A(f),$$

$$T'_n(f) := \sum_{\emptyset \subseteq A \subsetneq [n]} k_{n,A} T'_A(f).$$

Now, for $W = f(X_1, \dots, X_n)$ such that $\mathbb{E}[W] = 0$, $0 < \sigma^2 = \mathbb{E}[W^2] < \infty$, and assuming all the expectations below are finite, [2, Theorem 2.2] gives the bound

$$d_W(\sigma^{-1}W, \mathcal{N}) \leq \frac{1}{\sigma^2} \sqrt{\text{Var}(\mathbb{E}[T_n(f)|X])} + \frac{1}{2\sigma^3} \sum_{j=1}^n \mathbb{E}|\Delta_j f(X)|^3, \tag{2}$$

while [14, Theorem 4.2] yields

$$d_K(\sigma^{-1}W, \mathcal{N}) \leq \frac{1}{\sigma^2} \sqrt{\text{Var}(\mathbb{E}[T_n(f)|X])} + \frac{1}{\sigma^2} \sqrt{\text{Var}(\mathbb{E}[T'_n(f)|X])}$$

$$+ \frac{1}{4\sigma^3} \sum_{j=1}^n \sqrt{\mathbb{E}|\Delta_j f|^6} + \frac{\sqrt{2\pi}}{16\sigma^3} \sum_{j=1}^n \mathbb{E}|\Delta_j f(X)|^3, \tag{3}$$

where in both cases \mathcal{N} is now a standard normal random variable.

Our main abstract result generalizes (2) and (3) to the case when X is generated by a hidden Markov model. It is as follows.

Proposition 2.1. *Let (Z, X) be a hidden Markov model with Z an aperiodic time-homogeneous and irreducible Markov chain with finite state space \mathcal{S} , and X taking values in a non-empty finite \mathcal{A} . Let $W := f(X_1, \dots, X_n)$ with $\mathbb{E}[W] = 0$ and $0 < \sigma^2 = \mathbb{E}[W^2] < \infty$. Then there exist a finite sequence of independent random variables $R = (R_0, R_1, \dots, R_{|\mathcal{S}|(n-1)})$, with R_i taking values in $\mathcal{S} \times \mathcal{A}$ for $i = 0, \dots, |\mathcal{S}|(n-1)$, and a measurable function $h : (\mathcal{S} \times \mathcal{A})^{|\mathcal{S}|(n-1)+1} \rightarrow \mathbb{R}$ such that $h(R_0, \dots, R_{|\mathcal{S}|(n-1)})$ and $f(X_1, \dots, X_n)$ are identically distributed. Therefore,*

$$d_W(\sigma^{-1}W, \mathcal{N}) \leq \frac{1}{\sigma^2} \sqrt{\text{Var}(\mathbb{E}[T_{|R|}(h)|R])} + \frac{1}{2\sigma^3} \sum_{i=0}^{|\mathcal{S}|(n-1)} \mathbb{E}|\Delta_i h(R)|^3 \tag{4}$$

and

$$d_K(\sigma^{-1}W, \mathcal{N}) \leq \frac{1}{\sigma^2} \sqrt{\text{Var}(\mathbb{E}[T_{|R|}(h)|R])} + \frac{1}{\sigma^2} \sqrt{\text{Var}(\mathbb{E}[T'_{|R|}(h)|R])} + \frac{1}{4\sigma^3} \sum_{j=0}^{|R|-1} \sqrt{\mathbb{E}|\Delta_j h(R)|^6} + \frac{\sqrt{2\pi}}{16\sigma^3} \sum_{j=0}^{|R|-1} \mathbb{E}|\Delta_j h(R)|^3, \tag{5}$$

where \mathcal{N} is a standard normal random variable.

At first glance the above results might appear to be simple corollaries to (2) and (3). Indeed, as is well known, every Markov chain (in a Polish space) admits a representation via i.i.d. random variables U_1, \dots, U_n , uniformly distributed on $(0, 1)$, and the inverse distribution function. Therefore, $f(X_1, \dots, X_n) \stackrel{d}{=} h(U_1, \dots, U_n)$, for some function h , where, as usual, $\stackrel{d}{=}$ indicates equality in distribution. However, providing quantitative estimates for $\mathbb{E}|\Delta_j h(U_1, \dots, U_n)|$ via f seems to be out of reach, since passing from f to h involves the ‘unknown’ inverse distribution function. For this reason, we develop for our analysis a more amenable, although more restrictive, choice of i.i.d. random variables, which is described intuitively in the next paragraph and then again in greater detail in Section 2.1.

Consider $R = (R_0, \dots, R_{|\mathcal{S}|(n-1)})$ as stacks of independent random variables on the $|\mathcal{S}|$ possible states of the hidden chain that determine the next step in the process, with R_0 specifying the initial state. Each R_i takes values in $\mathcal{S} \times \mathcal{A}$ and is distributed according to the transition probability from the present hidden state. Then, one has $f(X_1, \dots, X_n) \stackrel{d}{=} h(R_0, \dots, R_{|\mathcal{S}|(n-1)})$, for $h = f \circ \gamma$, where the function γ translates between R and X . This construction is carried out in more detail in the next section. Further note that when $(X_i)_{i \geq 1}$ is a sequence of independent random variables, the hidden chain in the model consists of a single state, and then the function γ is the identity function and so $h = f$.

In order for Proposition 2.1 to be meaningful, a further quantitative study of the terms in the upper bounds is necessary. It turns out that the variance terms determine the order of decay; see Remark 5.2. Nevertheless, we obtain additional quantitative estimates for all the terms involved; the proofs are presented in Section 5.

Proposition 2.2. *With the notation as above, let f be Lipschitz with respect to the Hamming distance, i.e., let*

$$|f(x) - f(y)| \leq c \sum_{i=1}^n \mathbf{1}_{x_i \neq y_i}$$

for every $x, y \in \mathcal{A}^n$ and where $c > 0$. Then, for any $r > 0$,

$$\begin{aligned} \mathbb{E}|\Delta_i h(R)|^r &\leq C_1 (\ln n)^r, \\ \mathbb{E}|h(R) - \mathbb{E}[h(R)]|^r &\leq C_2 n^{r/2} (\ln n)^r, \end{aligned}$$

for n large enough and $C_1, C_2 > 0$ depending on r and the parameters of the model.

Let R' and R'' be independent copies of R , and let $\tilde{\mathbf{R}}$ be the random set of recombinations of R, R' , and R'' . The set $\tilde{\mathbf{R}}$ consists of $3^{|R|}$ random vectors of size $|R|$:

$$\tilde{\mathbf{R}} := \{Z = (Z_0, \dots, Z_{|R|-1}) : Z_i \in \{R_i, R'_i, R''_i\}\}.$$

Let

$$\begin{aligned}
 B_{|R|}(h) &:= \sup_{Y, Y', Z, Z' \in \tilde{\mathbf{R}}} \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \right], \\
 B_{|R|}^{(k)}(h) &:= \sup_{Y, Z, Z' \in \tilde{\mathbf{R}}} \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} \Delta_i h(Z)^2 \Delta_j h(Z')^2 \right], \\
 B_{|R|}^{(j)}(h) &:= \sup_{Y, Z, Z' \in \tilde{\mathbf{R}}} \mathbb{E} \left[\mathbf{1}_{\Delta_{i,k}h(Y) \neq 0} \Delta_i h(Z)^2 \Delta_k h(Z')^2 \right].
 \end{aligned}
 \tag{6}$$

Then the following general bound for the variance terms holds.

Proposition 2.3. *With the notation as above and for $U = T_{|R|}(h)$ or $U = T'_{|R|}(h)$, we have*

$$\begin{aligned}
 \sqrt{\text{Var}(\mathbb{E}[U|R])} &\leq \frac{1}{\sqrt{2}} \sum_{\emptyset \subseteq A \subsetneq [|R|]} k_{|R|,A} \left(\sum_{i=0}^{|R|-1} \sum_{j,k \notin A} \left(\mathbf{1}_{i=j=k} \mathbb{E} |\Delta_i h(R)|^4 + \mathbf{1}_{i \neq j \neq k} B_{|R|}(h) \right. \right. \\
 &\quad \left. \left. + (\mathbf{1}_{i \neq j=k} + \mathbf{1}_{i=k \neq j}) B_{|R|}^{(k)}(h) + (\mathbf{1}_{i \neq j=k} + \mathbf{1}_{i=j \neq k}) B_{|R|}^{(j)}(h) \right) \right)^{1/2},
 \end{aligned}$$

where $[|R|] := \{1, \dots, |R|\}$.

Note that in Proposition 2.3, the underlying function f is not assumed to be Lipschitz. Moreover, the function h is not symmetric, and therefore the expression above cannot be simplified further, in contrast to the similar results in [14]. The proofs of Propositions 2.2 and 2.3 are technical and are delayed to Sections 5.1 and 5.2 respectively.

2.1. Additional details on the construction of R

Let (Z, X) be a hidden Markov model with Z an aperiodic time-homogeneous and irreducible Markov chain on a finite state space \mathcal{S} , and with X taking values in an alphabet \mathcal{A} . Let P be the transition matrix of the hidden chain, and let Q be the $|\mathcal{S}| \times |\mathcal{A}|$ probability matrix for the observations; i.e., Q_{ij} is the probability of seeing output j if the latent chain is in state i . Let the initial distribution of the hidden chain be μ . Then

$$\begin{aligned}
 \mathbb{P}((Z_1, \dots, Z_n; X_1, \dots, X_n) = (z_1, \dots, z_n; x_1, \dots, x_n)) \\
 = \mu(z_1) Q_{z_1, x_1} P_{z_1, z_2} \dots P_{z_{n-1}, z_n} Q_{z_n, x_n}.
 \end{aligned}$$

Next we introduce a sequence of independent random variables $R_0, \dots, R_{|\mathcal{S}|(n-1)}$ taking values in $\mathcal{S} \times \mathcal{A}$ and a function γ such that $\gamma(R_0, \dots, R_{|\mathcal{S}|(n-1)}) = (Z_1, \dots, Z_n; X_1, \dots, X_n)$. For any $s, s' \in \mathcal{S}, x \in \mathcal{A}$, and $i \in \{0, \dots, n-1\}$, let

$$\begin{aligned}
 \mathbb{P}(R_0 = (s, x)) &= \mu(s) Q_{s,x}, \\
 \mathbb{P}(R_{i|\mathcal{S}|+s'} = (s, x)) &= P_{s',s} Q_{s,x}.
 \end{aligned}$$

The random variables R_i are well defined since $\sum_x Q_{s,x} = 1$ for any $s \in \mathcal{S}$ and $\sum_s P_{s',s} = \sum_s \mu(s) = 1$ for any $s' \in \mathcal{S}$. One can think of the variables R_i as a set of instructions indicating where the hidden Markov model goes next. The function γ reconstructs the realization

$(Z_i, X_i)_{i \geq 1}$ sequentially from the sequence $(R_i)_{i \geq 0}$. In particular, γ captures the following relations:

$$\begin{aligned} (Z_1, X_1) &= R_0, \\ (Z_{i+1}, X_{i+1}) &= R_{i|S|+s} \quad \text{if } Z_i = s \text{ for } i \geq 1. \end{aligned}$$

One can also think of the sequence $(R_i)_{i \geq 0}$ as $|S|$ stacks of random variables on the $|S|$ possible states of the latent Markov chain, and the values being rules for the next step in the model. Note that only one variable on the i th level of the stack will be used to determine the $(i + 1)$ th hidden and observed pair. Furthermore, the distribution of the random variables R_i for $i \geq 1$ encodes the transition and output probabilities in the P and Q matrices of the original model.

Thus one can write $f(X_1, \dots, X_n) = h(R_0, \dots, R_{|S|(n-1)})$, for $h := f \circ \gamma$, where the function γ does the translation from $(R_i)_{i \geq 0}$ to $(Z_i, X_i)_{i \geq 1}$ as described above.

Let $R' = (R'_0, \dots, R'_{|S|(n-1)})$ be an independent copy of R . Let $A \subseteq \{0, 1, \dots, |S|(n-1)\}$, and let the change R^A of R be defined as follows:

$$R_i^A = \begin{cases} R'_i & \text{if } i \in A, \\ R_i & \text{if } i \notin A, \end{cases} \tag{7}$$

where, as before, when $A = \{j\}$ we write R^j instead of $R^{\{j\}}$.

Recall that the ‘discrete derivative’ of h with a perturbation A is

$$\Delta_i h^A = h(R^A) - h(R^{A \cup \{i\}}).$$

Then (4) and (5) follow from (2) and (3), respectively, since when (Z, X) is a hidden Markov model one writes

$$W = f(X_1, \dots, X_n) \stackrel{d}{=} h(R_0, \dots, R_{|S|(n-1)}),$$

where the sequence $(R_i)_{i \geq 0}$ is a sequence of independent random variables.

Remark 2.1. (i) The idea of using stacks of independent random variables to represent a hidden Markov model is somewhat reminiscent of Wilson’s cycle popping algorithm for generating a random directed spanning tree; see [17]. The algorithm has also been related to loop-erased random walks in [9].

(ii) If S consists of a single state, making the hidden chain redundant, there is a single stack of instructions. This corresponds to the independent setting of [2] and [14], and then γ is just the identity function.

(iii) The same approach, via the use of instructions, is also applicable when \mathcal{A} and \mathcal{S} are infinite countable. The $Q_{s,x}$ no longer form a finite matrix, but the same definition holds as long as $\sum_{x \in \mathcal{A}} Q_{s,x} = 1$ for all $s \in \mathcal{S}$. We need a countably infinite number of independent instructions to encode $(Z_i, X_i)_{1 \leq i \leq n}$. In particular, let R_0 and $(R_{i,s})_{1 \leq i \leq n, s \in \mathcal{S}}$ be such that

$$\begin{aligned} \mathbb{P}(R_0 = (s, x)) &= \mu(s)Q_{s,x}, \\ \mathbb{P}(R_{i,s'} = (s, x)) &= P_{s',s}Q_{s,x}. \end{aligned}$$

Then the function γ reconstructs $(Z_i, X_i)_{1 \leq i \leq n}$ from R_0 and $(R_{i,s})_{1 \leq i \leq n, s \in \mathcal{S}}$ via

$$\begin{aligned} (Z_1, X_1) &= R_0, \\ (Z_{i+1}, X_{i+1}) &= R_{i,s} \quad \text{if } Z_i = s \text{ for } i \geq 1. \end{aligned}$$

(iv) It is possible to obtain bounds on the various terms involved in Proposition 2.2 and Proposition 2.3 in the case when $|\mathcal{S}|$ is a function of n . Assuming that there exists a general deterministic upper bound $g_r(n)$ on $\mathbb{E}|\Delta_i h|^r$, one can bound the non-variance terms in Proposition 2.2 by $C|\mathcal{S}|(\sqrt{g_6(n)} + g_3(n))/\sigma^3$. The variance terms, using Proposition 2.3, will then be bounded by $C\sqrt{|\mathcal{S}|g_4(n) + |\mathcal{S}|^3(A)/\sigma^2}$, where A is a complicated expression that depends on the particular problem as well as on $\mathbb{E}|\Delta_i h|^r$ and $|\mathcal{S}|$.

The next crucial part in the analysis is the key result we use everywhere: if a change in an instruction propagates X levels (a random variable), then $\mathbb{P}(X > K) \leq (1 - \epsilon)^K$, for some absolute $0 < \epsilon < 1$. If $|\mathcal{S}|$ is finite, this holds under some standard assumptions on the model. If $|\mathcal{S}|$ grows with n , then under some minor additional restrictions in the hidden Markov model, we have $\mathbb{P}(X > K) \leq (1 - 1/|\mathcal{S}|)^K$. Then K can be chosen to be at most a small power of n (we take $K = \ln n$ in the finite case, which explains the logarithmic factors in our bounds). For meaningful bounds, $|\mathcal{S}|$ could grow at most like $\ln n$. However, much more fine-tuning is necessary in the actual proof.

For some recent results (different from ours) on normal approximation for functions of general Markov chains we refer the reader to [5].

3. Covering process

Although our framework was initially motivated by [11] and the problem of finding a normal approximation result for the length of the longest common subsequences in dependent random words, some applications to stochastic geometry are presented below. Our methodology can be applied to other related settings, in particular to the variant of the occupancy problem introduced in the recent article [10] (see Remark 4.1).

Let (K, \mathcal{K}) be the space of compact subsets of \mathbb{R}^d , endowed with the hit-and-miss topology. Let E_n be a cube of volume n , and let C_1, \dots, C_n be random variables in E_n , called *germs*. In the i.i.d. setting of [14], each C_i is sampled uniformly and independently in E_n ; i.e., if $T \subseteq E_n$ with measure $|T|$, then

$$\mathbb{P}(C_i \in T) = \frac{|T|}{n},$$

for all $i \in \{1, \dots, n\}$.

Here, we consider C_1, \dots, C_n , generated by a hidden Markov model in the following way. Let Z_1, \dots, Z_n be an aperiodic irreducible Markov chain on a finite state space \mathcal{S} . Each $s \in \mathcal{S}$ is associated with a measure m_s on E_n . Then for each measurable $T \subseteq E_n$,

$$\mathbb{P}(C_i \in T | Z_i = s) = m_s(T).$$

Assume that there are constants $0 < c_m \leq c_M$ such that for any $s \in \mathcal{S}$ and measurable $T \subseteq E_n$,

$$\frac{c_m |T|}{n} \leq m_s(T) \leq \frac{c_M |T|}{n}.$$

Note that $c_m = c_M = 1$ recovers the setting of [14].

Let K_1, \dots, K_n be compact sets (*grains*) with $\text{Vol}(K_i) \in (V_1, V_2)$ (absolute constants) for $i = 1, \dots, n$. Let $X_i = C_i + K_i$ for $i = 1, \dots, n$ be the *germ-grain* process. Consider the closed set formed by the union of the germs translated by the grains

$$F_n = \left(\bigcup_{k=1}^n X_k \right) \cap E_n.$$

We are interested in the volume covered by F_n ,

$$f_V(X_1, \dots, X_n) = \text{Vol}(F_n),$$

and the number of isolated grains,

$$f_I(X_1, \dots, X_n) = \#\{k : X_k \cap X_j \cap E_n = \emptyset, k \neq j\}.$$

Theorem 3.1. *Let \mathcal{N} be a standard normal random variable. Then, for all $n \in \mathbb{N}$,*

$$d_K \left(\frac{f_V - \mathbb{E}f_V}{\sqrt{\text{Var}f_V}}, \mathcal{N} \right) \leq C \left(\frac{n(\ln n)^3}{\sqrt{\text{Var}(f_V)^3}} + \frac{n^{1/2}(\ln n)^4}{\text{Var}(f_V)} \right), \tag{8}$$

$$d_K \left(\frac{f_I - \mathbb{E}f_I}{\sqrt{\text{Var}f_I}}, \mathcal{N} \right) \leq C \left(\frac{n(\ln n)^3}{\sqrt{\text{Var}(f_I)^3}} + \frac{n^{1/2}(\ln n)^4}{\text{Var}(f_I)} \right), \tag{9}$$

for some constant $C > 0$ independent of n .

The study of the order of growth of $\text{Var}f_I$ and $\text{Var}f_V$ is not really part of the scope of the present paper. In the independent case, there are constants $0 < c_V \leq C_V$ such that $c_V n \leq \text{Var}f_V \leq C_V n$ and $c_V n \leq \text{Var}f_I \leq C_V n$ for n sufficiently large (see [13, Theorem 4.4]). In our dependent setting, a variance lower bound of order n will thus provide a rate of order $(\log n)^4 / \sqrt{n}$.

The proofs of the normal approximations for f_V and f_I are carried out in Sections 3.1 and 3.2. Many of the more technical computations are carried out in Section 5.

3.1. Normal approximation for f_V

Write $f_V(X_1, \dots, X_n) = h(R_0, \dots, R_{|\mathcal{S}|(n-1)})$ for a set of instructions R defined as in Section 2.1. The volume of each grain is bounded by V_2 , so f_V is Lipschitz with respect to the Hamming distance, with constant V_2 . Proposition 2.1 holds, and from Proposition 2.2, the non-variance terms in the bounds in Proposition 2.1 are bounded by $C(\ln n)^3 / \sqrt{n}$. Here and below, C is a constant, independent of n , which can vary from line to line. Indeed, for instance,

$$\begin{aligned} \frac{1}{4\sigma^3} \sum_{j=0}^{|R|-1} \sqrt{\mathbb{E}|\Delta_j h(R)|^6} &\leq C \text{Var}(f_V)^{-3/2} (|\mathcal{S}|(n-1) + 1) (\ln n)^3 \\ &\leq C n (\ln n)^3 / \text{Var}(f_V)^{3/2}. \end{aligned} \tag{10}$$

To analyze the bound on the variance terms given by Proposition 2.3, first note that

$$\sum_{i=0}^{|R|-1} \sum_{j, k \notin A} \mathbf{1}_{i=j=k} \mathbb{E}|\Delta_i h(R)|^4 \leq C n (\ln n)^4,$$

using Proposition 2.2. The other terms are bounded as follows.

Proposition 3.1. *Let $A \subsetneq [|R|]$, and let $B_{|R|}(h)$, $B_{|R|}^k(h)$, and $B_{|R|}^j(h)$ be as in (6). Then*

$$\sum_{i=0}^{|R|-1} \sum_{j,k \notin A} \mathbf{1}_{i \neq j \neq k} B_{|R|}(h) \leq Cn(\ln n)^8, \tag{11}$$

$$\sum_{i=0}^{|R|-1} \sum_{j,k \notin A} (\mathbf{1}_{i \neq j = k} + \mathbf{1}_{i = k \neq j}) B_{|R|}^{(k)}(h) \leq Cn(\ln n)^4, \tag{12}$$

$$\sum_{i=0}^{|R|-1} \sum_{j,k \notin A} (\mathbf{1}_{i \neq j = k} + \mathbf{1}_{i = j \neq k}) B_{|R|}^{(j)}(h) \leq Cn(\ln n)^4, \tag{13}$$

for some constant $C > 0$ that does not depend on n .

Proof. See Section 5.3 for the proof of the first bound. The others follow similarly. □

The bound on the variance terms in Proposition 2.3 becomes

$$\begin{aligned} \sqrt{\text{Var}(\mathbb{E}[U|R])} &\leq \frac{1}{\sqrt{2}} \sum_{A \subsetneq [|R|]} k_{|R|,A} (Cn(\ln n)^4 + Cn(\ln n)^8 + 2Cn(\ln n)^4)^{1/2} \\ &\leq C\sqrt{n}(\ln n)^4. \end{aligned} \tag{14}$$

Then (8) follows from (14), (10), and Proposition 2.1.

3.2. Normal approximation for f_I

The proof of (9) is more involved since the function f_I is not Lipschitz. Abusing notation, write $f_I(X_1, \dots, X_n) = h(R_0, \dots, R_{|\mathcal{S}|(n-1)})$ for a set of instructions R as in Section 2.1. Proposition 2.1 holds, and, as in our analysis for f_V , we proceed by estimating the non-variance terms in the bounds. The following holds.

Proposition 3.2. *For any $t = 1, 2, \dots$ and $i \in \{0, \dots, |\mathcal{S}|(n-1)\}$,*

$$\mathbb{E}|\Delta_i h|^t \leq C(\ln n)^t, \tag{15}$$

where $C = C(t) > 0$.

Proof. See Section 5.4. The approach is similar to the one employed for Proposition 2.2 and uses a graph representation. □

Therefore, for the non-variance term in Proposition 2.1, we have

$$\frac{1}{4\sigma^3} \sum_{j=0}^{|R|-1} \sqrt{\mathbb{E}|\Delta_j h(R)|^6} + \frac{\sqrt{2\pi}}{16\sigma^3} \sum_{j=0}^{|R|-1} \mathbb{E}|\Delta_j h(R)|^3 \leq Cn \left(\frac{\ln n}{\sqrt{\text{Var}(f_I)}} \right)^3. \tag{16}$$

We are left to analyze the bound on the variance terms given by Proposition 2.3. First, note that using Proposition 3.2,

$$\sum_{i=0}^{|R|-1} \sum_{j,k \notin A} \mathbf{1}_{i=j=k} \mathbb{E}|\Delta_i h(R)|^4 \leq Cn(\ln n)^4.$$

Proposition 3.3. *Let $A \subsetneq [|R|]$, and let $B_{|R|}(h)$, $B_{|R|}^k(h)$, and $B_{|R|}^j(h)$ be as in (6). Then the bounds (11), (12), and (13) hold in this setting as well.*

Proof. See Section 5.5. □

The bound on the variance terms in Proposition 2.3 becomes

$$\begin{aligned} \sqrt{\text{Var}(\mathbb{E}[U|R])} &\leq \frac{1}{\sqrt{2}} \sum_{A \subsetneq [|R|]} k_{|R|,A} \left(Cn(\ln n)^4 + Cn(\ln n)^8 + 2Cn(\ln n)^4 \right)^{1/2} \\ &\leq C\sqrt{n}(\ln n)^4. \end{aligned} \tag{17}$$

Then (9) follows from (17), (16), and Proposition 2.1.

4. Set approximation with random tessellations

Let $K \subseteq [0, 1]^d$ be compact, and let X be a finite collection of points in K . The Voronoi reconstruction, or the Voronoi approximation, of K based on X is given by

$$K^X := \{y \in \mathbb{R}^d : \text{the closest point to } y \text{ in } X \text{ lies in } K\}.$$

For $x \in [0, 1]^d$, denote by $V(x; X)$ the Voronoi cell with nucleus x within X , given by

$$V(x; X) := \{y \in [0, 1]^d : \|y - x\| \leq \|y - x'\| \text{ for any } x' \in (X, x)\},$$

where $(X, x) = X \cup \{x\}$, and where, as usual, $\|\cdot\|$ is the Euclidean norm in \mathbb{R}^d . The volume approximation of interest is

$$\varphi(X) := \text{Vol}(K^X) = \sum_i \mathbf{1}_{X_i \in K} \text{Vol}(V(X_i; X)).$$

In [14], $X = (X_1, \dots, X_n)$ is a vector of n i.i.d. random variables uniformly distributed on $[0, 1]^d$. Here, we consider X_1, \dots, X_n generated by a hidden Markov model in the following way. Let Z_1, \dots, Z_n be an aperiodic irreducible Markov chain on a finite state space \mathcal{S} . Each $s \in \mathcal{S}$ is associated with a measure m_s on $[0, 1]^d$. Then for each measurable $T \subseteq [0, 1]^d$,

$$\mathbb{P}(X_i \in T | Z_i = s) = m_s(T).$$

Assume, moreover, that there are constants $0 < c_m \leq c_M$ such that for any $s \in \mathcal{S}$ and measurable $T \subseteq [0, 1]^d$,

$$c_m \frac{|T|}{n} \leq m_s(T) \leq c_M \frac{|T|}{n}.$$

Recall the notions of Lebesgue boundary of K given by

$$\partial K := \{x \in [0, 1]^d : \text{Vol}(B(x, \epsilon) \cap K) > 0 \text{ and } \text{Vol}(B(x, \epsilon) \cap K^c) > 0, \text{ for any } \epsilon > 0\}$$

and

$$\partial K^r := \{x : d(x, \partial K) \leq r\}, \partial K^r_+ := K^c \cap \partial K^r,$$

where $d(x, A)$ is the Euclidean distance from $x \in \mathbb{R}^d$ to $A \subseteq \mathbb{R}^d$.

Now, for $\beta > 0$, let

$$\gamma(K, r, \beta) := \int_{\partial K_r^+} \left(\frac{\text{Vol}(B(x, \beta r) \cap K)}{r^d} \right)^2 dx.$$

Next, recall that K is said to satisfy the (weak) rolling ball condition if

$$\gamma(K, \beta) := \liminf_{r>0} \text{Vol}(\partial K^r)^{-1} (\gamma(K, r, \beta) + \gamma(K^c, r, \beta)) > 0. \quad (18)$$

Our main result is as follows.

Theorem 4.1. *Let $K \subseteq [0, 1]^d$ satisfy the rolling ball condition. Moreover, assume that there exist $S_-(K)$, $S_+(K)$, $\alpha > 0$ such that*

$$S_+(K)r^\alpha \leq \text{Vol}(\partial K^r) \leq S_+(K)r^\alpha \quad \text{for every } r > 0.$$

Then, for $n \geq 1$,

$$d_K \left(\frac{\varphi(X) - \mathbb{E}\varphi(X)}{\sqrt{\text{Var}(\varphi(X))}}, \mathcal{N} \right) \leq C \frac{(\ln n)^3}{n^{1/2-\alpha/d}}, \quad (19)$$

where $C > 0$ is a constant not depending on n .

As in [14], we split Theorem 4.1 into two results. The first one establishes a central limit theorem.

Proposition 4.1. *Let $0 < \sigma^2 = \text{Var}(\varphi(X))$. Assume that $\text{Vol}(\partial K^r) \leq S_+(K)r^\alpha$ for some $S_+(K)$, $\alpha > 0$. Then, for $n \geq 1$,*

$$d_K \left(\frac{\varphi(X) - \mathbb{E}\varphi(X)}{\sigma}, \mathcal{N} \right) \leq C \left(\frac{(\ln n)^2}{\sigma^2 n^{3/2+\alpha/2d}} + \frac{(\ln n)^3}{\sigma^3 n^{2+\alpha/2d}} \right), \quad (20)$$

where $C > 0$ is a constant not depending on n .

The second result introduces bounds on the variance under some additional assumptions.

Proposition 4.2. *Let $K \subseteq [0, 1]^d$ satisfy the rolling ball condition. Moreover, assume that there exist $S_-(K)$, $S_+(K)$, $\alpha > 0$ such that*

$$S_+(K)r^\alpha \leq \text{Vol}(\partial K^r) \leq S_+(K)r^\alpha \quad \text{for every } r > 0.$$

Then, for n sufficiently large,

$$C_d^- S_-(K) \gamma(K) \leq \frac{\text{Var}(\varphi(X))}{n^{-1-\alpha/d}} \leq C_d^+ S_+(K) \quad (21)$$

for some $C_d^-, C_d^+ > 0$.

It is clear that Theorem 4.1 will be proved once Proposition 4.1 and Proposition 4.2 are established.

4.1. Proof of Proposition 4.1

Again, as before, we introduce a set of instructions R and a function h such that $h(R) = \varphi(X)$. We apply Proposition 2.1, and the initial step is to bound $\mathbb{E}|\Delta_i h(R)|^r$, where $r > 0$. In fact, the following holds.

Proposition 4.3. *Under the assumptions of Proposition 4.1,*

$$\mathbb{E}|\Delta_i h(R)|^r \leq c_{d,r,\alpha} S_+(K) (\ln n)^r n^{-r-\alpha/d}, \tag{22}$$

where $c_{d,r,\alpha}$ depends on the parameters of the model and the dimension d , as well as on r and α . Moreover, for $n, q \geq 1$,

$$\mathbb{E}|\varphi(X) - \mathbb{E}\varphi(X)|^r \leq C_{d,r,\alpha} S_+(K) (\ln n)^r n^{-r/2-\alpha/d} \tag{23}$$

for some $C_{d,r,\alpha} > 0$.

Before presenting the proof, we introduce some notation. Recall that $x, y \in [0, 1]^d$ are said to be Voronoi neighbors within the set X if $V(x; X) \cap V(y; X) \neq \emptyset$. In general, the Voronoi distance $d_V(x, y; X)$ between x and y within X is given by the smallest $k \geq 1$ such that there exist $x = x_0, x_1 \in X, \dots, x_{k-1} \in X, x_k = y$ such that x_i, x_{i+1} are Voronoi neighbors for $i = 0, \dots, k - 1$.

Denote by $v(x, y; X) = \text{Vol}(V(y; X) \cap V(x; (y, X)))$ the volume that $V(y; X)$ loses when x is added to X . Then, for $x \notin X$,

$$\varphi(X, x) - \varphi(X) = \mathbf{1}_{x \in K} \sum_{y \in X \cap K^c} v(x, y; X) - \mathbf{1}_{x \in K^c} \sum_{y \in X \cap K} v(x, y; X).$$

Let $R_k(x; X)$ be the distance from x to the farthest point in the cell of a k th-order Voronoi neighbor in X ; i.e., for $X = (X_1, \dots, X_n)$,

$$R_k(x; X) = \sup \{ \|y - x\| : y \in V(X_i; X), d_V(x, X_i; X) \leq k \},$$

with $R(x; X) := R_1(x; X)$. If x does not have k th-order neighbors, take $R_k(x; X) = \sqrt{d}$. Then

$$\text{Vol}(V(x; X)) \leq \kappa_d R(x; X)^d,$$

where $\kappa_d = \pi^{d/2} / \Gamma(d/2 + 1)$ is the volume of the unit ball in \mathbb{R}^d .

Proof of Proposition 4.3. The proof relies on two results. First, we have the following technical lemma which is established in Section 5.6.

Lemma 4.1. *Assume there exist $S_+(K), \alpha > 0$ such that $\text{Vol}(\partial K^r) \leq S_+(K)r^\alpha$ for all $r > 0$. Let*

$$U_k(i) = \mathbf{1}_{d(X_i, \partial K) \leq R_k(X_i; X)} R_k(X_i; X)^d.$$

Then, for some $c_{d,qd+\alpha,k} > 0$,

$$\mathbb{E}U_k^q(i) \leq S_+(K)c_{d,qd+\alpha,k} n^{-q-\alpha/d}$$

for all $n \geq 1, q \geq 1$.

Second, within our framework, we have the following version of [14, Proposition 6.4] where $S(R)$ is the original set of points generated by R and $S(R^i)$ is the set of points generated after the change in the instruction R_i .

Proposition 4.4. (i) If, for every $s \in S(R) \setminus S(R^i)$, the set $R_1(s, S(R))$, which contains s and all its neighbors, is either entirely in K or entirely in K^c , then $\Delta_i h(R) = 0$. A similar result holds for $s \in S(R^i) \setminus S(R)$ and the set $R_1(s, S(R^i))$.

(ii) Assume $|i - j|$ is large enough so that $(S(R^i) \setminus S(R)) \cup (S(R^j) \setminus S(R)) = S(R^{ij}) \setminus S(R)$, where $S(R^{ij})$ is the set of points generated after the changes in both R_i and R_j . Suppose that for every $s_1 \in S(R^i) \Delta S(R)$ and $s_2 \in S(R^j) \Delta S(R)$, at least one of the following holds:

1. $d_V(s_1, s_2; S(R^{ij}) \cap S(R)) \geq 2$, or
2. $d_V(s_1, \partial K; S(R^{ij}) \cap S(R)) \geq 2$ and $d_V(s_2, \partial K; S(R^{ij}) \cap S(R)) \cap S(R) \geq 2$.

Then $\Delta_{i,j} h(R) = 0$.

Now, write

$$|\Delta_i h(R)| \leq \sum_{s \in S(R) \setminus S(R^i)} \mathbf{1}_{d_{S(R)}(s, \partial K) \leq R_1(s; S(R))} k_d R_1(s; S(R))^d + \sum_{s \in S(R^i) \setminus S(R)} \mathbf{1}_{d_{S(R^i)}(s, \partial K) \leq R_1(s; S(R^i))} k_d R_1(s; S(R^i))^d.$$

As before, for some $T > 0$, there exist an event E and $\epsilon > 0$ such that, conditioned on E , $|S(R^i) \setminus S(R)| = |S(R) \setminus S(R^i)| \leq T$ and $\mathbb{P}(E^c) \leq (1 - \epsilon)^T$. Then, from Lemma 4.1, there exist $S_+(K)$, $\alpha > 0$ such that

$$\mathbb{E}|\Delta_i h(R)|^r \leq c_{d,r,\alpha}(1 - \epsilon)^T + c_{d,r,\alpha} S_+(K) T^r n^{-r-\alpha/d},$$

where $c_{d,r,\alpha}$ depends on the parameters of the model and the dimension d , as well as on r and α . If $T = c \ln n$ for a suitable $c > 0$, then

$$\mathbb{E}|\Delta_i h(R)|^r \leq c_{d,r,\alpha} S_+(K) (\ln n)^r n^{-r-\alpha/d},$$

as desired. An application of the r th-moment Efron–Stein inequality (see [12, 16]) then yields (23). □

For the non-variance term in Theorem 2.1, we have

$$\frac{1}{4\sigma^3} \sum_{j=0}^{|R|-1} \sqrt{\mathbb{E}|\Delta_j h(R)|^6} + \frac{\sqrt{2\pi}}{16\sigma^3} \sum_{j=0}^{|R|-1} \mathbb{E}|\Delta_j h(R)|^3 \leq C\sigma^{-3} (\ln n)^3 n^{-2-\alpha/2d}. \tag{24}$$

To analyze the bound on the variance terms given by Proposition 2.3, first note that

$$\sum_{i=0}^{|R|-1} \sum_{j, k \notin A} \mathbf{1}_{i=j=k} \mathbb{E}|\Delta_i h(R)|^4 \leq C(\ln n)^4 n^{-3-\alpha/d},$$

again using Proposition 4.3. The other terms are bounded as follows.

Proposition 4.5. *Let $A \subsetneq [|R|]$, and let $B_{|R|}(h)$, $B_{|R|}^k(h)$ and $B_{|R|}^j(h)$ be as in (6). Then, for $\epsilon > 0$,*

$$\begin{aligned} \sum_{i=0}^{|R|-1} \sum_{j,k \notin A} \mathbf{1}_{i \neq j \neq k} B_{|R|}(h) &\leq C_\epsilon \left(n^{-3-2\alpha/d} (\ln n)^{10+4\epsilon} \right), \\ \sum_{i=0}^{|R|-1} \sum_{j,k \notin A} (\mathbf{1}_{i \neq j=k} + \mathbf{1}_{i=k \neq j}) B_{|R|}^{(k)}(h) &\leq C_\epsilon \left(n^{-3-2\alpha/d} (\ln n)^{10+4\epsilon} \right), \\ \sum_{i=0}^{|R|-1} \sum_{j,k \notin A} (\mathbf{1}_{i \neq j=k} + \mathbf{1}_{i=j \neq k}) B_{|R|}^{(j)}(h) &\leq C_\epsilon \left(n^{-3-2\alpha/d} (\ln n)^{10+4\epsilon} \right), \end{aligned}$$

for some constant $C_\epsilon > 0$ that does not depend on n .

Proof. See Section 5.7 for the proof of the first bound. The others follow similarly. □

The bounds on the variance terms in Proposition 2.3 become

$$\begin{aligned} \sqrt{\text{Var}(\mathbb{E}[U|R])} &\leq \frac{1}{\sqrt{2}} \sum_{A \subsetneq [|R|]} k_{|R|,A} \left(C \frac{(\ln n)^4}{n^{3+\alpha/d}} + C_\epsilon \frac{(\ln n)^{10+4\epsilon}}{n^{3+2\alpha/d}} \right)^{1/2} \\ &\leq C \frac{(\ln n)^2}{n^{3/2+\alpha/2d}}. \end{aligned} \tag{25}$$

Then (20) follows from (25), (24), and Proposition 2.1.

4.2. Proof of Proposition 4.2

The upper bound follows immediately from Proposition 4.3.

Recall the following result ([14, Corollary 2.4]) concerning the variance. Let $X := (X_1, \dots, X_n) \in E^n$, where E is a Polish space. If X' is an independent copy of X , and $f : E^n \rightarrow \mathbb{R}$ is measurable, with $\mathbb{E}[f(X)^2] < \infty$,

$$\text{Var}(f(X)) \geq \sum_{i=1}^n \mathbb{E} \left[\left(\mathbb{E}[\Delta_i f(X', X) | X] \right)^2 \right]. \tag{26}$$

In our setting we take $f = \varphi$. Unlike in [14], the function φ is not symmetric, and the right-hand side of (26) cannot be simplified. The lower bound provided by [14, Corollary 2.4] recovers the correct order of growth of the variance, which is enough for our purposes. However, more precise generic results are also available; see e.g. [1].

Let H be the realization of the hidden chain for X . By the law of total variance, $\text{Var}(\varphi(X)) \geq \text{Var}(\varphi(X)|H)$. Let X' be an independent copy of X , given H . Note that, given H , $(X_i)_{i=1, \dots, n}$ and $(X'_i)_{i=1, \dots, n}$ are independent random variables which are *not* identically distributed.

Applying (26) to $\varphi(X|H)$, we obtain

$$\text{Var}(\varphi(X)|H) \geq \sum_{i=1}^n \mathbb{E}_{X'_i}^H \left(\mathbb{E}_X^H [\varphi(X^i) - \varphi(X)] \right)^2,$$

where $X^i = (X_1, \dots, X_{i-1}, X'_i, X_{i+1}, \dots, X_n)$, and \mathbb{E}^H indicates that H is given. (To simplify notation in what follows we drop the H .) The difference from the proof in [14] is that now the variables are no longer identically distributed. Write

$$\mathbb{E}_X[\varphi(X^i) - \varphi(X)] = \mathbb{E}_X[\varphi(X^i) - \varphi(X^{(i)})] - \mathbb{E}_X[\varphi(X) - \varphi(X^{(i)})],$$

where $X^{(i)} = (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$. By Lemma 4.1,

$$\mathbb{E}_X[\varphi(X) - \varphi(X^{(i)})] \leq c_{d,\alpha} n^{-1-\alpha/d}. \tag{27}$$

We are left to study $\mathbb{E}[\varphi(X^i) - \varphi(X^{(i)})]$. Recall that

$$\begin{aligned} \varphi(X^i) - \varphi(X^{(i)}) &= \mathbf{1}_{\{X'_i \in K\}} \sum_{j \neq i} \mathbf{1}_{\{X_j \in K^C\}} v(X'_i, X_j; X^{(i,j)}) \\ &\quad - \mathbf{1}_{\{X'_i \in K^C\}} \sum_{j \neq i} \mathbf{1}_{\{X_j \in K\}} v(X'_i, X_j; X^{(i,j)}). \end{aligned}$$

Now, for the case $X'_i \in K^C$ (the other case being equivalent), we have

$$\begin{aligned} &\left| \mathbb{E}_{X, X'_i} \left[- \mathbf{1}_{\{X'_i \in K^C\}} \sum_{j \neq i} \mathbf{1}_{\{X_j \in K\}} v(X'_i, X_j; X^{(i,j)}) \right] \right| \\ &\geq \mathbb{E}_{X'_i} \left[\mathbf{1}_{\{X'_i \in \partial K_+^{n-1/d}\}} \sum_{j \neq i} \mathbb{E}_X \left[\mathbf{1}_{\{X_j \in K\}} v(X'_i, X_j; X^{(i,j)}) \right] \right], \end{aligned}$$

since $v(X'_i, X_j; X^{(i,j)}) \geq 0$. Then

$$\begin{aligned} &\mathbb{E}_X \left[\mathbf{1}_{\{X_j \in K\}} v(x, X_j; X^{(i,j)}) \right] \\ &\geq \mathbb{E}_{X^{(i,j)}} \left[c_1 \int_{y \in K} v(x, y; X^{(i,j)}) dy \right] \\ &\geq c_1 \text{Vol}(B(x, \beta n^{-1/d}) \cap K) \inf_{y: \|x-y\| \leq \beta n^{-1/d}} \mathbb{E}_{X^{(i,j)}} \left[v(x, y; X^{(i,j)}) \right], \end{aligned}$$

using the independence after conditioning on H and the properties of the model. We want to find an event that implies that $v(x, y; X^{(i,j)}) \geq cn^{-1}$. One instance is when no point of $X^{(i,j)}$ falls in $B(y, 6\beta n^{-1/d})$. Indeed, then $B(y, 3\beta n^{-1/d}) \subset V(y, X^{(i,j)})$. The distance between y and x is less than $\beta n^{-1/d}$, and so there is $z \in B(y, 3\beta n^{-1/d})$, namely $z = x + \beta n^{-1/d}(x - y)/\|x - y\|$, such that

$$B(z, \beta n^{-1/d}) \subset V(x, (X^{(i,j)}, y)) \subset B(y, 3\beta n^{-1/d}) \subset V(y; X^{(i,j)}).$$

Then $v(x, y; X^{(i,j)}) \geq \text{Vol}(B(z, \beta n^{-1/d})) = \kappa_d \beta^d n^{-1}$. Finally,

$$\begin{aligned} &\inf_{y: \|x-y\| \leq \beta n^{-1/d}} \mathbb{E}_{X^{(i,j)}} \left[v(x, y; X^{(i,j)}) \right] \\ &\geq \kappa_d \beta^d n^{-1} \mathbb{P}\left(X^{(i,j)} \cap B(y, 6\beta n^{-1/d}) = \emptyset\right) \\ &\geq \kappa_d \beta^d n^{-1} (1 - c_2 \beta^d n^{-1})^n \\ &\geq c_d \beta^n n^{-1}, \end{aligned}$$

for some $c_{d,\beta} > 0$ depending on the parameters of the model, the dimension d , and β . Then

$$\mathbb{E}_X \left[\mathbf{1}_{\{X_j \in K\}} v(x, X_j; X^{(i,j)}) \right] \geq c_{d,\beta} \text{Vol}(B(x, \beta n^{-1/d})) n^{-1}.$$

Therefore, by the very definition of $\gamma(K, r, \beta)$ and since the case $X'_i \in K$ is symmetric,

$$\begin{aligned} \mathbb{E}_{X'_i} \mathbb{E}_X \left[(\varphi(X^i) - \varphi(X^{(i)}))^2 \right] &\geq c_{d,\beta} \left(c_1 \int_{\partial K_+^{n-1/d}} \text{Vol}(B(x, \beta n^{-1/d}) \cap K)^2 dx \right. \\ &\quad \left. + c_1 \int_{\partial K_-^{n-1/d}} \text{Vol}(B(x, \beta n^{-1/d}) \cap K^c)^2 dx \right) \\ &= c_{d,\beta} (n^{-2} \gamma(K, n^{-1/d}, \beta) + n^{-2} \gamma(K^c, n^{-1/d}, \beta)). \end{aligned}$$

If the rolling ball condition (18) and the lower bound on $\partial K^{n-1/d}$ both hold, then

$$\mathbb{E}_{X'_i} \mathbb{E}_X \left[(\varphi(X^i) - \varphi(X^{(i)}))^2 \right] \geq c_{d,\beta} S_-(K) \gamma(K, \beta) n^{-2-\alpha/d},$$

which dominates the contribution (27) from $\mathbb{E}[\varphi(X) - \varphi(X^{(i)})]$. Therefore, finally,

$$\text{Var}(\varphi(X)) \geq c_{d,\beta}^- S_-(K) \gamma(K, \beta) n^{-1-\alpha/d},$$

as desired.

Remark 4.1. Let us expand a bit on another potential application of our generic framework, namely the occupancy problem as studied in [10]. To set up the notation, (Z_1, \dots, Z_n) is an aperiodic, irreducible, and time-homogeneous (hidden) Markov chain that transitions between different alphabets. To each alphabet is associated a distribution over the collection of all possible letters, giving rise to the observed letters (X_1, \dots, X_n) . We assume that the number of alphabets is finite but that the number of total letters is $\lfloor \alpha n \rfloor$, for some fixed $\alpha > 0$. One studies $W := f(X_1, \dots, X_n)$ —the number of letters that have not appeared among the X_1, \dots, X_n . Then an analysis as in the proof of Theorem 3.1 leads to the following:

$$d_K \left(\frac{W - \mathbb{E}W}{\sqrt{\text{Var}(W)}}, \mathcal{N} \right) \leq C \left(\frac{n(\ln n)^3}{\sqrt{\text{Var}(W)^3}} + \frac{n^{1/2}(\ln n)^4}{\text{Var}(W)} \right),$$

where $\text{Var}(W)$ is a function of n , \mathcal{N} is the standard normal distribution, and $C > 0$ is a constant depending on the parameters of the model but not on n . As mentioned at the beginning of the section, the study of the precise order of growth of the variance of W , in our dependent framework, is not within the scope of the current paper. For the i.i.d. case one can show (see e.g. [8]) that $\text{Var}(W) \sim (\alpha e^{-1/\alpha} - (1 + \alpha)e^{-2/\alpha})n$ as $n \rightarrow \infty$.

5. Technical results

5.1. Bounds on terms involving $\Delta_i h$

Recall the setting of Proposition 2.2 and Proposition 2.3. Let (Z, X) be a hidden Markov model and let the latent chain Z be irreducible and aperiodic, with finite state space \mathcal{S} . Assume that Z is started at the stationary distribution.

We start by establishing a technical result regarding Z .

First, note that there exist $K \geq 1$ and $\epsilon \in (0, 1)$ such that

$$\mathbb{P}(Z_n = s, Z_{n+K} = s') \geq \epsilon,$$

and thus,

$$\mathbb{P}(Z_{n+K} = s') \geq \epsilon, \quad \mathbb{P}(Z_{n+K} = s' | Z_n = s) \geq \epsilon \tag{28}$$

for all $n \geq 1$ and $s, s' \in \mathcal{S}$.

Lemma 5.1. *Let $K \geq 1$ and $\epsilon \in (0, 1)$ be as in (28), and let $(Z_i)_{i \geq 1}$ be an irreducible and aperiodic Markov chain with finite state space \mathcal{S} . Then*

$$\mathbb{P}(Z_{j+K} \neq s_1, Z_{j+2K} \neq s_2, \dots, Z_{j+tK} \neq s_t) \leq (1 - \epsilon)^t \tag{29}$$

for any $t \geq 1, j \geq 1$, and $(s_1, \dots, s_t) \in \mathcal{S}^t$.

Proof. We show (29) by induction. The case $t=1$ follows from (28). Next, for $(s_1, \dots, s_{t+1}) \in \mathcal{S}^{t+1}$,

$$\begin{aligned} & \mathbb{P}(Z_{j+K} \neq s_1, Z_{j+2K} \neq s_2, \dots, Z_{j+(t+1)K} \neq s_{t+1}) \\ &= \sum_{s'_1 \neq s_1, \dots, s'_{n+1} \neq s_{t+1}} \mathbb{P}(Z_{j+K} = s'_1, \dots, Z_{n+1} = s'_{t+1}) \\ &= \sum_{s'_1 \neq s_1, \dots, s'_{n+1} \neq s_{t+1}} \mathbb{P}(Z_{j+(t+1)K} = s'_{t+1} | Z_{j+K} = s'_1, \dots, Z_{j+tK} = s'_t) \\ & \quad \cdot \mathbb{P}(Z_1 = s'_1, \dots, Z_{j+tK} = s'_t) \\ &= \sum_{s'_1 \neq s_1, \dots, s'_{n+1} \neq s_{t+1}} \mathbb{P}(Z_{j+(t+1)K} = s'_{t+1} | Z_{j+tK} = s'_t) \mathbb{P}(Z_{j+K} = s'_1, \dots, Z_{j+tK} = s'_t) \\ &= \sum_{s'_1 \neq s_1, \dots, s'_t \neq s_t} \mathbb{P}(Z_{j+(t+1)K} \neq s_{t+1} | Z_{j+tK} = s'_t) \mathbb{P}(Z_{j+K} = s'_1, \dots, Z_{j+tK} = s'_t) \\ &\leq (1 - \epsilon) \sum_{s'_1 \neq s_1, \dots, s'_t \neq s_t} \mathbb{P}(Z_{j+K} = s'_1, \dots, Z_{j+tK} = s'_t) \\ &= (1 - \epsilon) \mathbb{P}(Z_{j+K} \neq s_1, \dots, Z_{j+tK} \neq s_t) \\ &\leq (1 - \epsilon)^{t+1}, \end{aligned}$$

where we have used the Markov property, (28), and finally the induction hypothesis. This suffices for the proof of (29), and thus the proof of the lemma is complete. \square

Let $f : \mathcal{A}^n \rightarrow \mathbb{R}$ be Lipschitz, i.e., be such that $|f(x) - f(y)| \leq c \sum_{i=1}^n \mathbf{1}_{x_i \neq y_i}$ for every $x, y \in \mathcal{A}^n$, where $c > 0$. Let $R = (R_0, \dots, R_{|\mathcal{S}|(n-1)})$ be a vector of independent random variables, and let h be the function such that

$$f(X_1, \dots, X_n) \stackrel{d}{=} h(R_0, \dots, R_{|\mathcal{S}|(n-1)}).$$

Let R' be an independent copy of R . The next result provides a tail inequality that is key for Proposition 2.2.

Proposition 5.1. *Let $K > 0$ and $\epsilon > 0$ be as in (28). Then*

$$\mathbb{P}(|\Delta_t h(R)| \geq cx) \leq C(1 - \epsilon)^{x/K} \tag{30}$$

for any $x \in \mathbb{N}$, where $C > 0$ depends on the parameters of the model but neither on n nor on x .

Proof. The sequence of instructions $R^i := (R_0, \dots, R'_i, \dots, R_{|\mathcal{S}(n-1)|})$ may give rise to a different realization (Z', X') of the hidden Markov model, as compared to (Z, X) , the one generated by R . The two models are not independent. In particular, if instruction R_i determines (Z_j, X_j) and R'_i determines (Z'_j, X'_j) , then $(Z_k, X_k) = (Z'_k, X'_k)$ for $k < j$. Let s be the smallest nonnegative integer (possibly $s = \infty$) such that $Z_{j+s} = Z'_{j+s}$. Then for any $k > j + s$, $(Z_k, X_k) = (Z'_k, X'_k)$ as well. Finally, if $k \in \{j, \dots, j + s - 1\}$, the pairs (Z_k, X_k) and (Z'_k, X'_k) are independent. We show next that for $K \geq 1$ as in (28), and any $t \in \mathbb{N}$,

$$\mathbb{P}(s \geq tK) \leq (1 - \epsilon)^t. \tag{31}$$

Indeed,

$$\begin{aligned} \mathbb{P}(s > tK) &\leq \mathbb{P}\left(Z_{j+K} \neq Z'_{j+K}, Z_{j+2K} \neq Z'_{j+2K}, \dots, Z_{j+tK} \neq Z'_{j+tK}\right) \\ &= \sum_{(s_1, \dots, s_t) \in \mathcal{S}^t} \mathbb{P}\left(Z_{j+K} \neq s_1, Z'_{j+K} = s_1, \dots, Z_{j+tK} \neq s_t, Z'_{j+tK} = s_t\right). \end{aligned}$$

By independence,

$$\begin{aligned} &\mathbb{P}\left(Z_{j+K} \neq s_1, Z'_{j+K} = s_1, \dots, Z_{j+tK} \neq s_t, Z'_{j+tK} = s_t\right) \\ &= \mathbb{P}\left(Z_{j+K} \neq s_1, \dots, Z_{j+tK} \neq s_t\right) \mathbb{P}\left(Z'_{j+K} = s_1, \dots, Z'_{j+tK} = s_t\right), \end{aligned}$$

and thus by Lemma 5.1

$$\begin{aligned} \mathbb{P}(s > tK) &\leq \sum_{(s_1, \dots, s_t)} (1 - \epsilon)^t \mathbb{P}\left(Z'_{j+K} = s_1, \dots, Z'_{j+tK} = s_t\right) \\ &\leq (1 - \epsilon)^t, \end{aligned}$$

as desired.

Let $E(t)$ be the event

$$E(t) := \left\{X_{j+K} \neq X'_{j+K}, X_{j+2K} \neq X'_{j+2K}, \dots, X_{j+tK} \neq X'_{j+tK}\right\}.$$

Note that $\mathbb{P}(E(t)) \leq \mathbb{P}(s \geq tK) \leq (1 - \epsilon)^t$. In particular, if $|h(R) - h(R^i)| \geq cx$, where $c > 0$ is the Lipschitz constant of the associated function f , then $s \geq x$, as there are at least x positions k such that $X_k \neq X'_k$. Thus,

$$\begin{aligned} \mathbb{P}\left(|h(R) - h(R^i)| \geq cx\right) &\leq \mathbb{P}(E(\lfloor x/K \rfloor)) \\ &\leq C(1 - \epsilon)^{x/K}, \end{aligned} \tag{32}$$

where $C > 0$ depends on the parameters of the model but not on x . This suffices for the proof of (30). □

We now turn to the proof of Proposition 2.2.

Proof of Proposition 2.2. Let E_t be the event that $|h(R) - h(R^i)| \geq tK$. Then

$$\mathbb{E}|h(R) - h(R^i)|^r = \mathbb{E}|h(R) - h(R^i)|^r \mathbf{1}_{E_t} + \mathbb{E}|h(R) - h(R^i)|^r \mathbf{1}_{E_t^c}.$$

Recall that $|g(x)| \leq cn$ for all $x \in \mathcal{A}^n$, and then $|h(R) - h(R^i)| \leq 2cn$. Using (32),

$$\begin{aligned} \mathbb{E}|h(R) - h(R^i)|^r &\leq (2cn)^r \mathbb{P}(E_t) + (ctK)^r \mathbb{P}(E_t^c) \\ &\leq (2cn)^r (1 - \epsilon)^t + (ctK)^r. \end{aligned} \tag{33}$$

Let $t = -r \ln n / (\ln(1 - \epsilon)) > 0$. Then

$$\mathbb{E}|h(R) - h(R^i)|^r \leq (2c)^r + \left(-\frac{crK}{\ln(1 - \epsilon)} \right)^r (\ln n)^r. \tag{34}$$

The order of the bound is optimal for t such that

$$(1 - \epsilon)^t \leq \left(\frac{\ln n}{n} \right)^r, \tag{35}$$

or

$$t \geq -\frac{r(\ln n - \ln(\ln n))}{\ln(1 - \epsilon)};$$

it follows that

$$\mathbb{E}|h(R) - h(R^i)|^r \leq (2c)^r + \left(-\frac{crK}{\ln(1 - \epsilon)} \right)^r (\ln n - \ln(\ln n))^r,$$

and the right-hand side has the same order of growth as (34).

If the growth order of $(1 - \epsilon)^t$ is larger than the one in (35), the bound on the second term in (33) is of larger order as well.

Then $\mathbb{E}|\Delta_i h(R)|^r \leq C_1 (\ln n)^r$, for $C_1 > 0$ depending on the parameters of the model and r . The first part of Proposition 2.2 is established.

For the upper bound on the central moments of $f(X)$, recall the following generalizations of the Efron–Stein inequality (see [12, 16]): for $r \geq 2$,

$$\left(\mathbb{E}|h(R) - \mathbb{E}h(R)|^r \right)^{1/r} \leq \frac{r-1}{2^{1/r}} \left(\sum_{i=0}^{|R|-1} \left(\mathbb{E}|h(R) - h(R^i)|^r \right)^{2/r} \right)^{1/2},$$

and for $r \in (0, 2)$,

$$\left(\mathbb{E}|h(R) - \mathbb{E}h(R)|^r \right)^{1/r} \leq \frac{1}{\sqrt{2}} \left(\sum_{i=0}^{|R|-1} \mathbb{E}|h(R) - h(R^i)|^2 \right)^{1/2}.$$

Then, for all $r > 0$,

$$\begin{aligned} \mathbb{E}|h(R) - \mathbb{E}h(R)|^r &\leq \left(\max \left\{ \frac{1}{\sqrt{2}}, \frac{r-1}{2^{1/r}} \right\} \right)^r \left((|\mathcal{S}|(n-1) + 1)C(\ln n)^2 \right)^{r/2} \\ &\leq C_2 n^{r/2} (\ln n)^r, \end{aligned}$$

where $C_2 > 0$ is a function of $|\mathcal{S}|$ and r . □

Remark 5.1. (i) Recall that in the independent setting, there is a single stack, or equivalently, the state space of the latent chain consists of a single element. Then for s as defined in the first paragraph of the above proof, $\mathbb{P}(s > 1) = 0$. Thus we can take $tK = 2$, and since $\mathbb{P}(E_t) \leq \mathbb{P}(s \geq tk) = 0$, (33) becomes

$$\mathbb{E}|h(R) - h(R^i)|^r \leq (2c)^r,$$

which recovers the independent case.

(ii) Note that the bound on the central moments also follows from using an exponential bounded difference inequality for Markov chains proved by Paulin [15]. This holds for the general case when X is a Markov chain (not necessarily time-homogeneous) taking values in a Polish space $\Lambda = \Lambda_1 \times \dots \times \Lambda_n$, with mixing time τ_{min} . Then, for any $t \geq 0$,

$$\mathbb{P}(|f(X) - \mathbb{E}[f(X)]| \geq t) \leq 2 \exp \left(\frac{-2t^2}{\|c^*\|^2 \tau_{min}} \right),$$

where f is such that

$$|f(x) - f(y)| \leq \sum_{i=1}^n c_i \mathbf{1}_{x_i \neq y_i},$$

for any $x, y \in \mathbb{R}^n$ and some $c^* = (c_1, \dots, c_n) \in \mathbb{R}^n$, and where $\|c^*\|^2 = \sum_{i=1}^n c_i^2$.

5.2. Proof of Proposition 2.3

In this section we no longer require the underlying function f to be Lipschitz.

Let $U := \sum_{\emptyset \subseteq A \subsetneq [R]} k_{|R|,A} U_A / 2$ for a general family of square-integrable random variables $U_A(R, R')$. From [2, Lemma 4.4],

$$\begin{aligned} \sqrt{\text{Var}(\mathbb{E}[U|R])} &\leq \frac{1}{2} \sum_{\emptyset \subseteq A \subsetneq [R]} \sqrt{\text{Var}(\mathbb{E}[U_A|R])} \\ &\leq \frac{1}{2} \sum_{\emptyset \subseteq A \subsetneq [R]} \sqrt{\mathbb{E}[\text{Var}(U_A|R)]}. \end{aligned}$$

As in [14], this inequality will be used for both $U_A = T_A(h)$ and $U_A = T'_A(h)$. A major difference from the setting in [14, Section 5] is that the function h is not symmetric; i.e., if σ is a permutation of $\{0, \dots, |\mathcal{S}|(n-1)\}$, it is not necessarily the case that $h(R_0, \dots, R_{|\mathcal{S}|(n-1)}) = h(R_{\sigma(0)}, \dots, R_{\sigma(|\mathcal{S}|(n-1))})$. Indeed, each variable in R is associated with a transition at a particular step and from a particular state. Fix $A \subsetneq [R]$ and let \tilde{R} be another independent copy of R . Introduce the substitution operator

$$\tilde{S}_i(R) = (R_0, \dots, \tilde{R}_i, \dots, R_{|R]}).$$

Recall that from the Efron–Stein inequality,

$$\text{Var}(U_A|R') \leq \frac{1}{2} \sum_{i=0}^{|R|-1} \mathbb{E}[(\tilde{\Delta}_i U_A(R))^2|R'],$$

where $\tilde{\Delta}_i U_A(R) = U_A(\tilde{S}_i(R)) - U_A(R)$.

Then,

$$\sqrt{\text{Var}(\mathbb{E}[U|R])} \leq \frac{1}{\sqrt{8}} \sum_{\emptyset \subseteq A \subsetneq [R]} k_{|R|,A} \sqrt{\sum_{i=0}^{|R|-1} \mathbb{E}[\tilde{\Delta}_i U_A]^2}. \tag{36}$$

Recall also that $U_A = \sum_{j \notin A} \Delta_j h(R) a(\Delta_j h(X^A))$, where the function a is either the identity, or $a(\cdot) = |\cdot|$. Then

$$\begin{aligned} \sum_{i=0}^{|R|-1} \mathbb{E}[\tilde{\Delta}_i U_A]^2 &= \sum_{i=0}^{|R|-1} \sum_{j,k \notin A} \mathbb{E}[|\tilde{\Delta}_i(\Delta_j h(R) a(\Delta_j h(R^A)))| \\ &\quad \times |\tilde{\Delta}_i(\Delta_k h(R) a(\Delta_k h(R^A)))|]. \end{aligned} \tag{37}$$

Fix $0 \leq i \leq |R| - 1$, and note that for $j \notin A$,

$$\begin{aligned} &\tilde{\Delta}_i(\Delta_j h(R) - a(\Delta_j h(R^A))) \\ &= \tilde{\Delta}_i(\Delta_j h(R)) a(\Delta_j h(R^A)) + \Delta_j h(\tilde{S}_i(R)) \tilde{\Delta}_i(a(\Delta_j h(R^A))). \end{aligned} \tag{38}$$

Then, using $|\tilde{\Delta}_i a(\cdot)| \leq |\tilde{\Delta}_i(\cdot)|$, the summands in (37) are bounded by

$$4 \sup_{Y, Y', Z, Z'} \mathbb{E}|\tilde{\Delta}_i(\Delta_j h(Y)) \Delta_j h(Y') \tilde{\Delta}_i(\Delta_k h(Z)) \Delta_k h(Z')|, \tag{39}$$

where Y, Y', Z, Z' are recombinations of R, R', \tilde{R} ; i.e., $Y_i \in \{R_i, R'_i, \tilde{R}_i\}$, for $i \in [0, |R| - 1]$.

Next, as in [14], we bound each type of summand appearing in (37).

If $i = j = k$, and using $\tilde{\Delta}_i(\Delta_i(\cdot)) = \Delta_i(\cdot)$, (39) is bounded by

$$4 \sup_{Y, Y', Z, Z'} \mathbb{E}|\Delta_i h(Y) \Delta_i h(Y') \Delta_i h(Z) \Delta_i h(Z')| \leq 4 \mathbb{E}|\Delta_i h(R)|^4.$$

If $i \neq j \neq k$, we can switch \tilde{R}_i and R'_i , and Y is still a recombination. Then (39) is equal to

$$\begin{aligned} &4 \sup_{Y, Y', Z, Z'} \mathbb{E}[\Delta_i(\Delta_j h(Y)) \Delta_j h(Y') \Delta_i(\Delta_k h(Z)) \Delta_k h(Z')] \\ &\leq 4 \sup_{Y, Y', Z, Z'} \mathbb{E}[\mathbf{1}_{\Delta_{i,j} h(Y) \neq 0} (|\Delta_j h(Y)| + |\Delta_j h(Y^i)|) |\Delta_j h(Y')| \\ &\quad \times \mathbf{1}_{\Delta_{i,k} h(Z) \neq 0} (|\Delta_k h(Z)| + |\Delta_k h(Z^i)|) |\Delta_k h(Z')|] \\ &\leq 16 \sup_{Y, Y', Z, Z'} \mathbb{E}[\mathbf{1}_{\Delta_{i,j} h(Y) \neq 0, \Delta_{j,k} h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2], \end{aligned} \tag{40}$$

where the last step follows from the Cauchy–Schwarz inequality.

If $i \neq j = k$, (39) is equal to

$$\begin{aligned}
 & 4 \sup_{Y, Y', Z, Z'} \mathbb{E} |\tilde{\Delta}_i(\Delta_j(h(Y))\Delta_j h(Y')) \tilde{\Delta}_i(\Delta_j(h(Z))\Delta_j h(Z'))| \\
 &= 4 \sup_{Y, Z} \mathbb{E} |\tilde{\Delta}_i(\Delta_j(h(Y))^2 \Delta_j h(Z)^2)| \\
 &= 4 \sup_{Y, Z} \mathbb{E} |\Delta_j(\Delta_i(h(Y))^2 \Delta_j h(Z)^2)| \\
 &\leq 16 \sup_{Y, Z, Z'} \mathbb{E} |\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} \Delta_i h(Z)^2 \Delta_j h(Z')^2|, \tag{41}
 \end{aligned}$$

where we have exchanged \tilde{R}_i and R'_i and used the Cauchy–Schwarz inequality as in (40).

Similarly, if $i = j \neq k$, the bound is

$$\begin{aligned}
 & 4 \sup_{Y, Y', Z, Z'} \mathbb{E} |\tilde{\Delta}_i(\Delta_i(h(Y))\Delta_i h(Y')) \tilde{\Delta}_i(\Delta_k(h(Z))\Delta_k h(Z'))| \\
 &= 4 \sup_{Y, Y', Z, Z'} \mathbb{E} |\Delta_i h(Y)\Delta_i h(Y') \Delta_i(\Delta_k(h(Z))\Delta_k h(Z'))| \\
 &= 4 \sup_{Y, Z, Z'} \mathbb{E} |\Delta_i h(Y)^2 \Delta_i(\Delta_k(h(Z))\Delta_k h(Z'))| \\
 &\leq 8 \sup_{Y, Z, Z'} \mathbb{E} |\mathbf{1}_{\Delta_{i,k}h(Y) \neq 0} \Delta_i h(Z)^2 \Delta_k h(Z')^2|. \tag{42}
 \end{aligned}$$

Finally, if $i = k \neq j$, the bound is by symmetry

$$\begin{aligned}
 & 4 \sup_{Y, Y', Z, Z'} \mathbb{E} |\tilde{\Delta}_i(\Delta_j(h(Y))\Delta_j h(Y')) \tilde{\Delta}_i(\Delta_i(h(Z))\Delta_i h(Z'))| \\
 &\leq 8 \sup_{Y, Z, Z'} \mathbb{E} |\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} \Delta_i h(Z)^2 \Delta_j h(Z')^2|. \tag{43}
 \end{aligned}$$

Combining (40), (41), (42), and (43) in (37), we finally get

$$\begin{aligned}
 & \sum_{i=0}^{|R|-1} \mathbb{E} [\tilde{\Delta}_i U_A]^2 \\
 &\leq 16 \sum_{i=0}^{|R|-1} \sum_{j, k \notin A} \left(\mathbf{1}_{i=j=k} \mathbb{E} |\Delta_i h(R)|^4 + \mathbf{1}_{i \neq j \neq k} B_{|R|}(h) \right. \\
 &\quad \left. + (\mathbf{1}_{i \neq j = k} + \mathbf{1}_{i = k \neq j}) B_{|R|}^{(k)}(h) + (\mathbf{1}_{i \neq j = k} + \mathbf{1}_{i = j \neq k}) B_{|R|}^{(j)}(h) \right),
 \end{aligned}$$

where

$$\begin{aligned}
 B_{|R|}(h) &:= \sup_{Y, Y', Z, Z'} \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{i,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \right], \\
 B_{|R|}^{(k)}(h) &:= \sup_{Y, Z, Z'} \mathbb{E} |\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} \Delta_i h(Z)^2 \Delta_j h(Z')^2|, \\
 B_{|R|}^{(j)}(h) &:= \sup_{Y, Z, Z'} \mathbb{E} |\mathbf{1}_{\Delta_{i,k}h(Y) \neq 0} \Delta_i h(Z)^2 \Delta_k h(Z')^2|.
 \end{aligned}$$

This suffices for the proof of Proposition 2.3.

Remark 5.2. As observed in [6], the terms involving $\Delta_i h(R)$ in (4) and (5) can be removed, leaving only the variance terms. Here is a different way to establish this fact for the particular framework we consider in Section 3. Recall that the expressions on the right-hand sides of (4) and (5) are bounds on terms of the form

$$\mathbb{E}|g'_t(W) - g'_t(W)T| + |\mathbb{E}[g_t(W)W - g'_t(W)T]|,$$

where $|g'_t| \leq 1$ and $|g_t(W)W - g'_t(W)| = |\mathbf{1}_{W \leq t} - \mathbb{P}(\mathcal{N} \leq t)| \leq 1$ (see [14] and [3]). First, note that

$$|g'_t(W) - g'_t(W)T| \geq |g'_t(W)T| - 1$$

and

$$1 \geq |g_t(W)W - g'_t(W)| \geq |g_t(W)W| - 1.$$

Then, by the triangle inequality and the above,

$$|g_t(W)W - g'_t(W)T| \leq |g_t(W)W| + |g'_t(W)T| \leq |g'_t(W) - g'_t(W)T| + 3.$$

Let $\mathbb{E}|g'_t(W) - g'_t(W)T| \leq f(n)$ for some function f , with $f(n) \rightarrow \infty$ and such that for $\sigma^2 = \sigma^2(n), f(n)/\sigma^2 \rightarrow 0$. Then

$$|\mathbb{E}[g_t(W)W - g'_t(W)T]|/\sigma^3 \leq Cf(n)/\sigma^2,$$

for some constant $C > 0$ that does not depend on n . Therefore, the asymptotic behavior of the bounds in (4) and (5) is given by the terms corresponding to $\mathbb{E}|g'_t(W) - g'_t(W)T|$, i.e., the terms involving the variance. This modification of the method is also valid in our framework and would ‘improve’ our results. However, it does not have a significant effect on the rates obtained in our applications in Section 3, and so we will not pursue it any further here.

5.3. Proof of Proposition 3.1

Recall that

$$B_{|R|}(h) := \sup_{Y, Y', Z, Z'} \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \right], \tag{44}$$

where the supremum is taken over recombinations of R and its independent copies R' and R'' .

Let E be the event that at least one of the perturbations of the instructions in (44) yields a difference in more than K points. By Proposition 5.1, there is $\epsilon > 0$ such that $\mathbb{P}(E) \leq (1 - \epsilon)^K$. Then, by the Lipschitz property of h ,

$$\begin{aligned} & \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \right] \\ &= \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \right] \\ & \quad + \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_{E^c} \right] \\ &\leq \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_{E^c} \right] + Cn^4(1 - \epsilon)^K \\ &\leq CK^4 \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} \mathbf{1}_{E^c} \right] + Cn^4(1 - \epsilon)^K. \end{aligned} \tag{45}$$

If $S(Y)$ is the set of points generated by the instructions Y and $S(Y^i)$ —the set of points generated by Y after the perturbation of Y_i —let

$$S_1 := S(Y)\Delta S(Y^i),$$

where Δ is the symmetric difference operator. Similarly, let

$$\begin{aligned} S_2 &:= S(Y)\Delta S(Y^j), \\ S_3 &:= S(Y')\Delta S((Y')^i), \\ S_4 &:= S(Y')\Delta S((Y')^j). \end{aligned}$$

Note that, conditioned on E^c , $|S_i| \leq 2K$ for $i = 1, 2, 3, 4$. Furthermore, if $s_1 \cap s_2 = \emptyset$ for all $(s_1, s_2) \in (S_1, S_2)$, then $\Delta_{i,j}h(Y) = 0$. Then

$$\mathbf{1}_{\Delta_{i,j}h(Y)} \leq \sum_{(s_1, s_2) \in (S_1, S_2)} \mathbf{1}_{s_1 \cap s_2 \neq \emptyset}.$$

This bound is meaningful if the sets S_1 and S_2 are disjoint sets of random variables. Conditioned on E^c , this is the case if $|i - j| \geq |S|K$. We introduce events E_1, E_2 , and E_3 corresponding to 0, 1, or 2 of the conditions $\{|i - j| \geq |S|K, |j - k| \geq |S|K\}$ holding, respectively. The events E_1, E_2 , and E_3 are deterministic. Then we have

$$\mathbb{E}[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} \mathbf{1}_{E^c}] = \mathbb{E}[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} \mathbf{1}_{E^c} (\mathbf{1}_{E_1} + \mathbf{1}_{E_2} + \mathbf{1}_{E_3})].$$

First we use the trivial bound $\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} \leq 1$ to get

$$\mathbb{E}[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} \mathbf{1}_{E^c} \mathbf{1}_{E_1}] \leq \mathbf{1}_{E_1}. \tag{46}$$

Then, for the term with $\mathbf{1}_{E_3}$,

$$\mathbb{E}[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} \mathbf{1}_{E^c} \mathbf{1}_{E_3}] \leq \mathbf{1}_{E_3} \mathbb{E} \left[\sum_{(s_1, s_2) \in (S_1, S_2)} \sum_{(s_3, s_4) \in (S_3, S_4)} \mathbf{1}_{s_1 \cap s_2 \neq \emptyset, s_3 \cap s_4 \neq \emptyset} \right].$$

To bound $\mathbb{E}[\mathbf{1}_{s_1 \cap s_2 \neq \emptyset, s_3 \cap s_4 \neq \emptyset}]$, we condition on s_2, s_3 , and the values of all hidden variables H . Then, since S_1 and S_4 are disjoint, we have independence:

$$\begin{aligned} \mathbb{E}[\mathbf{1}_{s_1 \cap s_2 \neq \emptyset, s_3 \cap s_4 \neq \emptyset}] &= \mathbb{E}[\mathbb{E}[\mathbf{1}_{s_1 \cap s_2 \neq \emptyset, s_3 \cap s_4 \neq \emptyset} | s_2, s_3, H]] \\ &= \mathbb{E}[\mathbb{E}[\mathbf{1}_{s_1 \cap s_2 \neq \emptyset} | s_2, s_3, H] \mathbb{E}[\mathbf{1}_{s_3 \cap s_4 \neq \emptyset} | s_2, s_3, H]]] \\ &\leq \left(\frac{c_M V_2}{n} \right)^2. \end{aligned}$$

Therefore,

$$\mathbb{E}[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} \mathbf{1}_{E^c} \mathbf{1}_{E_3}] \leq \mathbf{1}_{E_3} CK^4/n^2 \tag{47}$$

for some $C > 0$ independent of K and n , where we have used that $|S_i| \leq 2K$ for $i = 1, 2, 3, 4$.

Finally, for the term with E_2 , we may assume that $|i - j| \geq |\mathcal{S}|K$, since the case $|j - k| \geq |\mathcal{S}|K$ is identical. Write, using the trivial bound on $\mathbf{1}_{\Delta_{j,k}h(Y')} \neq 0$,

$$\mathbb{E}[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} \mathbf{1}_{E^c} \mathbf{1}_{E_2}] \leq \mathbf{1}_{E_3} \mathbb{E} \left[\sum_{(s_1, s_2) \in (S_1, S_2)} \mathbf{1}_{s_1 \cap s_2 \neq \emptyset} \right].$$

Next, as before,

$$\mathbb{E}[\mathbf{1}_{s_1 \cap s_2 \neq \emptyset}] = \mathbb{E}[\mathbb{E}[\mathbf{1}_{s_1 \cap s_2 \neq \emptyset} | s_2, H]] \leq \frac{c_M V_2}{n}.$$

Then

$$\mathbb{E}[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} \mathbf{1}_{E^c} \mathbf{1}_{E_2}] \leq \mathbf{1}_{E_2} c_M K^2 / n. \tag{48}$$

Combining (45), (46), (48), and (47), we get the following bound on (44):

$$B_{|R|}(h) \leq C \left(\mathbf{1}_{E_1} K^4 + \mathbf{1}_{E_2} K^6 / n + \mathbf{1}_{E_3} K^8 / n^2 + n^4 (1 - \epsilon)^K \right).$$

Then,

$$\begin{aligned} & \sum_{i=0}^{|R|-1} \sum_{j, k \notin A} \mathbf{1}_{i \neq j \neq k} B_{|R|}(h) \\ & \leq C \left(nK^5 + n^2 K^7 / n + n^3 K^8 / n^2 + n^7 (1 - \epsilon)^K \right) \\ & \leq Cn(\ln n)^8, \end{aligned}$$

when we choose $K = c \ln n$ for a suitable $c > 0$, independent of n .

5.4. Proof of Proposition 3.2

As in the proof of Proposition 5.1, the sequence of instructions R^i may give rise to a different realization (Z', X') . Indeed, if the instruction R_i determines (Z_j, X_j) and R'_i determines (Z'_j, X'_j) , it is possible that $(Z_j, X_j) \neq (Z'_j, X'_j)$. Let $s \geq 0$ be the smallest integer (possibly $s = \infty$) such that $Z_{j+s} = Z'_{j+s}$. Then, as in (31), there is $\epsilon > 0$ such that for $K \in \mathbb{N}$,

$$\mathbb{P}(s \geq K) \leq (1 - \epsilon)^K.$$

Fix K , and let E be the event corresponding to $\{s \geq K\}$. Using the trivial bound $|h(R)| \leq n$, and thus $|\Delta_i h(R)| \leq 2n$, we have

$$\begin{aligned} \mathbb{E}|\Delta_i h|^t &= \mathbb{E}[|\Delta_i h|^t \mathbf{1}_E] + \mathbb{E}[|\Delta_i h|^t \mathbf{1}_{E^c}] \\ &\leq (2n)^t (1 - \epsilon)^K + \mathbb{E}[|\Delta_i h|^t \mathbf{1}_{E^c}]. \end{aligned} \tag{49}$$

Let $S(R)$ be the set of points generated by the sequence of instructions R , and let $S(R^j)$ be the points generated by R after the perturbation of R_j . Set $S = S(R) \Delta S(R^j)$ for the symmetric difference and $S^c = S(R) \cap S(R^j)$. Note that E^c implies that $|\mathcal{S}| \leq 2K$. Furthermore,

$$|\Delta_i h| \leq \sum_{s \in S} \sum_{x \in S^c} \mathbf{1}_{s \cap x \neq \emptyset},$$

and

$$|\Delta_i h|^t \leq \sum_{(s_1, \dots, s_t) \in S^t} \sum_{(x_1, \dots, x_t) \in (S^c)^t} \prod_{j, \ell=1}^t \mathbf{1}_{s_j \cap x_\ell \neq \emptyset},$$

To estimate (49), we need to evaluate

$$\mathbb{E} \left[\prod_{j, \ell=1}^t \mathbf{1}_{s_j \cap x_\ell \neq \emptyset} \right],$$

and to do so we proceed as in [14] by studying the shape of the relations of $(s_j, x_\ell)_{j, \ell \in \{1, \dots, t\}}$.

Identify the set $(s_j, x_\ell)_{j, \ell \in \{1, \dots, t\}}$ with the edges of the graph G whose vertices correspond to $(s_j)_{j \in \{1, \dots, t\}}$ and $(x_\ell)_{\ell \in \{1, \dots, t\}}$. In particular, if $s_{j_1} = s_{j_2}$ for some $j_1 \neq j_2$, we identify them with the same point in the graph G . Conditioned on the realization of the hidden chain Z , we have independence. Then, if G is a tree, fix a root and condition recursively on vertices at different distances from the root. By the restrictions on the volume of the grain and the sampling distribution,

$$\mathbb{E} \left[\prod_{j, \ell=1}^t \mathbf{1}_{s_j \cap x_\ell \neq \emptyset} \middle| Z = z^n \right] \leq \left(\frac{c_M V_2}{n} \right)^{|E(G)|},$$

where $|E(G)|$ is the number of edges in the graph G . Furthermore,

$$\mathbb{E} \left[\prod_{j, \ell=1}^t \mathbf{1}_{s_j \cap x_\ell \neq \emptyset} \right] \leq \left(\frac{c_M V_2}{n} \right)^{|E(G)|}.$$

Note that the same result holds if G is a graph without cycles, i.e., a collection of disjoint trees. In general, G might have cycles. Let T be a subgraph of G that contains no cycles. Then

$$\prod_{j, \ell=1}^t \mathbf{1}_{s_j \cap x_\ell \neq \emptyset} \leq \prod_{e=(e_1, e_2) \in E(T)} \mathbf{1}_{e_1 \cap e_2 \neq \emptyset},$$

where the product on the right-hand side runs over the edges $e = (e_1, e_2)$ of the graph T with $e_1 \in S$ and $e_2 \in S^c$. Let $|S|$ be the number of distinct vertices in (s_1, \dots, s_t) , and similarly let $|x|$ be the number for (x_1, \dots, x_t) . The graph G is complete bipartite with $|S| + |x|$ vertices. We can find a subgraph T of G , also with $|S| + |x|$ vertices and no cycles. Then

$$\begin{aligned} \mathbb{E}[|\Delta_i h|^t \mathbf{1}_{E^c}] &\leq \mathbb{E} \left[\mathbf{1}_E^c \sum_{(s_1, \dots, s_t) \in S^t} \sum_{(x_1, \dots, x_t) \in (S^c)^t} \prod_{j, \ell=1}^t \mathbf{1}_{s_j \cap x_\ell \neq \emptyset} \right] \\ &= \mathbb{E} \left[\mathbf{1}_E^c \sum_{a, b=1}^t \sum_{\substack{(s_1, \dots, s_t) \in S^t \\ |S|=a}} \sum_{\substack{(x_1, \dots, x_t) \in (S^c)^t \\ |x|=b}} \prod_{j, \ell=1}^t \mathbf{1}_{s_j \cap x_\ell \neq \emptyset} \right] \\ &\leq \mathbb{E} \left[\mathbf{1}_E^c \sum_{a, b=1}^t C_t |S|^a |S^c|^b \left(\frac{c_M V_2}{n} \right)^{a+b-1} \right] \\ &\leq C_t K^r, \end{aligned}$$

where $C_t > 0$ is a constant depending on t , and where we have used that $|\mathcal{S}| \leq 2K$ and $|\mathcal{S}^c| \leq 2n$.

Let $K = \text{cln } n$ for a suitable $c > 0$; then (49) implies (15) as desired.

5.5. Proof of Proposition 3.3

As before, let E be the event that all perturbations of instructions in (44) propagate at most K levels. We have that $\mathbb{P}(E^c) \leq (1 - \epsilon)^K$ for some $\epsilon \in (0, 1)$. Using the trivial bound $|h(Y)| \leq n$,

$$\begin{aligned} & \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \right] \\ &= \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \right] \\ & \quad + \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_{E^c} \right] \\ & \leq \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \right] + 4n^4(1 - \epsilon)^K. \end{aligned} \tag{50}$$

Let $S(Y^i)$ be the set of points generated by the sequence of instructions Y after the perturbation of Y_i . Let S be the set of all points in the expectation above, and furthermore let

$$\begin{aligned} S_1 &:= S(Y)\Delta S(Y^i), & S_2 &:= S(Y)\Delta S(Y^j), \\ S_3 &:= S(Y')\Delta S((Y')^j), & S_4 &:= S(Y')\Delta S((Y')^k), \\ S_5 &:= S(Z)\Delta S(Z^i), & S_6 &:= S(Z')\Delta S(Z^k), \end{aligned}$$

where Δ is the symmetric difference operator. Conditioned on E , $|S_i| \leq 2K$ for $i = 1, \dots, 6$ and $|\mathcal{S}| \leq 10n$.

Conditioned on E , if $j - i \leq |\mathcal{S}|K$, the perturbation in i might be propagating past the position corresponding to the instruction j , leading to difficulties in the analysis of $\Delta_{i,j}h(Y)$. This is why we condition further on the events E_1, E_2, E_3 corresponding respectively to 0, 1, or 2 of the conditions $\{|i - j| \geq |\mathcal{S}|K, |j - k| \geq |\mathcal{S}|K\}$ holding true. Note that E_1, E_2 , and E_3 are deterministic.

If E_1 holds, we use the trivial bound $\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} \leq 1$, which leads to

$$\begin{aligned} & \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \mathbf{1}_{E_1} \right] \\ & \leq \mathbb{E} \left[|\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \mathbf{1}_{E_1} \right] \\ & \leq \mathbf{1}_{E_1} CK^4, \end{aligned} \tag{51}$$

using the Cauchy–Schwarz inequality.

Conditioned on E_3 , the sets $S_1, S_2 \cup S_3$, and S_4 are pairwise disjoint. Next, in similarity to an argument presented in [14], if $s_1 \cap s = \emptyset$ and $s_2 \cap s = \emptyset$ for all $(s_1, s_2, s) \in (S_1, S_2, S)$, then $\Delta_{i,j}h(Y) = 0$. Therefore,

$$\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} \leq \sum_{\substack{s_1 \in S_1 \\ s_2 \in S_2}} \sum_{s \in S} \mathbf{1}_{s_1 \cap s \neq \emptyset, s_2 \cap s \neq \emptyset},$$

and also

$$\mathbf{1}_{\Delta_{j,k}h(Y') \neq 0} \leq \sum_{\substack{s_3 \in S_3 \\ s_4 \in S_4}} \sum_{s \in S} \mathbf{1}_{s_3 \cap s \neq \emptyset, s_4 \cap s \neq \emptyset}.$$

Furthermore,

$$|\Delta_j h(Z)| \leq \sum_{s_5 \in S_5} \sum_{s \in S} \mathbf{1}_{s_5 \cap s \neq \emptyset}$$

and

$$|\Delta_k h(Z')| \leq \sum_{s_6 \in S_6} \sum_{s \in S} \mathbf{1}_{s_6 \cap s \neq \emptyset}.$$

Therefore,

$$\begin{aligned} & \mathbb{E}[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \mathbf{1}_{E_3}] \\ & \leq \mathbb{E} \left[\left(\sum_{\substack{(s_1, s_2, s_3, s_4) \in (S_1, S_2, S_3, S_4) \\ (s', s'') \in S^2}} \mathbf{1}_{s_1 \cap s' \neq \emptyset, s_2 \cap s'' \neq \emptyset, s_3 \cap s' \neq \emptyset, s_4 \cap s'' \neq \emptyset} \right) \right. \\ & \quad \cdot \left. \left(\sum_{s_5 \in S_5} \sum_{s \in S} \mathbf{1}_{s_5 \cap s \neq \emptyset} \right)^2 \left(\sum_{s_6 \in S_6} \sum_{s \in S} \mathbf{1}_{s_6 \cap s \neq \emptyset} \right)^2 \mathbf{1}_E \mathbf{1}_{E_3} \right] \\ & \leq \mathbb{E} \left[\sum_{\substack{(s_1, \dots, s_4) \in (S_1, \dots, S_4) \\ (s_5, \dots, s_8) \in S_6^4}} \sum_{\substack{(s', s'') \in S^2 \\ (s'_5, \dots, s'_8) \in S^4}} \mathbf{1}_{s_1 \cap s' \neq \emptyset, s_2 \cap s'' \neq \emptyset, s_3 \cap s' \neq \emptyset, s_4 \cap s'' \neq \emptyset} \prod_{a, b=5}^8 \mathbf{1}_{s_a \cap s'_b \neq \emptyset} \mathbf{1}_E \mathbf{1}_{E_3} \right], \end{aligned} \tag{52}$$

where $S_{56} = S_5 \cup S_6$ and $|S_{56}| \leq 4K$, conditioned on E .

To evaluate the summand expression we use the graph representation. Let E_ℓ be the event that there are ℓ distinct points among $s', s'', s'_5, \dots, s'_8$, different from s_1, \dots, s_8 . Note that $\ell \in [0, 6]$. Conditioned on E_ℓ , we can find a subgraph with no cycles and $\ell + 2$ edges, of the graph with edges $\{\{s_1, s'\}, \{s_2, s''\}, \{s_3, s'\}, \{s_4, s''\}\} \cup \{\{s_a, s'_b\} : a, b \in [5, 8]\}$. Indeed, note that there are at least 3 different points among s_1, \dots, s_4 . Next, if there are x points present among s', s'' and $\ell - x$ points among s'_5, \dots, s'_8 , we can find a subgraph with no cycles with at least $\ell - x$ edges among $\{\{s_a, s'_b\} : a, b \in [5, 8]\}$ and $x + 2$ edges among $\{\{s_1, s'\}, \{s_2, s''\}, \{s_3, s'\}, \{s_4, s''\}\}$.

Then, if we further condition on the values of the hidden variables H , we get, by independence,

$$\mathbb{E} \left[\mathbf{1}_{\substack{s_1 \cap s' \neq \emptyset, s_2 \cap s'' \neq \emptyset, \\ s_3 \cap s' \neq \emptyset, s_4 \cap s'' \neq \emptyset}} \prod_{a, b=5}^8 \mathbf{1}_{s_a \cap s'_b \neq \emptyset} \mathbf{1}_E \mathbf{1}_{E_3} \mathbf{1}_{E_\ell} \middle| H \right] \leq \mathbf{1}_{E_3} \left(\frac{c_M V_2}{n} \right)^{\ell+2}.$$

Then (52) is further bounded by

$$\mathbf{1}_{E_3} \sum_{\ell=0}^6 (4K)^8 \binom{6}{\ell} (10n)^\ell \left(\frac{c_M V_2}{n} \right)^{\ell+2} \leq \mathbf{1}_{E_3} C K^8 n^{-2}, \tag{53}$$

for some $C > 0$ independent of n and K .

Finally, assume that E_2 holds and that $|i - j| \geq |S|K$. The case $|j - k| \geq |S|K$ is identical. As above, using the trivial bound $\mathbf{1}_{\Delta_{j,k}h(Y') \neq 0} \leq 1$,

$$\begin{aligned} & \mathbb{E}[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0, \Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \mathbf{1}_{E_2}] \\ & \leq \mathbb{E}[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \mathbf{1}_{E_2}] \\ & \leq \mathbb{E} \left[\left(\sum_{\substack{(s_1, s_2) \in (S_1, S_2) \\ s' \in S}} \mathbf{1}_{s_1 \cap s' \neq \emptyset, s_2 \cap s' \neq \emptyset} \right) \left(\sum_{s_5 \in S_5} \sum_{s \in S} \mathbf{1}_{s_5 \cap s \neq \emptyset} \right)^2 \left(\sum_{s_6 \in S_6} \sum_{s \in S} \mathbf{1}_{s_6 \cap s \neq \emptyset} \right)^2 \mathbf{1}_E \mathbf{1}_{E_2} \right] \\ & \leq \mathbb{E} \left[\sum_{\substack{(s_1, s_2) \in (S_1, S_2) \\ (s_5, \dots, s_8) \in S_{56}^4}} \sum_{\substack{s' \in S \\ (s'_5, \dots, s'_8) \in S^4}} \mathbf{1}_{s_1 \cap s' \neq \emptyset, s_2 \cap s' \neq \emptyset} \prod_{a,b=5}^8 \mathbf{1}_{s_a \cap s'_b \neq \emptyset} \mathbf{1}_E \mathbf{1}_{E_2} \right]. \end{aligned} \tag{54}$$

If we condition on E_ℓ and the values of the hidden variables H , we get

$$\mathbb{E} \left[\mathbf{1}_{s_1 \cap s' \neq \emptyset, s_2 \cap s' \neq \emptyset} \prod_{a,b=5}^8 \mathbf{1}_{s_a \cap s'_b \neq \emptyset} \mathbf{1}_E \mathbf{1}_{E_2} \mathbf{1}_{E_\ell} | H \right] \leq \mathbf{1}_{E_2} \left(\frac{c_M V_2}{n} \right)^{\ell+1},$$

since in this case s_1 and s_2 are distinct and we can find a subgraph with $\ell + 1$ edges and no cycles.

Then (54) is bounded by

$$\mathbf{1}_{E_2} \sum_{\ell=0}^6 (4K)^6 \binom{6}{\ell} (10n)^\ell \left(\frac{c_M V_2}{n} \right)^{\ell+1} \leq \mathbf{1}_{E_2} C K^6 n^{-1}, \tag{55}$$

for some $C > 0$.

We get the following bound on $B_{|R|}(h)$ using (50), (51), (55), and (53):

$$B_{|R|}(h) \leq C \left(\mathbf{1}_{E_1} K^4 + \mathbf{1}_{E_2} K^6/n + \mathbf{1}_{E_3} K^8/n^2 + n^4(1 - \epsilon)^K \right).$$

Then,

$$\begin{aligned} & \sum_{i=0}^{|R|-1} \sum_{j,k \notin A} \mathbf{1}_{i \neq j \neq k} B_{|R|}(h) \\ & \leq C \left(nK^6 + n^2 K^7/n + n^3 K^8/n^2 + n^7(1 - \epsilon)^K \right) \\ & \leq Cn(\ln n)^8, \end{aligned}$$

where we have chosen $K = c \ln n$ for a suitable $c > 0$, independent of n .

5.6. Proof of Lemma 4.1

To simplify computations, we introduce the process X' defined as

$$X' = \bigcup_{m \in \mathbb{Z}^d} (X + m).$$

Unlike in the independent setting in [14], here the law of X' is invariant only under integer-valued translations. Note that, almost surely, X' has exactly n points in any cube $[t, t + 1]^d$, where $t \in \mathbb{R}$. Let $T_x = \{[y, y + 1]^d : y \in \mathbb{R}^d, x \in [y, y + 1]^d\}$. Define $\overline{R}_k(x; X)$ as

$$\overline{R}_k(x; X) := \sup_{T \in T_x} R_k(x; X' \cap T).$$

Note that if $x \in [0, 1]^d$, then $[0, 1]^d \in T_x$ and so $\overline{R}_k(x; X') \geq R_k(x; X)$. When the X_i are sampled independently and uniformly, as in [14], it is the case that $\overline{R}_k(x; X')$ does not depend on the position of x . However, in the hidden Markov model case we need to find a further bound on $\overline{R}_k(x; X')$.

For that purpose, consider the cube $K_0 := [-1/2, 1/2]^d$ of volume 1 centered at $\mathbf{0} \in \mathbb{R}^d$. Let B_A be the open ball of \mathbb{R}^d centered at $\mathbf{0}$ and of volume $A < 1$, to be chosen later. Next, let $\tilde{X} = (0, \tilde{X}_1, \dots, \tilde{X}_{n-1})$ be such that $\tilde{X}_i \in K_0$ for all $i = 1, \dots, n - 1$. Furthermore, for any Lebesgue-measurable $T \subseteq K_0$, set

$$\mathbb{P}(\tilde{X}_i \in T) = c_m |T \cap B_A| + c_M |T \cap B_A^c|$$

for all $i \in 1, \dots, n - 1$, where $|\cdot|$ now denotes the Lebesgue measure of the corresponding sets. If $A = (c_M - 1)/(c_M - c_m)$, then the above is a well-defined positive measure on K_0 . From the restrictions of the hidden Markov model, if $\tilde{R}_k = R_k(0; \tilde{X})$,

$$\overline{R}_k(x; X) \leq \tilde{R}_k.$$

Indeed, \tilde{R}_k represents the worst-case scenario where the remaining points of X are least likely to be distributed in the volume closest to x .

Then,

$$\mathbb{E}U_k^q(i) \leq \mathbb{E}_{X_i, \tilde{X}} \left[\mathbf{1}_{d(X_i; \partial K) \leq \tilde{R}_k} \tilde{R}_k^{qd} \right] \leq S_+(K) \mathbb{E}_{\tilde{X}} \left[\tilde{R}_k^{qd+\alpha} \right], \tag{56}$$

where we have used the upper bound on $Vol(\partial K^r)$.

To estimate $\mathbb{E}[\tilde{R}_k^{qd+\alpha}]$, note that if $\tilde{R}_k \geq r$, there will be an open ball of radius $r/2k$ in K_0 containing no points of \tilde{X} . Moreover, there will be $s_d \in (0, 1)$, depending only on the dimension d , such that every ball of radius $2k$ contains a cube of side length $s_d r/k$ of the form $[g - s_d r/2k, g + s_d r/2k]$, where $g \in (s_d r/k)\mathbb{Z}^d$. Then, if $s_d r/k < 1$,

$$\begin{aligned} \mathbb{P}(\tilde{R}_k \geq r) &\leq \mathbb{P}(\exists g \in (s_d r/k)\mathbb{Z}^d : \tilde{X} \cap [g - s_d r/2k, g + s_d r/2k] = \mathbf{0}) \\ &\leq \#\{g : g \in (s_d r/k)\mathbb{Z}^d \cap [-r, r]^d\} \mathbb{P}(\tilde{X} \cap [-s_d r/2k, s_d r/2k] = \mathbf{0}) \\ &\leq \frac{k^d}{(s_d)^d} (1 - c_m (s_d r/k)^d)^{n-1}. \end{aligned}$$

If, on the other hand, $s_d r/k \geq 1$, then $\tilde{X} \cap [g - s_d r/2k, g + s_d r/2k] = \tilde{X}$ and $\mathbb{P}(\tilde{R}_k \geq r) = 0$. Using $1 - x \leq e^{-x}$, for any $u > 0$ we have

$$\begin{aligned} \mathbb{E}[\tilde{R}(0, \tilde{X})^u] &= \int_0^\infty \mathbb{P}(\tilde{R}(0, \tilde{X}) \geq r^{1/u}) dr \\ &\leq c_{d,k} \int_0^\infty (1 - c_m (s_d r^{1/u}/k)^d)^{n-1} dr \\ &\leq c_{d,k} \int_0^\infty \exp(-c_m(n-1)(s_d r^{1/u}/k)^d) dr \\ &\leq c_{d,k,u} (n-1)^{u/d} \int_0^\infty \exp(-r^{d/u}) dr. \end{aligned}$$

Applying the above in (56) yields

$$\mathbb{E}U_k^q(i) \leq c_{d,k,qd+\alpha} S_+(K)n^{-q-\alpha/d},$$

where $c_{d,k,qd+\alpha} > 0$ depends only on the parameters of the transition probabilities of the hidden chain and on d, k , and $qd + \alpha$, but neither on n nor on i .

5.7. Proof of Proposition 4.5

We analyze

$$B_{|R|}(h) := \sup_{Y, Y', Z, Z'} \mathbb{E}[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} \mathbf{1}_{\Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2], \tag{57}$$

where as before the supremum is taken over recombinations Y, Y', Z, Z' of R, R', R'' . Let E be the event that all perturbations of the instructions in (57) propagate at most T levels. There is $\epsilon > 0$, depending only on the parameters of the models, such that $\mathbb{P}(E^c) \leq (1 - \epsilon)^T$.

As before, conditioned on E , if $|j - i| \leq |S|K$, the perturbation in i might be propagating past the position corresponding to the instruction j , leading to difficulties in the analysis of $\Delta_{i,j}h(Y)$. This is the reason for conditioning further on the events E_1, E_2, E_3 corresponding respectively to 0, 1, or 2 of the conditions $\{|i - j| \geq |S|K, |j - k| \geq |S|K\}$ holding. Note that E_1, E_2 , and E_3 are deterministic.

In this setting, we also study the event that all Voronoi cells are small. For that purpose, as in [14], we introduce the event $\Omega_n(X)$, given by

$$\Omega_n(X) := \left(\max_{1 \leq j \leq n} R(X_j; X) \leq n^{-1/d} \rho_n \right),$$

where $\rho_n = (\ln n)^{1/d+\epsilon'}$ for ϵ' sufficiently small. Then, after conditioning on the realization of the hidden chain, a proof as in [14, Lemma 6.8] leads to

$$n^\eta (1 - \mathbb{P}(\Omega_n(X))) \rightarrow 0 \tag{58}$$

as $n \rightarrow \infty$, for all $\eta > 0$.

We now estimate $B_{|R|}(h)$. Write

$$\begin{aligned} & \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} \mathbf{1}_{\Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \right] \\ &= \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} \mathbf{1}_{\Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_{E^c} \right] \\ & \quad + \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} \mathbf{1}_{\Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \mathbf{1}_{\Omega_n^c} \right] \\ & \quad + \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} \mathbf{1}_{\Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \mathbf{1}_{\Omega_n} \mathbf{1}_{E_1} \right] \\ & \quad + \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} \mathbf{1}_{\Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \mathbf{1}_{\Omega_n} \mathbf{1}_{E_2} \right] \\ & \quad + \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} \mathbf{1}_{\Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \mathbf{1}_{\Omega_n} \mathbf{1}_{E_3} \right]. \end{aligned} \tag{59}$$

Using $|\Delta_j h(Z)|, |\Delta_k h(Z')| \leq 1$, we get that the first two terms in (59) are bounded by $\mathbb{P}(E^c) + \mathbb{P}(\Omega_n^c)$. Next,

$$\begin{aligned} & \mathbb{E}\left[\mathbf{1}_{\Delta_{i,j}h(Y)\neq 0}\mathbf{1}_{\Delta_{j,k}h(Y')\neq 0}|\Delta_j h(Z)|^2|\Delta_k h(Z')|^2\mathbf{1}_E\mathbf{1}_{\Omega_n}\mathbf{1}_{E_1}\right] \\ & \leq \mathbf{1}_{E_1}\mathbb{E}\left[|\Delta_j h(Z)|^2|\Delta_k h(Z')|^2\mathbf{1}_E\mathbf{1}_{\Omega_n}\right] \\ & \leq C\mathbf{1}_{E_1}T^4n^{-4-2\alpha/d}\rho_n^{4d}, \end{aligned} \tag{60}$$

where we have used the Cauchy–Schwarz inequality.

Next, define as before

$$\begin{aligned} S_1 & := S(Y)\Delta S(Y^i), \quad S_2 := S(Y)\Delta S(Y^j), \\ S_3 & := S(Y')\Delta S((Y')^j), \quad S_4 := S(Y')\Delta S((Y')^k). \end{aligned}$$

Further, let $S_0 = S(Y) \cap S(Y^i) \cap S(Y^j)$ and $S'_0 = S(Y') \cap S((Y')^j) \cap S((Y')^k)$. By Proposition 4.4(ii), it follows that conditioned on Ω_n ,

$$\mathbf{1}_{\Delta_{i,j}h(Y)\neq 0} \leq \sum_{s_1 \in S_1, s_2 \in S_2} \mathbf{1}_{d_{S_0}(s_1, \partial K) \leq 2n^{-1/d}\rho_n} \mathbf{1}_{d_{S_0}(s_2, \partial K) \leq 2n^{-1/d}\rho_n} \mathbf{1}_{d_{S_0}(s_1, s_2) \leq 2n^{-1/d}\rho_n}.$$

Conditioned on E_3 , the sets $S_1, S_2 \cup S_3$, and S_4 are pairwise disjoint:

$$\begin{aligned} & \mathbb{E}\left[\mathbf{1}_{\Delta_{i,j}h(Y)\neq 0}\mathbf{1}_{\Delta_{j,k}h(Y')\neq 0}|\Delta_j h(Z)|^2|\Delta_k h(Z')|^2\mathbf{1}_E\mathbf{1}_{\Omega_n}\mathbf{1}_{E_3}\right] \\ & \leq C\mathbf{1}_{E_3}T^4n^{-4-2\alpha/d}\rho_n^{4d}\mathbb{E}\left[\mathbf{1}_{\Delta_{i,j}h(Y)\neq 0}\mathbf{1}_{\Delta_{j,k}h(Y')\neq 0}\mathbf{1}_E\mathbf{1}_{\Omega_n}\right]. \end{aligned}$$

By conditioning on the realization of all hidden chains H , we obtain

$$\begin{aligned} & \mathbb{E}\left[\mathbf{1}_{\Delta_{i,j}h(Y)\neq 0}\mathbf{1}_{\Delta_{j,k}h(Y')\neq 0}\mathbf{1}_E\mathbf{1}_{\Omega_n}\right] \\ & = \mathbb{E}\left[\mathbb{E}\left[\mathbf{1}_{\Delta_{i,j}h(Y)\neq 0}\mathbf{1}_{\Delta_{j,k}h(Y')\neq 0}\mathbf{1}_E\mathbf{1}_{\Omega_n} \mid H\right]\right] \\ & \leq \mathbb{E}\left[\mathbb{E}\left[\sum_{\substack{s_1 \in S_1, s_2 \in S_2 \\ s'_1 \in S_3, s'_2 \in S_4}} \mathbf{1}_{d_{S_0}(s'_1, \partial K) \leq 2n^{-1/d}\rho_n} \mathbf{1}_{d_{S_0}(s_1, s_2) \leq 2n^{-1/d}\rho_n} \mathbf{1}_{d_{S'_0}(s'_1, s'_2) \leq 2n^{-1/d}\rho_n} \mathbf{1}_E\mathbf{1}_{\Omega_n} \mid H\right]\right] \\ & \leq \mathbb{E}\mathbb{E}\left[\sum_{s_2 \in S_2, s'_1 \in S_1} \mathbf{1}_{d_{S_0}(s'_1, \partial K) \leq 2n^{-1/d}\rho_n} \mathbf{1}_E\mathbf{1}_{\Omega_n}\right. \\ & \quad \left.\mathbb{E}\left[\sum_{s_1 \in S_1, s'_2 \in S_4} \mathbf{1}_{d_{S_0}(s_1, s_2) \leq 2n^{-1/d}\rho_n} \mathbf{1}_{d_{S'_0}(s'_1, s'_2) \leq 2n^{-1/d}\rho_n} \mid s'_1, s_2\right] \mid H\right]. \end{aligned}$$

Now, conditioned on H, s'_1 , and s_2 , we have independence in the innermost expectation. Therefore, the above is bounded by

$$\mathbb{E}\left[\sum_{s_2 \in S_2, s'_1 \in S_1} \mathbf{1}_{d_{S_0}(s'_1, \partial K) \leq 2n^{-1/d}\rho_n} \mathbf{1}_E\mathbf{1}_{\Omega_n} 4T^2 2^d n^{-2} \rho_n^{2d}\right] \leq CT^4 n^{-2} \rho_n^{2d} n^{-\alpha/d} \rho_n^\alpha.$$

Then,

$$\begin{aligned} & \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} \mathbf{1}_{\Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \mathbf{1}_{\Omega_n} \mathbf{1}_{E_3} \right] \\ & \leq C \mathbf{1}_{E_3} T^8 n^{-6-3\alpha/d} \rho_n^{6d+\alpha}. \end{aligned} \quad (61)$$

Finally, for the event E_2 , assuming that $|i-j| \geq |S|K$, the other case being identical, we have

$$\begin{aligned} & \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} \mathbf{1}_{\Delta_{j,k}h(Y') \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \mathbf{1}_{\Omega_n} \mathbf{1}_{E_2} \right] \\ & \leq \mathbb{E} \left[\mathbf{1}_{\Delta_{i,j}h(Y) \neq 0} |\Delta_j h(Z)|^2 |\Delta_k h(Z')|^2 \mathbf{1}_E \mathbf{1}_{\Omega_n} \mathbf{1}_{E_2} \right] \\ & \leq C \mathbf{1}_{E_2} T^6 n^{-5-3\alpha/d} \rho_n^{5d+\alpha}. \end{aligned} \quad (62)$$

Using (59), (60), (62), and (61) leads to

$$\begin{aligned} B_{|R|}(h) & \leq C \left((1-\epsilon)^T + \mathbb{P}(\Omega_n^c) + \mathbf{1}_{E_1} T^4 n^{-4-2\alpha/d} \rho_n^{4d} \right. \\ & \quad \left. + \mathbf{1}_{E_2} T^6 n^{-5-3\alpha/d} \rho_n^{5d+\alpha} + \mathbf{1}_{E_3} T^8 n^{-6-3\alpha/d} \rho_n^{6d+\alpha} \right). \end{aligned}$$

Then,

$$\begin{aligned} & \sum_{i=0}^{|R|-1} \sum_{j,k \notin A} \mathbf{1}_{i \neq j \neq k} B_{|R|}(h) \\ & \leq C \left(n^3 (1-\epsilon)^T + n^3 \mathbb{P}(\Omega_n^c) + T^6 n^{-3-2\alpha/d} \rho_n^{4d} + T^7 n^{-3-3\alpha/d} \rho_n^{5d+\alpha} + T^8 n^{-3-3\alpha/d} \rho_n^{6d+\alpha} \right) \\ & \leq C (n^{-3-2\alpha/d} (\ln n)^{10+4\epsilon'}), \end{aligned}$$

where we have chosen $K = c \ln n$, for a suitable $c > 0$, independent of n , using also (58) and the definition of ρ_n .

Funding information

The first author's research is supported in part by grant no. 524678 from the Simons Foundation. The second author was partially supported by a TRIAD NSF grant (award 1740776) and the FNR grant APOGee at Luxembourg University (R-AGR-3585-10-C).

Competing interests

There were no competing interests to declare which arose during the preparation or publication process for this article.

References

- [1] BOUSQUET, O. AND HOUDRÉ, C. (2019). Iterated jackknives and two-sided variance inequalities. In *High Dimensional Probability VIII: The Oaxaca Volume* (Progress in Probability 74), Birkhäuser, Cham, pp. 33–40.
- [2] CHATTERJEE, S. (2008). A new method for normal approximation. *Ann. Prob.* **36**, 1584–1610.
- [3] CHATTERJEE, S. (2014). A short survey of Stein's method. In *Proceedings of the International Congress of Mathematicians: Seoul 2014*, Vol. IV, Kyung Moon Sa, Seoul, pp. 1–24.
- [4] CHEN, L. H. Y., GOLDSTEIN, L. AND SHAO, Q.-M. (2014). *Normal Approximation by Stein's Method*. Springer, Berlin, Heidelberg.

- [5] CHEN, P., SHAO, Q.-M. AND XU, L. (2020). *A universal probability approximation method: Markov process approach*. Preprint. Available at <https://arxiv.org/abs/2011.10985>.
- [6] CHU, D., SHAO, Q.-M. AND ZHANG, Z. (2019). Berry–Esseen bounds for functionals of independent random variables. Presented at the Symposium in Memory of Charles Stein (1920–2016). Available at <https://projecteuclid.org/journals/annals-of-probability/volume-47/issue-1/BerryEsseen-bounds-of-normal-and-nonnormal-approximation-for-unbounded-exchangeable/10.1214/18-AOP1255.full>.
- [7] DURBIN, R., EDDY, S., KROGH, A. AND MITCHISON, G. (1998). *Biological Sequence Analysis*. Cambridge University Press.
- [8] ENGLUND, G. (1981). A remainder term estimate for the normal approximation in classical occupancy. *Ann. Prob.* **9**, 684–692.
- [9] GORODEZKY, I. AND PAK, I. (2012). Generalized loop-erased random walks and approximate reachability. *Random Structures Algorithms* **44**, 201–223.
- [10] GRABCHAK, M., KELBERT, M. AND PARIS, Q. (2020). On the occupancy problem for a regime-switching model. *J. Appl. Prob.* **57**, 53–77.
- [11] HOUDRÉ, C. AND KERCHEV, G. (2019). On the rate of convergence for the length of the longest common subsequences in hidden Markov models. *J. Appl. Prob.* **56**, 558–573.
- [12] HOUDRÉ, C. AND MA, J. (2016). On the order of the central moments of the length of the longest common subsequences in random words. In *High Dimensional Probability VII: The Cargèse Volume* (Progress in Probability **71**), Birkhäuser, Cham, pp. 105–136.
- [13] KENDALL, W. S. AND MOLCHANOV, I. (eds) (2010). *New Perspectives in Stochastic Geometry*. Oxford University Press.
- [14] LACHIÈZE-REY, R. AND PECCATI, G. (2017). New Berry–Esseen bounds for functionals of binomial point process. *Ann. Appl. Prob.* **27**, 1992–2031.
- [15] PAULIN, D. (2015). Concentration inequalities for Markov chains by Marton couplings and spectral methods. *Electron. J. Prob.* **20**, 32 pp.
- [16] RHEE, W. AND TALAGRAND, M. (1986). Martingale inequalities and the jackknife estimate of the variance. *Statist. Prob. Lett.* **4**, 5–6.
- [17] WILSON, D. B. (1996). Generating random spanning trees more quickly than the cover time. In *STOC '96: Proceedings of the Twenty-Eighth Annual ACM Symposium on Theory of Computing*, Association for Computing Machinery, New York, pp. 296–303.