Mehdi Molkaraie, Member, IEEE and Hans-Andrea Loeliger, Fellow, IEEE

Abstract—The paper proposes Monte Carlo algorithms for the computation of the information rate of two-dimensional source/channel models. The focus of the paper is on binary-input channels with constraints on the allowed input configurations. The problem of numerically computing the information rate, and even the noiseless capacity, of such channels has so far remained largely unsolved. Both problems can be reduced to computing a Monte Carlo estimate of a partition function. The proposed algorithms use tree-based Gibbs sampling and multilayer (multitemperature) importance sampling. The viability of the proposed algorithms is demonstrated by simulation results.

Index Terms—Two-dimensional channels, constrained channels, partition function, Gibbs sampling, importance sampling, factor graphs, sum-product message passing, capacity, information rate.

I. INTRODUCTION

Numerically computing the Shannon information rate for source/channel models with memory can be difficult. In many cases of practical interest, analytical results are not available or hard to evaluate numerically. For a large class of channels, however, Monte Carlo methods as proposed in [1]–[3] have been shown to yield reliable numerical results.

In this paper, we consider the extension of such Monte Carlo methods to source/channel models with a two-dimensional (2-D) structure. The focus of the paper is on 2-D binary-input channels with constraints on the allowed input configurations; for example, we consider the channel where adjacent channel input symbols must not both have the value 1. Variations of such channels are of interest in magnetic and optical storage, where the constraints are imposed, e.g., in order to reduce the intersymbol interference or to help in timing control [4]–[8]. We will consider both noiseless and noisy versions of such channels. With suitable modifications (simplifications), the methods of this paper can also be applied to other 2-D source/channel models such as channels with intersymbol interference.

In the one-dimensional (1-D) case, computing the capacity of noiseless constrained channels was addressed and solved by Shannon [9], see also [4]. For the noisy case, the Monte Carlo

Preliminary versions of the material of this paper were presented in [16]-[18]. methods of [1]–[3] can be used to compute the information rate. The 2-D case is harder. Even the noiseless capacity is hard to compute numerically: while very tight analytical results are available for a number of special cases (e.g., [10]–[15]), other cases have remained open problems. The noisy case has remained largely unsolved.

1

The capacity of a noiseless constrained channel is essentially the logarithm of the partition function of the corresponding indicator function (see Section II). Moreover, computing the information rate of noisy source/channel models can also be reduced to the computation of a partition function (see Section VI). The heart of the paper, therefore, are Monte Carlo algorithms for the computation of partition functions. Several such algorithms are well known [19]–[21], see also [22], [23], but some modifications will be necessary for the problems of interest in this paper. In particular, we will find tree-based Gibbs sampling (due to Hamze and de Freitas [24]) extremely useful. We will observe that Monte Carlo estimates of a partition function may actually be obtained as a by-product of tree-based Gibbs sampling, which does not seem to have been noticed before.

In related prior work, Monte Carlo algorithms have been used to compute bounds on, or approximations of, the information rate of 2-D source/channels with memory [25], [26]. Some of this work uses generalized belief propagation [27], which appears to yield very good approximations to the information rate [25], [18], [28], [29].

In contrast to all this prior work, the Monte Carlo methods of this paper are asymptotically unbiased, i.e., in the limit of infinitely many samples, the estimates are guaranteed to converge to the desired quantity (the information rate). Moreover, the focus of this paper is on constrained channels, for which these computational problems are harder than for intersymbol interference channels (cf. Section VI-B).

The empirical success of the proposed algorithms is epitomized by Fig. 8, which shows the uniform-input information rate of a binary-input channel with input constraints and additive white Gaussian noise (AWGN). As far as known to the authors, no such figure (for such a channel) has been presented before.

If the reader is not familiar with Gibbs sampling, the following comments on the general nature of this work may be in order. First, Gibbs sampling is easily proved (under very mild conditions) to yield samples that are *eventually* distributed according to the desired distribution and *asymptotically* independent [23] (i.e., deleting the first t samples

Mehdi Molkaraie was with the Dept. of Information Technology and Electrical Engineering, ETH Zürich, CH-8092 Zürich, Switzerland. He is now with the Dept. of Statistics and Actuarial Science, University of Waterloo, Waterloo N2L 3G1, Canada. Email: mmolkaraie@uwaterloo.ca. Hans-Andrea Loeliger is with the Dept. of Information Technology and Electrical Engineering, ETH Zürich, CH-8092 Zürich, Switzerland. Email: loeliger@isi.ee.ethz.ch.

and decimating the remaining sample sequence by a factor m results in an i.i.d. sequence in the limit $t, m \to \infty$). However, the dependencies among the samples may decay extremely slowly, which is the pivotal issue with Gibbs sampling and makes naive Gibbs sampling perfectly useless for the problems of this paper (and for many other problems). The challenge, therefore, is to speed up Gibbs sampling (i.e., to decrease the dependencies of the samples) by various additional tricks and insights so that it becomes useful.

Second, the class of problems for which the methods proposed in this paper will work is not easily expressed in exact terms. Again, the issue is not formal applicability (which is quite sweeping), but the speed of convergence, which strongly depends on the particulars of the case and is not easily predicted.

The paper is organized as follows. In Section II, we review partition functions and noiseless 2-D constrained channels, and we introduce the corresponding notation. In Section III, we review several Monte Carlo algorithms that will be used in this paper. However, additional ideas are necessary to make these algorithms work for our applications. In particular, we will need tree-based Gibbs sampling as described in Section IV. The application to noiseless constrained channels is demonstrated in Section V. The application to noisy channels is described and demonstrated in Section VI. The appendix reviews sampling from Markov chains, which is needed in Section IV.

II. PARTITION FUNCTION OF 2-D GRAPHICAL MODELS

Let $\mathcal{X}_1, \mathcal{X}_2, \ldots, \mathcal{X}_N$ be finite sets, let \mathcal{X} be the Cartesian product $\mathcal{X} \stackrel{\scriptscriptstyle \triangle}{=} \mathcal{X}_1 \times \mathcal{X}_2 \times \ldots \times \mathcal{X}_N$, and let f be a nonnegative function $f : \mathcal{X} \to \mathbb{R}$. We are interested in computing (exactly or approximately) the *partition function*

$$Z_f \stackrel{\scriptscriptstyle \Delta}{=} \sum_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) \tag{1}$$

for cases where

• N is large and

• f has a suitable factorization (as will be detailed below). We will usually assume $Z_f \neq 0$. Then

$$p_f(\mathbf{x}) \triangleq \frac{f(\mathbf{x})}{Z_f}$$
 (2)

is a probability mass function on \mathcal{X} . We also define the support of f (and of p_f) as

$$\mathcal{S}_f \stackrel{\triangle}{=} \{ \mathbf{x} \in \mathcal{X} : f(\mathbf{x}) > 0 \}.$$
(3)

If $f(\mathbf{x})$ has a cycle-free factor graph representation (and if $|\mathcal{X}_1|, |\mathcal{X}_2|, \ldots, |\mathcal{X}_N|$ are not too large), then Z_f can be computed efficiently by sum-product message passing [30], [31]. In this paper, however, we consider factor graphs with cycles. In particular, we are interested in examples of the following type.

Example: Simple 2-D Constrained Channel

Consider a grid of $N = M \times M$ binary (i.e., $\{0,1\}$ -valued) variables x_1, \ldots, x_N with the constraint that no two

Fig. 1. Forney factor graph of the indicator function (4). The unlabeled

boxes represent factors as in (5).

(horizontally or vertically) adjacent variables have both the value 1. Let $f : \{0, 1\}^N \to \{0, 1\}$ be the indicator function of this constraint, which can be factored into

$$f(x_1, \dots, x_N) = \prod_{k, \, \ell \text{ adjacent}} \kappa(x_k, x_\ell), \tag{4}$$

where the product runs over all adjacent pairs (k, ℓ) and with factors

$$\kappa(x_k, x_\ell) = \begin{cases} 0, & \text{if } x_k = x_\ell = 1\\ 1, & \text{otherwise.} \end{cases}$$
(5)

The corresponding Forney factor graph of f is shown in Fig. 1, where the boxes labeled "=" are equality constraints [31]. (Note that, in Forney factor graphs, variables are represented by edges. Fig. 1 may also be viewed as a factor graph as in [30] where the boxes labeled "=" are the variable nodes.)

Note that, in this example, $Z_f = |\mathcal{S}_f|$.

This example is known as the 2-D $(1, \infty)$ run-length limited constrained channel [4]. The quantity

$$C_M \stackrel{\triangle}{=} \frac{1}{N} \log_2 Z_f \tag{6}$$

is known as the (finite-size) noiseless capacity of the channel.

For this particular example, upper and lower bounds on the infinite-size noiseless capacity

$$C_{\infty} \stackrel{\triangle}{=} \lim_{M \to \infty} C_M \tag{7}$$

were first proposed in [10] and improved in [11] and [32]. The bounds in [32] agree on nine decimal digits, which far exceeds the accuracy that can be achieved with the Monte Carlo methods of the present paper. However, the methods proposed in this paper work also for various generalizations of this example for which no tight bounds are known. \Box

Later on, in Section VI, we will consider noisy versions of such channels. As it turns out, the computation of the information rates of such channels also requires the computation of partition functions as in (1).



III. MONTE CARLO METHODS FOR PARTITION FUNCTION ESTIMATION

One well-known method to estimate $\Gamma_f \stackrel{\triangle}{=} 1/Z_f$ (and thus Z_f itself) goes as follows.

Ogata-Tanemura Method [19], [21]:

- 1) Draw samples $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(K)}$ from \mathcal{S}_f according to $p_f(\mathbf{x})$ as in (2).
- 2) Compute

$$\hat{\Gamma}_f = \frac{1}{K|\mathcal{S}_f|} \sum_{k=1}^K \frac{1}{f(\mathbf{x}^{(k)})} \tag{8}$$

It is easy to verify that $E(\hat{\Gamma}_f) = 1/Z_f$.

However, there are two major issues with this method. First, there is the problem of generating the samples $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(K)}$ according to $p_f(\mathbf{x})$. Ideally, we would like these samples to be independent, but (for the purposes of this paper) this is asking too much. In particular, we will use Gibbs sampling [22], [33], which produces dependent samples. However, with naive Gibbs sampling, the dependencies among the samples decay far too slowly for the estimate (8) to be useful for us (cf. the remarks in the Introduction). We will see in Section IV, how this issue is eased by tree-based Gibbs sampling as proposed by Hamze and de Freitas [24].

Second, it is usually assumed that f is strictly positive. In this case, $S_f = \mathcal{X}$ and $|S_f| = |\mathcal{X}|$ is known. However, this assumption excludes applications to constrained channels as in the example in Section II. Indeed, in that example, we would have $f(\mathbf{x}^{(k)}) = 1$ for all samples $\mathbf{x}^{(k)}$, and $|S_f| = Z_f$ is the desired unknown quantity. We will address this issue in Section IV-B.

With suitable modifications, which will address the mentioned issues, the Ogata-Tanemura method will turn out to work well for the capacity of noiseless constrained 2-D channels.

However, for our second application—the information rate of noisy 2-D constrained source/channel models—the Ogata-Tanemura method turns out to be inadequate. We will therefore resort to multilayer importance sampling as described below. We first describe the use of standard (single-layer) importance sampling to estimate Z_f .

Importance Sampling [22], [34]:

- 1) Draw samples $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(K)}$ from \mathcal{X} according to some auxiliary probability distribution $q(\mathbf{x}) = \frac{1}{Z_g}g(\mathbf{x})$, where $g(\mathbf{x})$ is a nonnegative function defined on \mathcal{X} satisfying $g(\mathbf{x}) \neq 0$ whenever $f(\mathbf{x}) \neq 0$.
- 2) Compute

$$\hat{R} = \frac{1}{K} \sum_{k=1}^{K} \frac{f(\mathbf{x}^{(k)})}{g(\mathbf{x}^{(k)})} \tag{9}$$

It is easy to verify that $E(\hat{R}) = Z_f/Z_g$.

The key issue with importance sampling is to find a useful function $g(\mathbf{x})$ such that

- $q(\mathbf{x})$ is a good approximation of $p(\mathbf{x})$ (so that $f(\mathbf{x})/g(\mathbf{x})$ does not wildly fluctuate),
- sampling from $q(\mathbf{x})$ is feasible, and



Fig. 2. Partition of Fig. 1 into two cycle-free parts (one part inside the two ovals, the other part outside the ovals).

• computing Z_g is feasible.

An obvious choice for $g(\mathbf{x})$ (and thus $q(\mathbf{x})$) is

$$g(\mathbf{x}) \stackrel{\scriptscriptstyle \triangle}{=} f(\mathbf{x})^{\alpha} \tag{10}$$

for $0 < \alpha < 1$. With this choice, any factorization of $f(\mathbf{x})$ results in a factorization of $g(\mathbf{x})$ that preserves the topology of the corresponding factor graph. (Note, however, that this choice of $g(\mathbf{x})$ is not helpful if $f(\mathbf{x})$ is $\{0, 1\}$ -valued.)

In order to sample from $q(\mathbf{x})$, we will again use tree-based Gibbs sampling (see Section IV-A).

In a variation of the algorithm, the estimator (9) of the ratio Z_f/Z_g could be replaced by Bennett's acceptance ratio method [35], which is also known as bridge sampling [36].

A function $g(\mathbf{x})$ with all the required properties may be hard to find, or it may not exist. This problem is addressed by multilayer importance sampling, which uses several auxiliary distributions.

Multilayer (Multi-Temperature) Importance Sampling:

Multilayer importance sampling, as described here, is a variation of annealed importance sampling as proposed by Neal [37], [34]; see also [38]. We use J parallel versions of importance sampling as follows. For $j = 0, 1 \dots, J$, let

$$g_j(\mathbf{x}) \stackrel{\scriptscriptstyle \Delta}{=} f(\mathbf{x})^{\alpha_j}$$
 (11)

with $0 \leq \alpha_J < \cdots < \alpha_1 < \alpha_0 = 1$. Note that $Z_{g_0} = Z_f$ and

$$\frac{Z_f}{Z_{g_J}} = \frac{Z_{g_0}}{Z_{g_1}} \frac{Z_{g_1}}{Z_{g_2}} \cdots \frac{Z_{g_{J-1}}}{Z_{g_J}}$$
(12)

For j = 1, 2, ..., J, compute the ratio $Z_{g_{j-1}}/Z_{g_j}$ by importance sampling as described before:

1) Draw samples $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(K)}$ from $q_i(\mathbf{x}) \propto q_i(\mathbf{x})$.

2) Compute

$$\hat{R}_{j} = \frac{1}{K} \sum_{k=1}^{K} \frac{g_{j-1}(\mathbf{x}^{(k)})}{g_{j}(\mathbf{x}^{(k)})}$$
(13)

$$= \frac{1}{K} \sum_{k=1}^{K} f(\mathbf{x}^{(k)})^{\alpha_{j-1} - \alpha_j}.$$
 (14)

Noting that $E(\hat{R}_j) = Z_{g_{j-1}}/Z_{g_j}$, we use $\prod_{j=1}^J \hat{R}_j$ as an estimate of Z_f/Z_{g_J} .

If the number of layers J is large, $g_j(\mathbf{x})$ is a good approximation of $g_{j-1}(\mathbf{x})$ at each layer j.

It remains to find an estimate of Z_{g_J} , which tends to be easier than the original problem of estimating Z_f . In particular, for $\alpha_J = 0$, we have $Z_{g_J} = |S_f|$, the cardinality of the support of f. In our application (Section VI), it turns out that Z_{g_J} can be computed efficiently by the tree-based Ogata-Tanemura method of Section IV-B.

IV. TREE-BASED GIBBS SAMPLING AND PARTITION FUNCTION ESTIMATION

A. Tree-Based Gibbs Sampling

Tree-based Gibbs sampling was proposed by Hamze and de Freitas [24]. It works as follows. Let (A, B) be a partition of the index set $\{1, 2, ..., N\}$ such that,

- for fixed \mathbf{x}_A , the factor graph of $f(\mathbf{x}) = f(\mathbf{x}_A, \mathbf{x}_B)$ is cycle-free, and
- for fixed \mathbf{x}_B , the factor graph of $f(\mathbf{x}) = f(\mathbf{x}_A, \mathbf{x}_B)$ is also cycle-free.

An example of such a partition is shown in Fig. 2. Starting from some initial configuration $\mathbf{x}^{(0)} = (\mathbf{x}_A^{(0)}, \mathbf{x}_B^{(0)})$, the samples $\mathbf{x}^{(k)} = (\mathbf{x}_A^{(k)}, \mathbf{x}_B^{(k)})$, $k = 1, 2, \ldots$, are created as follows: first, $\mathbf{x}_A^{(k)}$ is drawn from

$$p_f(\mathbf{x}_A | \mathbf{x}_B = \mathbf{x}_B^{(k-1)}) \propto f(\mathbf{x}_A, \mathbf{x}_B^{(k-1)});$$
(15)

second, $\mathbf{x}_B^{(k)}$ is drawn from

$$p_f(\mathbf{x}_B | \mathbf{x}_A = \mathbf{x}_A^{(k)}) \propto f(\mathbf{x}_A^{(k)}, \mathbf{x}_B).$$
(16)

The point is that drawing these samples is easy (by means of backward-filtering forward-sampling, see the appendix) since the corresponding factor graphs are cycle-free.

As in standard Gibbs sampling, the samples $(\mathbf{x}_A^{(k)}, \mathbf{x}_B^{(k)})$ are eventually (i.e., for $k \to \infty$) drawn from p_f (provided that the corresponding Markov chain is irreducible and aperiodic [23]). However, tree-based Gibbs sampling mixes faster than naive Gibbs sampling.

B. Application to Partition Function Estimation

With A and B as above, let

$$f_A(\mathbf{x}_A) \stackrel{\scriptscriptstyle \Delta}{=} \sum_{\mathbf{x}_B} f(\mathbf{x}_A, \mathbf{x}_B),$$
 (17)

and

$$f_B(\mathbf{x}_B) \stackrel{\scriptscriptstyle \triangle}{=} \sum_{\mathbf{x}_A} f(\mathbf{x}_A, \mathbf{x}_B).$$
 (18)

We then note that

$$Z_{f_A} = \sum_{\mathbf{x}_A} f_A(\mathbf{x}_A) \tag{19}$$

$$= Z_f, \tag{20}$$

i.e., f_A (and analogously f_B) has the same partition function as f itself.

We also note that samples $\mathbf{x}_{A}^{(1)}$, $\mathbf{x}_{A}^{(2)}$, ..., from

$$p_{f_A}(\mathbf{x}_A) \stackrel{\scriptscriptstyle \Delta}{=} \frac{f_A(\mathbf{x}_A)}{Z_f} = \sum_{\mathbf{x}_B} p_f(\mathbf{x}_A, \mathbf{x}_B)$$
 (21)

can be obtained by tree-based Gibbs sampling as in Section IV-A simply by dropping the second component (the *B*-part) of the samples $(\mathbf{x}_A^{(1)}, \mathbf{x}_B^{(1)}), (\mathbf{x}_A^{(2)}, \mathbf{x}_B^{(2)}), \dots$

We can now use the Ogata-Tanemura method (Section III) to estimate $\Gamma_f = 1/Z_f$ in two different ways. One way is to directly use the estimator (8) with samples $\mathbf{x}^{(k)} = (\mathbf{x}_A^{(k)}, \mathbf{x}_B^{(k)})$ as in Section IV-A. Another way is to apply the estimator (8) to f_A , i.e., to compute

$$\hat{\Gamma}_{f_A} \stackrel{\triangle}{=} \frac{1}{K|\mathcal{S}_{f_A}|} \sum_{k=1}^K \frac{1}{f_A(\mathbf{x}_A^{(k)})}$$
(22)

which is also an estimate of $1/Z_f$. The computation of

$$f_A(\mathbf{x}_A^{(k)}) = \sum_{\mathbf{x}_B} f(\mathbf{x}_A^{(k)}, \mathbf{x}_B),$$
(23)

which is required in (22), is easy since the corresponding factor graph is a tree; in fact, the result of this computation is obtained for free as a by-product of tree-based Gibbs sampling as is explained in the appendix. By symmetry, we also have an analogous estimate $\hat{\Gamma}_{f_B}$ defined as

$$\hat{\Gamma}_{f_B} \stackrel{\scriptscriptstyle \triangle}{=} \frac{1}{K|\mathcal{S}_{f_B}|} \sum_{k=1}^{K} \frac{1}{f_B(\mathbf{x}_B^{(k)})}$$
(24)

With the same samples $(\mathbf{x}_A^{(1)}, \mathbf{x}_B^{(1)})$, $(\mathbf{x}_A^{(2)}, \mathbf{x}_B^{(2)})$, ..., estimating $1/Z_f$ from (22) and (24) tends to converge faster than (8). More importantly for this paper, the direct Ogata-Tanemura method (8) cannot be used for constrained channels (cf. the example in Section II) where $|\mathcal{S}_f| = Z_f$ is the desired quantity. In contrast, the computation of $|\mathcal{S}_{f_A}|$ in (22) and $|\mathcal{S}_{f_B}|$ in (24) may be easy in such cases as will be exemplified below.

V. APPLICATION TO THE CAPACITY OF NOISELESS 2-D CONSTRAINED CHANNELS

We demonstrate the estimators (22) and (24) by their application to the example in Section II, the 2-D $(1, \infty)$ runlengthlimited constrained channel.

We will use factor graphs as in Fig. 1 with a partitioning as in Fig. 2. In this example, the quantities $|S_{f_A}|$ and $|S_{f_B}|$, which are needed in (22) and (24), respectively, are given by

$$|\mathcal{S}_{f_A}| = \sum_{\mathbf{x}_A} f(\mathbf{x}_A, \mathbf{0})$$
 (25)

$$|\mathcal{S}_{f_B}| = \sum_{\mathbf{x}_B} f(\mathbf{0}, \mathbf{x}_B), \qquad (26)$$



Fig. 3. Estimated noiseless capacity (in bits/symbol) vs. the number of samples K for a 10×10 grid with a $(1, \infty)$ constraint. The plot shows 10 different sample paths, each with two estimates, one from Γ_A and another from Γ_B . The horizontal dotted line shows the infinite-size capacity for this constraint.



Fig. 4. Estimated noiseless capacity (in bits/symbol) vs. the number of samples K for a 60×60 grid with a $(1, \infty)$ constraint. The plot shows 10 different sample paths, each with two estimates, one from Γ_A and another from Γ_B . The horizontal dotted line shows the infinite-size capacity for this constraint.

since

$$f(\mathbf{x}_A, \mathbf{0}) = \begin{cases} 1, & \text{if } f_A(\mathbf{x}_A) > 0\\ 0, & \text{if } f_A(\mathbf{x}_A) = 0, \end{cases}$$
(27)

and likewise for $f(\mathbf{0}, \mathbf{x}_B)$. It follows that $|\mathcal{S}_{f_A}|$ and $|\mathcal{S}_{f_B}|$ are easy to compute by sum-product message passing in the cycle-free factor graphs of $f(\mathbf{x}_A, \mathbf{0})$ and $f(\mathbf{0}, \mathbf{x}_B)$, respectively.

Some experimental results are shown in Figs. 3 through 6. All figures refer to f as in (4) and (5) and show the estimates of the finite-size capacity C_M (6) obtained from (22) and (24) vs. K for several different experiments.

In Fig. 3, we have $N = 10 \times 10$ and we obtain $C_{10} \approx 0.6082$ bits/symbol. In Fig. 4, we have $N = 60 \times 60$, and there are issues with slow convergence.

In order to speed up the mixing and thus improving the convergence, we now partition the factor graph into vertical strips each containing $M \times 2$ or $M \times 3$ variables (rather than $M \times 1$



Fig. 5. Same conditions as in Fig. 4, but with strips of width two.



Fig. 6. Same conditions as in Fig. 4, but with strips of width three.

variables as in Fig. 2). Exact sum-product message passing is still possible on such strips, e.g., by converting the strip into an equivalent cycle-free factor graph. The computation time is exponential in the width of the strip, but the faster mixing results in a substantial reduction of total computation time for strips of moderate width.

Simulation results for strips of width 2 and 3 are shown in Figs. 5 and 6, respectively, both for a grid of size $N = 60 \times 60$. Note that the convergence is much better than in Fig. 4, and we obtain $C_{60} \approx 0.5914$ bits/symbol.

Also shown in these figures, as a horizontal dotted line, is the infinite-size capacity $C_{\infty} \approx 0.5879$ bits/symbol from [32], which (in this example) is a lower bound on the finite-size capacity.

VI. ESTIMATING THE INFORMATION RATE OF NOISY 2-D SOURCE / CHANNEL MODELS

A. The Problem

We now consider the problem of computing the information rate of noisy 2-D source/channel models. Although the focus of this paper is on constrained channels, the proposed approach



Fig. 7. Extension of Fig. 1 to a factor graph of $p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$ with $p(\mathbf{y}|\mathbf{x})$ as in (29).

can also be applied to other 2-D source/channel models such as 2-D channels with intersymbol interference.

For a 2-D grid of size $N = M \times M$ (as before), let $\mathbf{X} = \{X_1, X_2, \dots, X_N\}$ be the source process (i.e., the input process of the channel) and let $\mathbf{Y} = \{Y_1, Y_2, \dots, Y_N\}$ be the output process of the channel; \mathbf{X} takes values in \mathcal{X} (as defined in Section II) and Y takes values in \mathbb{R}^N . We wish to compute the mutual information rate

$$\frac{1}{N}I(\mathbf{X};\mathbf{Y}) = \frac{1}{N}(H(\mathbf{Y}) - H(\mathbf{Y}|\mathbf{X}))$$
(28)

for cases where $p(\mathbf{x}, \mathbf{y})$ has a suitable factor graph. In particular, we will focus on the case where the channel is memoryless, i.e.,

$$p(\mathbf{y}|\mathbf{x}) = \prod_{n=1}^{N} p(y_n|x_n),$$
(29)

and where the channel input distribution $p(\mathbf{x})$ has a factorization with a factor graph as in Fig. 1. It then follows that $p(\mathbf{x}, \mathbf{y})$ has a factor graph as in Fig. 7.

In many cases of practical interest, $H(\mathbf{Y}|\mathbf{X})$ is analytically available, see our numerical experiments in Section VI-C. In such cases, the problem of computing the mutual information rate (28) reduces to computing

$$H(\mathbf{Y}) = \mathbf{E} \left[-\log_2 p(\mathbf{Y}) \right]. \tag{30}$$

If $H(\mathbf{Y}|\mathbf{X})$ is not analytically available, it can be computed by the same method as $H(\mathbf{Y})$, see [2, Section III].

B. The Method

As in [2], we approximate the expectation in (30) by the empirical average

$$H(\mathbf{Y}) \approx -\frac{1}{L} \sum_{\ell=1}^{L} \log_2(p(\mathbf{y}^{(\ell)})), \qquad (31)$$

where samples $\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(L)}$ are drawn according to $p(\mathbf{y})$. The problem of estimating the mutual information rate is thus reduced to



Fig. 8. Uniform-input information rate (in bits/symbol) vs. SNR for a 24×24 channel with a $(1,\infty)$ constraint and additive white Gaussian noise. The horizontal dotted line shows the noiseless capacity of this channel.

- creating samples y^(l) and
 computing p(y^(l)) for each sample.

If $p(\mathbf{x}, \mathbf{y})$ has a cycle-free factor graph (and if $|\mathcal{X}_1|, |\mathcal{X}_2|, \ldots, |\mathcal{X}_N|$ are not too large), then both tasks can be carried out in a single-loop algorithm as in [2]. In this paper, however, we assume that no such factor graph exists and we propose a double-loop algorithm (with an outer loop and an inner loop) to carry out these tasks. In the outer loop, we create samples $\mathbf{y}^{(1)}, \ldots, \mathbf{y}^{(L)}$ as follows.

- 1) Draw samples $\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(L)}$ according to $p(\mathbf{x})$. In simple cases (such as unconstrained channels with i.i.d. input), this may be trivial; in general, however, we do this by tree-based Gibbs sampling (as in Section IV-A) using the factor graph of $p(\mathbf{x})$.
- 2) For $\ell = 1, \ldots, L$, draw $\mathbf{y}^{(\ell)}$ from $p_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}^{(\ell)})$, i.e., by simulating the channel with input $\mathbf{x}^{(\ell)}$.

In the inner loop, we compute an estimate of

$$p(\mathbf{y}^{(\ell)}) = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \, p_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}^{(\ell)}|\mathbf{x}) \tag{32}$$

as follows. Note that, for fixed ℓ , the right-hand side of (32) is the partition function $Z_{f_{\ell}}$ of

$$f_{\ell}(\mathbf{x}) \stackrel{\scriptscriptstyle \triangle}{=} p(\mathbf{x}) \, p_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}^{(\ell)}|\mathbf{x}), \tag{33}$$

which has a suitable factor graph (as, e.g., in Fig. 7). In principle, we can thus estimate (32) by any of the methods of Section III. It turns out, however, that only multilayer importance sampling is able to handle the more demanding cases (as will be explained in our numerical experiments in Section VI-C) while the other methods of Section III suffer from slow and erratic convergence.

C. Numerical Experiments

In our numerical experiments, we consider a noisy version of the example in Section II, i.e., a noisy version of the 2-D $(1,\infty)$ runlength-limited constrained channel. We assume that the channel input distribution $p(\mathbf{x})$ is uniform over the



Fig. 9. Estimated information rate (in bits/symbol) vs. the number of samples L for a noisy 24×24 $(1, \infty)$ constraint at 0 dB. The plot shows 12 different sample paths.



Fig. 10. Estimated information rate (in bits/symbol) vs. the number of samples L for a noisy 24×24 $(1, \infty)$ constraint at 6 dB. The plot shows 12 different sample paths.

allowed configurations, i.e., $p(\mathbf{x}) = p_f(\mathbf{x})$ with f as in (4), and we assume that the noise is additive white Gaussian (and independent of **X**), i.e., $p(\mathbf{y}|\mathbf{x})$ is a product as in (29) with factors

$$p(y_n|x_n) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} \left(y_n - (-1)^{x_n}\right)^2\right) \quad (34)$$

and thus

$$H(\mathbf{Y}|\mathbf{X}) = \frac{N}{2}\log_2(2\pi e\sigma^2).$$
 (35)

We will use the signal-to-noise ratio (SNR) defined as

$$\operatorname{SNR} \stackrel{\scriptscriptstyle \triangle}{=} \frac{1}{\sigma^2},$$
 (36)

which we will specify in dB, i.e., $10 \log_{10}(SNR)$.

The grid size in all the plots is $N = 24 \times 24$ and the parameters α_j in (11) are set to $\alpha_j = 2^{-j}$, for j = 1, 2, ..., J.

Fig. 8 shows the computed information rate vs. the SNR over the interval [-10, 8] dB. The horizontal dotted line

in Fig. 8 shows the capacity of the noiseless version of this channel, which is about 0.596 bits per symbol.

Figs. 9 and 10 illustrate the convergence of the outer loop of the proposed double-loop algorithm at 0 dB and at 6 dB, respectively. Both figures show the estimated information rate vs. the number of samples L in (31) for 12 different sample paths.

As for the inner loop, the required number of layers J in (12) depends on the SNR. As the SNR increases (or equivalently as σ^2 decreases), the function $f_{\ell}(\mathbf{x})$ in (33) tends to have more isolated modes. Therefore, in order to obtain a good approximation at each level of multilayer importance sampling, larger values of J are required for higher SNR. In our numerical experiments at 0 dB and 6 dB, J is set to 3 and 6, respectively.

Fig. 11 shows the convergence of $\log_2 \hat{R}_j$ as in (14) for j = 1, 2, ..., 6, for a fixed output sample at 6 dB.

The value of Z_{g_J} is estimated by the tree-based Ogata-Tanemura method of Section IV-B. Fig. 12 shows the convergence of the estimate of $\log_2(Z_{g_6})/N$ according to (22) for four different sample paths.

D. Remarks

In statistical physics, the partition function typically has the form

$$Z = \sum_{\mathbf{x} \in \mathcal{X}} e^{-E(\mathbf{x})/T},$$
(37)

where T is the temperature and $E(\mathbf{x})$ is the energy of the configuration x. We therefore point out that the noise variance σ^2 in (34) may be viewed as the temperature parameter of the partition function $Z_{f_{\ell}}$ of (33). It is well known in statistical physics that computing the partition function is harder at low temperatures than at high temperatures. Similarly, we observe that computing the partition function $Z_{f_{\ell}}$ of (33) is harder at high SNR than at low SNR; in particular, at high SNR, more layers (higher values of J) are required for multilayer (multi-temperature) importance sampling.

We also note that, in the examples of Section VI-C, the choice of the parameters $\alpha_j = 2^{-j}$ in (11) is somewhat arbitrary. It is possible that other choices of these parameters result in faster convergence.

VII. CONCLUSION

Monte Carlo methods have been highly succesful in computing the information rate of source/channel models with 1-D memory. The extension of such methods to source/channel models with 2-D memory has been an open research problem. In this paper, we develop such methods with a focus on the (difficult) case of channels with input constraints, with or without noise. In contrast to previous techniques, which either use generalized belief propagation or compute only bounds on the information rate, the Monte Carlo algorithms of this paper are guaranteed to converge (asymptotically) to the desired information rate. A key role in the proposed algorithms is played by tree-based Gibbs sampling by Hamze and de Freitas, which we have shown to yield an estimate of the partition function as a by-product. The success of the proposed methods



Fig. 11. Computed $\log_2 \hat{R}_j$ as in (14), for $j = 1, 2, \ldots, 6$ vs. the number of samples K for a noisy 24×24 $(1, \infty)$ runlength-limited constraint at 6 dB. The plot shows $\log_2 \hat{R}_6, \log_2 \hat{R}_5, \ldots, \log_2 \hat{R}_1$ from top to bottom.



Fig. 12. Estimated $\log_2(Z_{g_6})/N$ vs. the number of samples K for a noisy 24×24 $(1, \infty)$ runlength-limited constraint at 6 dB.

is exemplified by Fig. 8, which (to the best of our knowledge) is the first such plot for a noisy 2-D channel. We also note that the extension of the proposed methods to computing upper and lower bounds on the information rate as in [2, Section VI] is straightforward.

ACKNOWLEDGEMENT

The authors wish to thank Radford Neal for helpful discussions and advice on annealed importance sampling. The authors also wish to thank David MacKay and Iain Murray for pointing out to us [24], and the reviewers and the Associate Editor for helpful comments.

APPENDIX

SAMPLING FROM MARKOV CHAINS

We recall some pertinent facts about the simulation of Markov chains and cycle-free factor graphs. Let $p(\mathbf{x}) = p(x_1, \ldots, x_n)$ be the probability mass function of a Markov



Fig. 13. Forney factor graph of (39) with messages $\overleftarrow{\mu}_{X_k}$ (40).

chain. If $p(\mathbf{x})$ is given in the form

$$p(\mathbf{x}) = p(x_1) \prod_{k=2}^{n} p(x_k | x_{k-1}),$$
(38)

then it is obvious how to draw i.i.d. samples according to $p(\mathbf{x})$. Now consider the case where $p(\mathbf{x})$ is not given in the form (38), but in the more general form

$$p(\mathbf{x}) \propto \prod_{k=2}^{n} g_k(x_{k-1}, x_k)$$
(39)

with general factors g_k . It is then still easy to draw i.i.d. samples according to $p(\mathbf{x})$, which may be seen as follows. First, a probability mass function of the form (39) can be rewritten in the form (38) (which allows efficient simulation). Second, this reparameterization of $p(\mathbf{x})$ may be efficiently carried out by backward sum-product message passing, as will be detailed below. The resulting algorithm is known as "backward-filtering forward-sampling" (or, in a time-reversed version, as "forward-filtering backward-sampling") [39].

Specifically, let $\overleftarrow{\mu}_{X_k}$ be the backward sum-product message along the edge X_k in the factor graph of (39), as is illustrated in Fig. 13 (cf. [31]). We then have $\overleftarrow{\mu}_{X_n}(x_n) = 1$ and

$$\overleftarrow{\mu}_{X_k}(x_k) \stackrel{\scriptscriptstyle \triangle}{=} \sum_{x_{k+1}} g_{k+1}(x_k, x_{k+1}) \overleftarrow{\mu}_{X_{k+1}}(x_{k+1}) \quad (40)$$

$$=\sum_{x_{k+1},\dots,x_n}\prod_{m=k+1}^n g_m(x_{m-1},x_m)$$
(41)

for $k = n - 1, n - 2, \dots, 1$. Then

$$p(x_1) = \sum_{x_2, \dots, x_n} p(x_1, \dots, x_n)$$
(42)

$$\propto \overleftarrow{\mu}_{X_1}(x_1) \tag{43}$$

and

$$p(x_k|x_{k-1}) = \frac{g_k(x_{k-1}, x_k)\overleftarrow{\mu}_{X_k}(x_k)}{\overleftarrow{\mu}_{X_{k-1}}(x_{k-1})}$$
(44)

for k = 2, ..., n. The proof of (44) follows from noting that

$$p(x_{k-1}) = \gamma \,\overrightarrow{\mu}_{X_{k-1}}(x_{k-1}) \,\overleftarrow{\mu}_{X_{k-1}}(x_{k-1}) \tag{45}$$

and

$$p(x_{k-1}, x_k) = \gamma \overrightarrow{\mu}_{X_{k-1}}(x_{k-1})g_k(x_{k-1}, x_k)\overleftarrow{\mu}_{X_k}(x_k) \quad (46)$$

where $\overrightarrow{\mu}_{X_{k-1}}$ is the forward sum-product message along the edge X_{k-1} and where γ is the missing scale factor in (39).

We also note that

$$\sum_{x_1} \overleftarrow{\mu}_{X_1}(x_1) = \sum_{\mathbf{x}} g(\mathbf{x}), \tag{47}$$

where $g(\mathbf{x})$ is defined as the right-hand side of (39). In this paper, this fact is used to compute the marginals (23) as a by-product of the sampling.

The generalization of all this to arbitrary factor graphs without cycles is straightforward.

REFERENCES

- D. Arnold and H.-A. Loeliger, "On the information rate of binary-input channels with memory," *Proc. 2001 IEEE Int. Conf. on Communications*, Helsinki, Finland, June 11–14, 2001, pp. 2692–2695.
- [2] D. Arnold, H.-A. Loeliger, P. O. Vontobel, A. Kavčić, and W. Zeng, "Simulation-based computation of information rates for channels with memory," *IEEE Trans. Inf. Theory*, vol. 52, no. 8, August 2006, pp. 3498–3508.
- [3] H. D. Pfister, J.-B. Soriaga, and P. H. Siegel, "On the achievable information rates of finite-state ISI channels," *Proc. 2001 IEEE Globecom*, San Antonio, USA, Nov. 2001, pp. 2992–2996.
- [4] K. A. S. Immink, Codes for Mass Data Storage Systems. Eindhoven: Shannon Foundation Publishers, 2004.
- [5] R. Roth, Introduction to Coding Theory. Cambridge University Press, 2006.
- [6] B. Vasic and E. M. Kurtas, Coding and Signal Processing for Magnetic Recording Systems. CRC Press, 2005.
- [7] K. A. S. Immink, P. H. Siegel, and J. K. Wolf, "Codes for digital recorders," *IEEE Trans. Inf. Theory*, vol. 44, Oct. 1998, pp. 2260–2299.
- [8] R. Wood, M. Williams, A. Kavčić, and J. Miles, "The feasibility of magnetic recording at 10 terabits per square inch on conventional media," *IEEE Trans. Magnetics*, vol. 45, Feb. 2009, pp. 917–923.
- [9] C. E. Shannon, "A mathematical theory of communications," *Bell Sys. Tech. Journal*, vol. 27, July 1948, pp. 379–423.
- [10] N. J. Calkin and H. S. Wilf, "The number of independent sets in a grid graph," SIAM J. Discr. Math., vol. 11, Feb. 1998, pp. 54–60.
- [11] W. Weeks IV and R. E. Blahut, "The capacity and coding gain of certain checkerboard codes," *IEEE Trans. Inf. Theory*, vol. 44, May 1998, pp. 1193–1203.
- [12] K. Kato and K. Zeger, "On the capacity of two-dimensional run-length constrained channels," *IEEE Trans. Inf. Theory*, vol. 45, July 1999, pp. 1527–1540.
- [13] H. Ito, A. Kato, Z. Nagy, and K. Zeger, "Zero capacity region of multidimensional run length constraints," *The Electronic Journal of Combinatorics*, vol. 6(1), 1999.
- [14] K. Censor and T. Etzion, "The positive capacity region of twodimensional run-length-constrained channels," *IEEE Trans. Inf. Theory*, vol. 52, Nov. 2006, pp. 5128–5140.
- [15] I. Tal and R. M. Roth, "Concave programming upper bounds on the capacity of 2-D constraints," *IEEE Trans. Inf. Theory*, vol. 57, Jan. 2011, pp. 381–391.
- [16] H.-A. Loeliger and M. Molkaraie, "Simulation-based estimation of the partition function and the information rate of two-dimensional models," *Proc. 2008 IEEE Int. Symp. on Information Theory*, Toronto, Canada, July 6–11, 2008, pp. 1113–1117.
- [17] H.-A. Loeliger and M. Molkaraie, "Estimating the partition function of 2-D fields and the capacity of constrained noiseless 2-D channels using tree-based Gibbs sampling," *Proc. 2009 IEEE Information Theory Workshop*, Taormina, Italy, October 11–16, 2009, pp. 228–232.
- [18] M. Molkaraie and H.-A. Loeliger, "Estimating the information rate of noisy constrained 2-D channels," *Proc. 2010 IEEE Int. Symp. on Information Theory*, Austin, USA, June 13–18, 2010, pp. 1678–1682.

- [19] Y. Ogata and M. Tanemura, "Estimation of interaction potentials of spatial point patterns through the maximum likelihood procedure," *Ann. Inst. Statist. Math.*, vol. 22, 1981, pp. 315–338.
- [20] M. Jerrum and A. Sinclair, "Polynomial-time approximation algorithms for the Ising model," *SIAM J. Computing*, vol. 11, Oct. 1993, pp. 1087– 1116.
- [21] G. Potamianos and J. Goutsias, "Stochastic approximation algorithms for partition function estimation of Gibbs random fields," *IEEE Trans. Inf. Theory*, vol. 43, Nov. 1997, pp. 1984–1965.
- [22] D. J. C. MacKay, "Introduction to Monte Carlo methods," in *Learning in Graphical Models*, M. I. Jordan, ed., Kluwer Academic Press, 1998, pp. 175–204.
- [23] P. Brémaud, Markov Chains: Gibbs Fields, Monte Carlo Simulations, and Queues. Springer, 1999.
- [24] F. Hamze and N. de Freitas, "From fields to trees," Proc. 2004 Conf. on Uncertainty in Artificial Intelligence, Banff, Canada, July 7–11, 2004, pp. 243–250.
- [25] O. Shental, N. Shental, S. Shamai (Shitz), I. Kanter, A. J. Weiss, and Y. Weiss, "Discrete-input two-dimensional Gaussian channels with memory: estimation and information rates via graphical models and statistical mechanics," *IEEE Trans. Inf. Theory*, vol. 54, April 2008, pp. 1500–1513.
- [26] J. Chen, and P. H. Siegel, "On the symmetric information rate of twodimensional finite-state ISI channels," *IEEE Trans. Inf. Theory*, vol. 52, no. 1, Jan. 2006, pp. 227–236.
- [27] J. S. Yedidia, W. T. Freeman, and Y. Weiss, "Constructing free energy approximations and generalized belief propagation algorithms," *IEEE Trans. Inf. Theory*, vol. 51, July 2005, pp. 2282–2312.
- [28] G. Sabato and M. Molkaraie, "Generalized belief propagation algorithm to estimate the capacity of multi-dimensional run-length limited constraints," *Proc. 2010 IEEE Int. Symp. on Information Theory*, Austin, USA, June 13–18, 2010, pp. 1213–1217.
- [29] G. Sabato and M. Molkaraie, "Generalized belief propagation for the noiseless capacity and information rates of run-length limited constraints," *IEEE Trans. Comm.*, vol. 60, March 2012, pp. 669–675.
- [30] F. R. Kschischang, B. J. Frey, and H.-A. Loeliger, "Factor graphs and the sum-product algorithm," *IEEE Trans. Inf. Theory*, vol. 47, Feb. 2001, pp. 498–519.
- [31] H.-A. Loeliger, "An introduction to factor graphs," *IEEE Signal Proc. Mag.*, Jan. 2004, pp. 28–41.
- [32] Z. Nagy and K. Zeger, "Capacity bounds for the three-dimensional (0, 1) run-length limited channel," *IEEE Trans. Inf. Theory*, vol. 46, May 2000, pp. 1030–1033.
- [33] S. Geman and D. Geman, "Stochastic relaxation, Gibbs distribution, and Bayesian restoration of images," *IEEE Trans. Pattern Analys. and Machine Intell.*, vol. 6, 1984, pp. 721–741.
- [34] R. M. Neal, "Annealed importance sampling," *Statistics and Computing*, vol. 11, April 2001, pp. 125–139.
- [35] C. H. Bennett, "Efficient estimation of free energy differences from Monte Carlo data," *Journal of Computational Physics*, vol. 22, October 1976, pp. 245–268.
- [36] X.-L. Meng and H. W. Wong, "Simulating ratios of normalizing constants via a simple identity: A theoretical exploration," *Statistica Sinica* vol. 6, Oct. 1996, pp. 831–860.
- [37] R. M. Neal, "Annealed importance sampling," Techn. Report 9850, Dept. of Statistics, University of Toronto, 1998.
- [38] C. Jarzynski, "Nonequilibrium equality for free energy differences," *Phys. Rev. Lett.*, vol. 78, 1997, pp. 2690–2693.
- [39] D. Gamerman and H. F. Lopes, Markov Chain Monte Carlo: Stochastic Simulation for Bayesian Inference. 2nd ed., CRC Press, 2006.