

# Importance sampling for weighted binary random matrices with specified margins

Matthew T. Harrison\*

Jeffrey W. Miller\*

Division of Applied Mathematics

Brown University

Providence, RI 02912

January 18, 2013

## Abstract

A sequential importance sampling algorithm is developed for the distribution that results when a matrix of independent, but not identically distributed, Bernoulli random variables is conditioned on a given sequence of row and column sums. This conditional distribution arises in a variety of applications and includes as a special case the uniform distribution over zero-one tables with specified margins. The algorithm uses dynamic programming to combine hard margin constraints, combinatorial approximations, and additional non-uniform weighting in a principled way to give state-of-the-art results.

KEYWORDS: bipartite graph, conditional inference, permanent, Rasch model, uniform distribution

## 1 Introduction

Let  $\Omega^*$  denote the set of  $m \times n$  binary matrices with row sums  $\mathbf{r} = (r_1, \dots, r_m)$  and column sums  $\mathbf{c} = (c_1, \dots, c_n)$ , and let  $\mathbf{w} = (w_{ij}) \in [0, \infty)^{m \times n}$  be a given nonnegative matrix. Define the distribution  $P^*$  on  $\{0, 1\}^{m \times n}$  via

$$P^*(\mathbf{z}) = \frac{1}{\kappa} \prod_{ij} w_{ij}^{z_{ij}} \mathbb{1}\{\mathbf{z} \in \Omega^*\}, \quad \kappa = \sum_{\mathbf{z} \in \Omega^*} \prod_{ij} w_{ij}^{z_{ij}}, \quad (1)$$

where  $\mathbb{1}$  is the indicator function and where we assume  $\kappa > 0$ .  $P^*$  is the conditional distribution of an  $m \times n$  array of independent Bernoulli random variables, say  $\mathbf{B} = (B_{ij})$ ,

---

\*Matthew T. Harrison is Assistant Professor of Applied Mathematics, Brown University, Providence, RI 02912 (email: Matthew.Harrison@Brown.edu) and Jeffrey W. Miller is a graduate student of Applied Mathematics, Brown University, Providence, RI 02912 (email: Jeffrey\_Miller@Brown.edu). This work was supported in part by the National Science Foundation (NSF) grant DMS-1007593, the Defense Advanced Research Projects Agency (DARPA) contract FA8650-11-1-7151, and, while MTH was in the Department of Statistics at Carnegie Mellon University, by the NSF grant DMS-0240019 and the National Institutes of Health (NIH) grant NIMH-2RO1MH064537. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the NSF, the NIH, or DARPA. The authors thank Sam Kou for sharing his code for approximating  $\alpha$ -permanents.

with  $\mathbb{P}(B_{ij} = 1) = w_{ij}/(1 + w_{ij})$  given the margins  $\mathbf{r}$  and  $\mathbf{c}$ , where  $\mathbb{P}$  denotes probability. This paper describes an importance sampling algorithm that can be used for Monte Carlo approximation of probabilities and expectations under  $P^*$  and also for Monte Carlo approximation of  $\kappa$ . After a preprocessing step, sampling from our proposal distribution requires  $O(md)$  operations per matrix, where  $d = \sum_j c_j = \sum_i r_i$  is the total number of ones in the matrix.

We are not aware of any existing importance sampling algorithms that permit practical inference under  $P^*$ , although many special cases have been studied in the literature. For example, if  $\mathbf{w} \equiv 1$  then  $P^*$  is the uniform distribution over zero-one tables with specified margins, or equivalently, the uniform distribution over bipartite graphs with specified degree sequence. For square matrices, if  $\mathbf{w}$  is identically one except with a zero diagonal, then  $P^*$  corresponds to the uniform distribution over directed graphs with specified degrees. And if  $\mathbf{r} = \mathbf{c} \equiv 1$ , then  $P^*$  is a distribution over weighted permutation matrices and  $\kappa$  is the permanent of  $\mathbf{w}$ . Empirically, our algorithm outperforms all existing importance sampling algorithms in these special cases. Although our algorithm works well for most examples arising in practice, performance depends on  $\mathbf{r}$ ,  $\mathbf{c}$ , and  $\mathbf{w}$ . Highly irregular margins or highly variable  $\mathbf{w}$ , particularly many zero entries in  $\mathbf{w}$ , tend to cause poor performance.

$P^*$  factors in such a way that we need only focus on the distribution, say  $P$ , of the first column. The columns are sampled sequentially, with each successive column viewed as the first column of a smaller matrix with updated margins based on the previously sampled columns. We decompose the structure of  $P$  into margin constraints, combinatorial factors, and non-uniform weighting terms, combine approximations of these terms in a principled way, and then use a dynamic programming algorithm to exactly and efficiently sample from the resulting proposal distribution  $Q$  for the first column. Sequentially sampling columns in this way defines a proposal distribution  $Q^*$  for the whole matrix. This strategy for algorithm design works well for many similar problems, including symmetric matrices and nonnegative integer-valued matrices, each of which will be described elsewhere owing to space constraints. It seems likely that the design principles used for our approach are applicable much more broadly.

## 2 Motivating applications

### 2.1 Conditional inference for graphs and tables

Let  $\mathbf{B} \in \{0, 1\}^{m \times n}$  be a matrix of independent Bernoulli random variables with

$$\text{logit } \mathbb{P}(B_{ij} = 1) = \alpha_i + \beta_j + \sum_k \theta_k \xi_{kij}, \quad (2)$$

where  $\alpha$ ,  $\beta$ , and  $\theta$  are parameters, perhaps with constraints to ensure identifiability, and  $\xi$  is a collection of observed covariates. Models of this form arise, for example, in educational testing, where  $B_{ij}$  indicates whether or not subject  $i$  responded correctly to question  $j$ . If  $\theta \equiv 0$ , then the model reduces to the classical Rasch model (Rasch, 1960, 1961). Otherwise, it is an extension of the Rasch model to include item-specific covariate effects, such as each subject's prior exposure to the content being tested in each question. This model is also a simple version of models used for the analysis of network data (c.f., Holland & Leinhardt, 1981; Fienberg et al., 1985; Goldenberg et al., 2010) where  $\mathbf{B}$  is the adjacency matrix of a directed graph,  $\alpha$  and  $\beta$  allow for degree heterogeneity, and  $\xi$  is a collection of edge-specific

covariates. For example, for social network data we might have that  $B_{ij}$  indicates whether subject  $i$  reported subject  $j$  as a friend,  $\alpha_i$  controls the relative propensity for subject  $i$  to report friends,  $\beta_j$  controls the relative propensity for subject  $j$  to be reported as a friend, and  $\xi_{kij}$  indicates whether the relationship between subject  $i$  and  $j$  is of type  $k$ .

In both of these examples, if  $\theta$  is the only parameter of interest, then the nuisance parameters  $\alpha$  and  $\beta$  complicate inference and can be removed by conditioning on the row and column sums of  $\mathbf{B}$  (e.g., Cox, 1958; Holland & Leinhardt, 1981; Mehta & Patel, 1995; Harrison, 2012). Conditioning also results in inferential procedures that are robust to modeling assumptions, implicit in (2), about the distribution of the margins. The resulting conditional model that drives inference is exactly  $P^* = P_\theta^*$  with  $w_{ij} = \exp(\sum_k \theta_k \xi_{kij})$ , perhaps with the additional constraint that  $w_{ii} = 0$  in the case of network data. The conditional model is a natural exponential family in  $\theta$  with no nuisance parameters, but with an intractable normalization constant  $\kappa = \kappa_\theta$ . Harrison (2012, Example 4.2) provides details about an importance sampling approach to exact conditional inference for this model. The example there is based on a preliminary version of the algorithm presented here.

Other approaches to conditional inference in this setting include exhaustive enumeration, such as the algorithms for conditional logistic regression in Stata (StataCorp, 2009) and LogXact (Cytel, 2010), Markov chain Monte Carlo approaches, such as the `elrm` R package (Zamar et al., 2007), and analytic approximations, such as the `cond` R package (Brazzale, 2005; Brazzale & Davison, 2008), none of which are practical for larger matrices and/or multivariate  $\theta$ . Approximation of  $\kappa$  was considered in Barvinok (2010b).

## 2.2 The uniform distribution and model validation

If  $\mathbf{w} \equiv 1$ , or more generally, if  $w_{ij} = \exp(\alpha_i + \beta_j)$  for real-valued  $\alpha$  and  $\beta$ , then  $P^*$  is the uniform distribution on  $\Omega^*$ . The uniform distribution can be used for testing if  $\theta \equiv 0$  in model (2), or equivalently, for model validation of (2) specialized to the case of  $\theta \equiv 0$ . Prominent examples include goodness-of-fit tests for the Rasch model (e.g., Rasch, 1960, 1961; Ponocny, 2001; Chen & Small, 2005; Chen et al., 2005) and for random bipartite graphs and directed graphs without reciprocity (e.g., Wasserman, 1977; Holland & Leinhardt, 1981; Snijders, 1991).

The uniform distribution over  $\Omega^*$  also plays a central role in testing for the presence of interactions in co-occurrence tables, particularly co-occurrence tables arising in ecology, where  $B_{ij}$  indicates the existence of species  $i$  in location  $j$  (e.g., Connor & Simberloff, 1979; Snijders, 1991; Gotelli, 2000; Chen et al., 2005). In these contexts, the uniform distribution subject to the margin totals is taken as a null hypothesis of no interaction among species (without necessarily assuming model (2)). Our original motivation for developing these algorithms came from a similar problem in neuroscience where  $B_{ij}$  indicated whether neuron  $i$  produced an action potential in time bin  $j$ , and the uniform distribution was taken as a null hypothesis of a lack of interaction among the neurons. This can be viewed as an example of conditional testing for multivariate binary time series. In that example,  $mn \approx 10^8$  and existing algorithms in the literature were not practical.

Besides the many statistical applications, when  $\mathbf{w} \equiv 1$ , the normalization constant  $\kappa$  is the number of binary matrices with specified margins, a topic of enduring interest in theoretical computer science (e.g., Kannan et al., 1999; Jerrum et al., 2004; Bezáková et al., 2007) and combinatorial approximation (e.g., Békéssy et al., 1972; McKay, 1984; Greenhill et al., 2006; Canfield et al., 2008; Barvinok, 2010a). Importance sampling algorithms for  $P^*$  can

be used to provide efficient approximations of  $\kappa$  (Blanchet, 2009).

Monte Carlo sampling algorithms for the uniform distribution have been developed by many authors (e.g., Besag & Clifford, 1989; McKay & Wormald, 1990; Snijders, 1991; Rao et al., 1996; Chen et al., 2005; Blanchet, 2009; Bezáková et al., 2007; Chen, 2007; Verhelst, 2008; Bayati et al., 2010). The approach here was inspired by the importance sampling algorithm in Chen et al. (2005), but provides a more principled method for algorithm design that leads to substantial improvements in the uniform case and that also extends to the non-uniform case.

## 2.3 Permanents and permenental processes

If  $\mathbf{w}$  is square and  $\mathbf{r} = \mathbf{c} \equiv 1$ , that is, if  $\Omega^*$  is the set of permutation matrices, then  $\kappa$  is the permanent of  $\mathbf{w}$ , also of enduring interest in theoretical computer science (e.g., Valiant, 1979; Jerrum et al., 2004). A variety of generalizations of permanents and determinants can be expressed as  $\kappa\mu$ , where  $\mu = \mathbb{E}(h(\mathbf{Z}))$  for some function  $h$ , where  $\mathbb{E}$  denotes expected value, and where  $\mathbf{Z}$  has distribution  $P^*$  (e.g., Littlewood, 1950; Vere-Jones, 1988, 1997; Diaconis & Evans, 2000). In principle, the algorithms here could be used to approximate the value of any of these objects, but the practicality of this approach depends heavily on  $h$ . For example, the  $\alpha$ -permanent (Vere-Jones, 1988, 1997) is

$$\text{per}_\alpha(\mathbf{w}) = \kappa \mathbb{E}(\alpha^{\text{cyc}(\mathbf{Z})}), \quad (3)$$

where  $\alpha \in \mathbb{R}$  and  $\text{cyc}(\mathbf{z})$  is the number of disjoint cycles in the permutation corresponding to  $\mathbf{z}$ . The case  $\alpha = 1$  corresponds to the permanent of  $\mathbf{w}$ , and the case  $\alpha = -1$  corresponds to  $(-1)^n$  times the determinant of  $\mathbf{w}$ . Permanents and  $\alpha$ -permanents arise in probability, statistics, and statistical physics in connection to permenental processes and random fields (e.g., Macchi, 1975; Diaconis & Evans, 2000; Shirai & Takahashi, 2003; McCullagh & Møller, 2006; Kou & McCullagh, 2009) and the distribution of order statistics (Vaughan & Venables, 1972; Bapat & Beg, 1989). Our approach is often effective for approximating (3) when  $\alpha > 0$  and  $|\log \alpha|$  is small.

## 3 Algorithm design

### 3.1 The target distribution for the first column

For a matrix  $\mathbf{z} = (z_{ij})$  we use  $\mathbf{z}^j$  to denote the  $j$ th column of  $\mathbf{z}$ , we use  $\mathbf{z}^{j:k}$  to denote the submatrix formed from columns  $j, \dots, k$ , we use  $\mathbf{R}(\mathbf{z}) = (R_i(\mathbf{z}))$  to denote the column vector of row sums defined by  $R_i = \sum_j z_{ij}$ , and we use  $\mathbf{C}(\mathbf{z}) = (C_j(\mathbf{z}))$  to denote the row vector of column sums defined by  $C_j = \sum_i z_{ij}$ . Fix the size of the matrix,  $m \times n$ , the weights  $\mathbf{w}$ , and the margins,  $\mathbf{r}$  and  $\mathbf{c}$ , and let  $\mathbf{Z}$  have distribution  $P^*$  defined in (1) with  $\Omega^* = \{\mathbf{z} \in \{0, 1\}^{m \times n} : \mathbf{R}(\mathbf{z}) = \mathbf{r}, \mathbf{C}(\mathbf{z}) = \mathbf{c}\}$ .

To sample from  $P^*$  we need only design a generic algorithm (generic in  $m, n, \mathbf{r}, \mathbf{c}, \mathbf{w}$ ) for sampling from the distribution  $P$  of the first column, namely,

$$P(\mathbf{x}) = \mathbb{P}(\mathbf{Z}^1 = \mathbf{x}).$$

The reason is that the conditional distribution of  $\mathbf{Z}^{2:n}$  given  $\mathbf{Z}^1$  has the same form as  $P^*$  in (1), but with different parameters. The size of the matrix is now  $m \times (n - 1)$ , the row sums

are updated to  $\mathbf{r} - \mathbf{R}(\mathbf{Z}^1)$ , the column sums are updated to  $\mathbf{c}^{2:n}$ , and the weight matrix is updated to  $\mathbf{w}^{2:n}$ . Once we have sampled  $\mathbf{Z}^1$ , then we can update these parameters and effectively start over, treating the second column of  $\mathbf{Z}$  like it was the first column of the new, updated problem, and then continuing sequentially until we have sampled the entire matrix. This is the same sequential strategy suggested by Chen et al. (2005). The supplementary material (located at the end of this document) contains a more detailed description of this column-wise factorization.

Henceforth, our target distribution is  $P$ , the distribution of  $\mathbf{Z}^1$ . Let  $\mathbf{Y}$  be a random matrix chosen uniformly over  $\Omega^*$ , let  $\Omega$  denote the support of  $\mathbf{Y}^1$ , namely,

$$\Omega = \{\mathbf{x} \in \{0, 1\}^{m \times 1} : \mathbf{x} = \mathbf{z}^1, \mathbf{z} \in \Omega^*\},$$

and for  $\mathbf{x} \in \Omega$  define

$$U(\mathbf{x}) = \mathbb{P}(\mathbf{Y}^1 = \mathbf{x}), \quad V(\mathbf{x}) = \mathbb{E}\left(\prod_{ij} w_{ij}^{Y_{ij}} \mid \mathbf{Y}^1 = \mathbf{x}\right).$$

It is straightforward to verify that

$$P(\mathbf{x}) \propto U(\mathbf{x})V(\mathbf{x})\mathbb{1}\{\mathbf{x} \in \Omega\} \quad (4)$$

for  $\mathbf{x} \in \{0, 1\}^{m \times 1}$ . This factorization conceptually isolates the hard margin constraints,  $\Omega$ , the combinatorics,  $U$ , and the non-uniform weighting,  $V$ . Although the separation is clearly artificial, it is useful to treat each of these factors separately when developing a proposal distribution.

### 3.2 The proposal distribution for the first column

Motivated by the factorization in (4), we consider proposal distributions for the first column of the form

$$Q(\mathbf{x}) \propto \tilde{U}(\mathbf{x})\tilde{V}(\mathbf{x})\mathbb{1}\{\mathbf{x} \in \tilde{\Omega}\},$$

where  $\tilde{U}$  and  $\tilde{V}$  are approximations of  $U$  and  $V$ , respectively, that factor according to

$$\tilde{U}(\mathbf{x}) \propto \prod_{i=1}^m u_i^{x_i}, \quad \tilde{V}(\mathbf{x}) \propto \prod_{i=1}^m v_i^{x_i}$$

for some  $\mathbf{u}, \mathbf{v} \in [0, \infty)^{m \times 1}$ , and where  $\tilde{\Omega}$  is of the form

$$\tilde{\Omega} = \{\mathbf{x} \in \{0, 1\}^{m \times 1} : x_{\pi_i} \in \mathcal{A}_i, \sum_{\ell=1}^i x_{\pi_\ell} \in \mathcal{B}_i, i = 1, \dots, m\} \quad (5)$$

for some permutation  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_m)$  of  $(1, \dots, m)$  and some subsets  $\mathcal{A} = \mathcal{A}_1 \times \dots \times \mathcal{A}_m \subseteq \{0, 1\}^m$  and  $\mathcal{B} = \mathcal{B}_1 \times \dots \times \mathcal{B}_m \subseteq \{0, 1, \dots, c_1\}^m$ . Combining these approximations creates a proposal distribution of the form

$$Q(\mathbf{x}) \propto \prod_{i=1}^m u_i^{x_i} v_i^{x_i} \mathbb{1}\{x_{\pi_i} \in \mathcal{A}_i, \sum_{\ell=1}^i x_{\pi_\ell} \in \mathcal{B}_i\} \quad (\mathbf{x} \in \{0, 1\}^{m \times 1}). \quad (6)$$

Any proposal distribution of this form permits fast, exact sampling and evaluation using  $O(mc_1)$  operations; see Section 3.3. The challenge is to find easily computable choices of

$\mathbf{u}$ ,  $\mathbf{v}$ ,  $\boldsymbol{\pi}$ ,  $\mathcal{A}$ , and  $\mathcal{B}$  such that  $Q$  is a good approximation to the target  $P$ . Fortunately, this seems to be possible in many cases; see Section 4.

For importance sampling to work, the support of  $Q$ , which is a subset of  $\tilde{\Omega}$ , must contain the support of  $P$ , which is a subset of  $\Omega$ . When  $\mathbf{w}$  has no zero entries, we require  $\mathbf{u}$  and  $\mathbf{v}$  to be positive and we engineer  $\tilde{\Omega}$  to exactly coincide with  $\Omega$  so that  $P$  and  $Q$  both have support  $\Omega$ ; see Section 4.1. When  $\mathbf{w}$  does have zero entries, we modify  $\mathbf{v}$  and  $\tilde{\Omega}$  to exclude certain elements in  $\Omega$ , but only elements that are not in the support of  $P$ . This ensures that the support of  $Q$  contains the support of  $P$ , but the supports may no longer be identical. In this case, if the importance sampling algorithm generates a column that is not in the support of  $P$ , then as it sequentially generates additional columns it will eventually try to create a  $Q$  that is identically zero, indicating that no assignment of the current column simultaneously satisfies the margin constraints and has positive weight. At this point the algorithm can assign an importance weight of zero and terminate. Certain patterns of zero weights make the algorithm highly inefficient because the algorithm rarely terminates with a nonzero importance weight. In the supplementary material we discuss alternative choices of  $\mathbf{v}$  and  $\tilde{\Omega}$  that are more efficient for certain patterns of particular interest, including the special case of zeros only on the diagonal.

### 3.3 Efficient sampling and evaluation of the proposal

Let  $\mathbf{X} \in \{0, 1\}^{m \times 1}$  have a distribution  $Q$  that factors according to (6) above for some  $\mathbf{u}$ ,  $\mathbf{v}$ ,  $\boldsymbol{\pi}$ ,  $\mathcal{A}$ , and  $\mathcal{B}$ . Define the permuted partial sums  $\mathbf{S} \in \{0, \dots, c_1\}^{m \times 1}$  according to  $S_i = \sum_{\ell=1}^i X_{\pi_\ell}$  for each  $i$ , and note that  $\mathbf{X}$  and  $\mathbf{S}$  are in bijective correspondence. The distribution of  $\mathbf{S}$  factors according to

$$\mathbb{P}(\mathbf{S} = \mathbf{s}) \propto \prod_{i=1}^m h_i(s_{i-1}, s_i) \quad (7)$$

for  $h_i(s_{i-1}, s_i) = u_{\pi_i}^{s_i - s_{i-1}} v_{\pi_i}^{s_i - s_{i-1}} \mathbb{1}\{s_i - s_{i-1} \in \mathcal{A}_i, s_i \in \mathcal{B}_i\}$ , where here and below we define  $S_0 = s_0 = 0$  for notational convenience.

The factorization in (7) implies that  $\mathbf{S}$  is a Markov chain. If we were given the standard Markov chain representation

$$\mathbb{P}(\mathbf{S} = \mathbf{s}) = \prod_{i=1}^m \mathbb{P}(S_i = s_i | S_{i-1} = s_{i-1}), \quad (8)$$

then generating a random observation of  $\mathbf{S}$  would be trivial. It is known that dynamic programming can be used to convert from Gibbs random field representations like (7) into Bayesian network representations like (8); see, e.g., Frey (1998). The next theorem, which is straightforward to verify (cf. Harrison & Geman, 2009), summarizes dynamic programming in this context.

**Theorem 1.** *Let  $(S_0, S_1, \dots, S_m)$  be a sequence of random variables where each  $S_i$  takes values in the finite set  $D_i$  and where  $D_0 = \{0\}$ . Suppose there exists a sequence of functions  $h_i : D_{i-1} \times D_i \mapsto [0, \infty)$  for  $i = 1, \dots, m$  such that the distribution of  $(S_1, \dots, S_m)$  can be expressed as*

$$\mathbb{P}(S_1 = s_1, \dots, S_m = s_m) \propto \prod_{i=1}^m h_i(s_{i-1}, s_i).$$

Recursively define  $g_m(s_{m-1}, s_m) = h_m(s_{m-1}, s_m)$  and

$$g_i(s_{i-1}, s_i) = h_i(s_{i-1}, s_i) \sum_{s_{i+1} \in D_{i+1}} g_{i+1}(s_i, s_{i+1}) \quad (i = 1, \dots, m-1),$$

where each  $g_i$  is defined over  $D_{i-1} \times D_i$ . Then  $S_0, \dots, S_m$  is a Markov chain and

$$\mathbb{P}(S_i = s_i | S_{i-1} = s_{i-1}) = \frac{g_i(s_{i-1}, s_i)}{\sum_{t \in D_i} g_i(s_{i-1}, t)} \quad (i = 1, \dots, m).$$

In the present context,  $S_i \in \mathcal{B}_i \subseteq \{0, \dots, c_1\}$  for each  $i$ , so the algorithm described in Theorem 1 for converting from (7) to (8) requires at most  $O(mc_1^2)$  operations. In fact, since in the present situation we have  $h_i(s_{i-1}, s_i) = 0$  for  $s_i - s_{i-1} \notin \{0, 1\}$ , implying the same for  $g_i$ , this yields an algorithm that requires  $O(mc_1)$  operations. Instead of representing all  $(c_1 + 1)^2$  combinations of  $(s_{i-1}, s_i)$ , we represent only the  $2c_1 + 1$  feasible combinations. Once the representation in (8) is computed, generating a random observation  $\mathbf{X}$  from  $Q$  or evaluating  $Q(\mathbf{x})$  at any  $\mathbf{x}$  takes  $O(m)$  operations.

## 4 Specification of components

### 4.1 Margin constraints

Here we discuss the construction of  $\tilde{\Omega}$ . In particular, the next theorem shows how to ensure that  $\tilde{\Omega} = \Omega$  for easily computable choices of  $\boldsymbol{\pi}$ ,  $\mathcal{A}$ , and  $\mathcal{B}$ .

**Theorem 2.** (Chen et al., 2005) Assume  $\Omega^* \neq \emptyset$ . Choose  $\boldsymbol{\pi}$  so that  $r_{\pi_1} \geq \dots \geq r_{\pi_m}$ . For each  $i = 1, \dots, m$ , define

$$\mathcal{A}_i = \begin{cases} \{0\} & (r_{\pi_i} = 0); \\ \{0, 1\} & (0 < r_{\pi_i} < n); \\ \{1\} & (r_{\pi_i} = n), \end{cases} \quad \mathcal{B}_i = \begin{cases} \{\max\{0, b_i\}, \dots, c_1\} & (i < m); \\ \{c_1\} & (i = m), \end{cases}$$

for  $b_i = \sum_{\ell=1}^i (r_{\pi_\ell} - \sum_{j=2}^n \mathbb{1}\{c_j \geq \ell\})$ . Define  $\tilde{\Omega}$  according to (5). Then  $\tilde{\Omega} = \Omega$ .

It is instructive to see how these choices of  $\boldsymbol{\pi}$ ,  $\mathcal{A}$ , and  $\mathcal{B}$  ensure that  $\tilde{\Omega} \supseteq \Omega$ , which is the primary requirement for importance sampling. The Gale–Ryser conditions (Gale, 1957; Ryser, 1957) state that there is a binary matrix with margins  $\mathbf{r} \in \{0, \dots, n\}^{m \times 1}$  and  $\mathbf{c} \in \{0, \dots, m\}^{1 \times n}$  if and only if  $\sum_i r_i = \sum_j c_j$  and

$$\sum_{\ell=1}^i r_{\pi_\ell} \leq \sum_{\ell=1}^i \sum_j \mathbb{1}\{c_j \geq \ell\} \quad \text{for all } i = 1, \dots, m-1,$$

where the permutation  $\boldsymbol{\pi}$  is chosen so that  $r_{\pi_1} \geq \dots \geq r_{\pi_m}$ . It is straightforward to see that  $\mathbf{x} \in \{0, 1\}^{m \times 1}$  will be in  $\Omega$  exactly when there is a way to fill out the remaining  $n-1$  columns of the binary matrix that obey the updated margins after accounting for  $\mathbf{x}$ . In other words,  $\mathbf{x} \in \Omega$  exactly when  $\mathbf{r} - \mathbf{x}$  and  $\mathbf{c}^{2:n}$  satisfy the Gale–Ryser conditions for the margins of an  $m \times (n-1)$  binary matrix.

The set  $\mathcal{A}$  is chosen so that  $\mathbf{x} \in \mathcal{A}$  if and only if  $\mathbf{r} - \mathbf{x} \in \{0, \dots, n-1\}^{m \times 1}$ . The set  $\mathcal{B}_m$  is chosen so that  $\sum_{i=1}^m x_i \in \mathcal{B}_m$  if and only if  $\sum_{i=1}^m r_i - x_i = \sum_{j=2}^n c_j$ . This is

equivalent to enforcing the column sum  $\sum_{i=1}^m x_i = c_1$ . Choosing the permutation  $\phi$  so that  $r_{\phi_1} - x_{\phi_1} \geq \dots \geq r_{\phi_m} - x_{\phi_m}$ , the remaining Gale–Ryser conditions are that

$$\sum_{\ell=1}^i (r_{\phi_\ell} - x_{\phi_\ell}) \leq \sum_{\ell=1}^i \sum_{j \geq 2} \mathbb{1}\{c_j \geq \ell\} \quad \text{for all } i = 1, \dots, m-1, \quad (9)$$

which implies that

$$\sum_{\ell=1}^i (r_{\pi_\ell} - x_{\pi_\ell}) \leq \sum_{\ell=1}^i \sum_{j \geq 2} \mathbb{1}\{c_j \geq \ell\} \quad \text{for all } i = 1, \dots, m-1, \quad (10)$$

since the permutation  $\phi$  makes the left side as large as possible. Solving (10) for  $\sum_{\ell=1}^i x_{\pi_\ell}$  gives the bounds encoded in the  $\mathcal{B}_i$  and shows that  $\tilde{\Omega} \supseteq \Omega$ .

We cannot use the permutation  $\phi$  in the construction of  $Q$  because  $\phi$  depends on  $\mathbf{x}$ , however, Chen et al. (2005) further prove that (9) and (10) are in fact equivalent, which means we are in the ideal situation where  $\Omega = \tilde{\Omega}$ . (Although they made use of the factorization in Theorem 2, their proposal distributions were not of the form in (6), except in the special case where  $\Omega = \{\mathbf{x} \in \{0, 1\}^{m \times 1} : \sum_i x_i = c_1\}$ .) Furthermore, Chen (2007) provides an extension of Theorem 2 for the case where, in addition to the margin constraints,  $\Omega^*$  also enforces a fixed pattern of structural zeros for which there is at most one structural zero in each row and column. This includes the important special case of adjacency matrices of directed graphs; see supplementary material.

## 4.2 Combinatorial approximations

Here we discuss approximation of  $U$  by  $\tilde{U}$ . Define

$$N_{m,n}(\mathbf{r}, \mathbf{c}) = |\{\mathbf{z} \in \{0, 1\}^{m \times n} : \mathbf{R}(\mathbf{z}) = \mathbf{r}, \mathbf{C}(\mathbf{z}) = \mathbf{c}\}|$$

to be the number of  $m \times n$  binary matrices with row sums  $\mathbf{r}$  and column sums  $\mathbf{c}$ . We have

$$U(\mathbf{x}) = \mathbb{P}(\mathbf{Y}^1 = \mathbf{x}) = \frac{N_{m,n-1}(\mathbf{r} - \mathbf{x}, \mathbf{c}^{2:n})}{N_{m,n}(\mathbf{r}, \mathbf{c})},$$

(where as before,  $\mathbf{Y}$  is uniform over  $\Omega^*$ ) and we desire an approximation of the form

$$U(\mathbf{x}) = \frac{N_{m,n-1}(\mathbf{r} - \mathbf{x}, \mathbf{c}^{2:n})}{N_{m,n}(\mathbf{r}, \mathbf{c})} \approx \gamma \prod_{i=1}^m u_i^{x_i}, \quad (11)$$

where  $\gamma$  is an irrelevant positive constant.

Temporarily pretending that (11) is accurate for any  $\mathbf{x} \in \{0, 1\}^{m \times 1}$ , we have

$$u_i \approx \frac{U(\mathbf{I}^i)}{U(\mathbf{0})} = \frac{N_{m,n-1}(\mathbf{r} - \mathbf{I}^i, \mathbf{c}^{2:n})}{N_{m,n-1}(\mathbf{r}, \mathbf{c}^{2:n})},$$

where  $\mathbf{I}^i$  is the  $i$ th column of the  $m \times m$  identity matrix  $\mathbf{I}$ . We cannot use this directly, since it is trying to evaluate  $U$  outside of  $\Omega$ , and, furthermore, computationally efficient procedures for evaluating  $N_{m,n-1}$  are not available. Nevertheless, it suggests using

$$u_i = \tilde{N}_{m,n-1}(\mathbf{r} - \mathbf{I}^i, \mathbf{c}^{2:n}) / \tilde{N}_{m,n-1}(\mathbf{r}, \mathbf{c}^{2:n}) \quad (i = 1, \dots, m) \quad (12)$$

for an approximation  $\tilde{N}$  of  $N$  that extends smoothly to invalid margins.

Several asymptotic approximations for  $N$  have appeared in the literature and could be used for  $\tilde{N}$ . For example, Canfield et al. (2008, Theorem 1) suggest the following, which we write asymmetrically with respect to  $\mathbf{r}$  and  $\mathbf{c}$  in order to simplify (13) below:

$$\begin{aligned}\tilde{N}_{m,n}(\mathbf{r}, \mathbf{c}) &= \binom{mn}{\sum_{k=1}^n c_k}^{-1} \prod_{i=1}^m \binom{n}{r_i} \prod_{j=1}^n \binom{m}{c_j} \exp\left(-\frac{1}{2}(1 - \mu_{m,n}(\mathbf{r}, \mathbf{c}))(1 - \nu_{m,n}(\mathbf{c}))\right), \\ \mu_{m,n}(\mathbf{r}, \mathbf{c}) &= \eta_{m,n}(\mathbf{c}) \sum_{i=1}^m \left(r_i - \frac{1}{m} \sum_{k=1}^n c_k\right)^2, \quad \nu_{m,n}(\mathbf{c}) = \eta_{m,n}(\mathbf{c}) \sum_{j=1}^n \left(c_j - \frac{1}{n} \sum_{k=1}^n c_k\right)^2, \\ \eta_{m,n}(\mathbf{c}) &= \frac{mn}{(\sum_{k=1}^n c_k)(mn - \sum_{k=1}^n c_k)}.\end{aligned}$$

Substituting this into (12) and simplifying gives

$$u_i = \frac{r_i}{n - r_i} \exp\left[\eta_{m,n-1}(\mathbf{c}^{2:n})(1 - \nu_{m,n-1}(\mathbf{c}^{2:n}))\left(\frac{1}{2} - r_i + \frac{1}{m} \sum_{k=2}^n c_k\right)\right]. \quad (13)$$

If  $r_i = 0$  or  $r_i = n$ , then the value of  $X_i$  is determined by  $\tilde{\Omega}$ , and any choice of  $u_i > 0$  gives the same  $Q$ ; we use  $u_i = 1$  in these cases. We find that (13) works well over a large range of margins when  $P^*$  is uniform. It is excellent if the margins are approximately semi-regular, that is, if the row and column sums do not deviate substantially from their respective mean values.

For certain pathological cases with wildly varying margins, such as those in Bezáková et al. (2006), (13) does not work well. However, if the margins are such that the resulting matrices have a very low density of ones, even if the margins are highly irregular, then good performance can be obtained by instead using the asymptotic approximation of  $N$  from Greenhill et al. (2006, Theorem 1.3). Details are provided in the supplementary material. In fact, for the specific pathological cases in Bezáková et al. (2006) using this alternative approximation gives  $Q^* \equiv P^*$  in the uniform case. None of the computationally efficient combinatorial approximations that we have found in literature work well when the margins are both highly irregular and lead to a moderate density of ones, but we are hopeful that advances in asymptotic enumeration techniques will eventually lead to approximations that work well in almost all cases.

Chen et al. (2005) observed that combinatorial approximations could be used to find a good choice of  $\mathbf{u}$  and they mentioned an early asymptotic approximation from O’Neil (1969), which was explored further by Blanchet (2009) in an asymptotic analysis of the algorithm. The examples in Chen et al. (2005), however, use  $u_i = r_i/(n - r_i)$ , which is motivated by considering only the row margin constraints. Although there are several substantial differences between their proposal distribution and ours for the special case of the uniform distribution of  $\Omega^*$ , we suspect that much of the improved performance of our algorithm results from using more accurate combinatorial approximations.

In the next section we use  $V$  and  $\tilde{V}$  to account for the effects of  $\mathbf{w}$ , including the effects of zeros in  $\mathbf{w}$ . Since these zeros affect the size of the support of  $P^*$ , an alternative, perhaps more natural approach is to allow  $U$  and  $\tilde{U}$  to capture the effects of zeros in  $\mathbf{w}$ . The supplementary material contains more details.

### 4.3 Non-uniform weighting

Here we discuss approximation of  $V$  by  $\tilde{V}$ . To develop an approximation, we will ignore the column margins and consider only the row margins. This is similar to the approach used by Chen et al. (2005) to develop combinatorial approximations for the uniform case. Let  $\mathbf{B} \in \{0, 1\}^{m \times n}$  be a matrix of independent Bernoulli(1/2) random variables, so that

$$\begin{aligned} V(\mathbf{x}) &= \mathbb{E} \left( \prod_{ij} w_{ij}^{Y_{ij}} \mid \mathbf{Y}^1 = \mathbf{x} \right) = \mathbb{E} \left( \prod_{ij} w_{ij}^{B_{ij}} \mid \mathbf{B}^1 = \mathbf{x}, \mathbf{R}(\mathbf{B}) = \mathbf{r}, \mathbf{C}(\mathbf{B}) = \mathbf{c} \right) \\ &\approx \mathbb{E} \left( \prod_{ij} w_{ij}^{B_{ij}} \mid \mathbf{B}^1 = \mathbf{x}, \mathbf{R}(\mathbf{B}) = \mathbf{r} \right) = \prod_{i=1}^m \mathbb{E} \left( \prod_{j=1}^n w_{ij}^{B_{ij}} \mid B_{i1} = x_i, R_i(\mathbf{B}) = r_i \right) \\ &\propto \prod_{i=1}^m v_i^{x_i}, \end{aligned} \tag{14}$$

where

$$\begin{aligned} v_i &= \frac{\mathbb{E} \left( \prod_{j=1}^n w_{ij}^{B_{ij}} \mid B_{i1} = 1, R_i(\mathbf{B}) = r_i \right)}{\mathbb{E} \left( \prod_{j=1}^n w_{ij}^{B_{ij}} \mid B_{i1} = 0, R_i(\mathbf{B}) = r_i \right)} \quad (r_i = 1, \dots, n-1; i = 1, \dots, m) \\ &= \frac{w_{i1} \binom{n-1}{r_i-1}^{-1} \sum_{\mathbf{b} \in \{0,1\}^{n-1}} \mathbb{1}\{\sum_{j=1}^{n-1} b_j = r_i - 1\} \prod_{j=2}^n w_{ij}^{b_{j-1}}}{\binom{n-1}{r_i}^{-1} \sum_{\mathbf{b} \in \{0,1\}^{n-1}} \mathbb{1}\{\sum_{j=1}^{n-1} b_j = r_i\} \prod_{j=2}^n w_{ij}^{b_{j-1}}}. \end{aligned} \tag{15}$$

In the supplementary material we describe how to compute all possible  $\mathbf{v}$  for all columns using  $O(nd)$  operations (where  $d = \sum_i r_i = \sum_j c_j$ ) in a one-time preprocessing step that can be done prior to sampling. As with  $u_i$ , we always define  $v_i = 1$  if  $r_i = 0$  or  $r_i = n$ . For cases where  $\mathbf{w}$  has zeros, we can sometimes have a zero in the denominator of (15) for  $r_i > 0$ . This happens when fewer than  $r_i$  of the  $n-1$  remaining weights in the row are nonzero. Consequently, we need to force  $X_i = 1$ , which we do by setting the corresponding  $\mathcal{A}_\ell = \{1\}$  in Section 4.1.

An important observation that we have thus far neglected is that many different choices of  $\mathbf{w}$  give rise to the same  $P^*$ . Define

$$\Lambda(\mathbf{w}) = \{\mathbf{t} \in [0, \infty)^{m \times n} : \mathbf{t} = \boldsymbol{\alpha} \boldsymbol{\beta}^\top \circ \mathbf{w}, \boldsymbol{\alpha} \in (0, \infty)^{m \times 1}, \boldsymbol{\beta} \in (0, \infty)^{n \times 1}\},$$

where  $\top$  denotes transpose and  $\mathbf{y} \circ \mathbf{z}$  is the Hadamard product, that is, element-wise multiplication of matrices of the same size, defined by  $(\mathbf{y} \circ \mathbf{z})_{ij} = y_{ij} z_{ij}$ . Then for every  $\mathbf{t} \in \Lambda(\mathbf{w})$  it is straightforward to verify that the  $P^*$  defined with the weight matrix  $\mathbf{w}$  and the  $P^*$  defined with the weight matrix  $\mathbf{t}$  are identical. Similarly, the two versions of  $V$  differ only by an inconsequential constant of proportionality. Unfortunately, our approximation  $\tilde{V}(\mathbf{x}) \propto \prod_i v_i^{x_i}$  defined above does not share this invariance. Consequently, proposal distributions constructed with  $\mathbf{w}$  and  $\mathbf{t}$ , respectively, could differ, even though the target distribution does not differ. We find this unappealing and remedy it in a preprocessing step prior to the construction of  $Q^*$  by first transforming  $\mathbf{w}$  into an equivalent, canonical  $\bar{\mathbf{w}} \in \Lambda(\mathbf{w})$ . In particular,  $\bar{\mathbf{w}}$  is the unique element of  $\Lambda(\mathbf{w})$  with the property that its average nonzero entry over any row or column is one, namely,

$$\sum_{i=1}^m \bar{w}_{ij} = \sum_{i=1}^m \mathbb{1}\{\bar{w}_{ij} > 0\}, \quad \sum_{j=1}^n \bar{w}_{ij} = \sum_{j=1}^n \mathbb{1}\{\bar{w}_{ij} > 0\} \quad (i = 1, \dots, m; j = 1, \dots, n). \tag{16}$$

The solution to (16) over  $\Lambda(\mathbf{w})$  exists, is unique, and is easy to find numerically (Rothblum & Schneider, 1989); see supplementary material for details and more discussion. In the examples below, we always define  $\mathbf{v}$  in terms of  $\bar{\mathbf{w}}$ , not  $\mathbf{w}$ . Not only does this ensure that  $Q^*$  has the same invariance property as  $P^*$ , but we find that performance of the algorithm tends to improve.

If we know that  $P^*$  is uniform over  $\Omega^*$ , for example, if  $\mathbf{w} \equiv 1$  or  $\bar{\mathbf{w}} \equiv 1$ , then  $V$  is constant over  $\Omega$ , and we can ignore  $\mathbf{v}$  in the construction of  $Q$ .

## 5 Importance sampling

### 5.1 Algorithm summary

1. Preprocessing: Compute  $\bar{\mathbf{w}}$  from (16) and precompute all possible  $\mathbf{v}$  using (15) with  $\bar{\mathbf{w}}$  in place of  $\mathbf{w}$ ; see supplementary material for details. Compute  $\boldsymbol{\pi}$ ,  $\mathcal{A}$ , and  $\mathcal{B}$  according to Theorem 2 and  $\mathbf{u}$  according to (13) for the first column.
2. Generating a single observation,  $\tilde{\mathbf{Z}}$ , from  $Q^*$ : The matrix  $\tilde{\mathbf{Z}}$  is generated column-by-column as follows. Set  $q = 1$ . Sequentially, for each column:
  - (a) For the current  $m, n, \mathbf{r}, \mathbf{c}$ , compute  $\boldsymbol{\pi}, \mathcal{A}, \mathcal{B}, \mathbf{u}$  as above. After the preprocessing, updating these quantities based on the previously sampled column requires  $O(m)$  operations.
  - (b) Use Theorem 1 to compute the Markov chain representation for  $Q$ . For the  $j$ th column, this takes  $O(mc_j)$  operations. If  $Q \equiv 0$ , then  $\tilde{\mathbf{Z}}$  will not be in the support of  $P^*$ ; assign a final importance weight of zero and go to step 3.
  - (c) Generate a random observation  $\mathbf{X}$  from  $Q$  and evaluate  $Q(\mathbf{X})$ . This takes  $O(m)$  operations.
  - (d) Assign the current column of  $\tilde{\mathbf{Z}}$  to be  $\mathbf{X}$ . Update  $q \leftarrow qQ(\mathbf{X})$ ,  $\mathbf{c} \leftarrow \mathbf{c}^{2:n}$ ,  $n \leftarrow n - 1$ ,  $\mathbf{r} \leftarrow \mathbf{r} - \mathbf{X}$ . If  $n > 0$ , continue looping over columns; go to step 2a. If  $n = 0$ , the final matrix is  $\tilde{\mathbf{Z}}$ , and we have  $Q^*(\tilde{\mathbf{Z}}) = q$ ; go to step 3.
3. To generate additional independent observations from  $Q^*$ , reset all variables to their original values after step 1 and repeat step 2.

The same algorithm can be used to evaluate  $Q^*(\mathbf{z})$  for any  $\mathbf{z} \in \Omega^*$ . Simply assign  $\mathbf{X}$  to be the current column of  $\mathbf{z}$  in step 2c, instead of sampling a new column. (The algorithm can be applied for any ordering of the columns, and  $Q^*$  will depend on the chosen ordering. The supplementary material describes the heuristics that we use to choose a column ordering.)

### 5.2 Monte Carlo approximation and diagnostics

Let  $\mathbf{Z}$  have distribution  $P^*$ , let  $\mathbf{Z}_1, \dots, \mathbf{Z}_T$  be random sample from  $Q^*$  generated as above, and let  $h$  be a function over  $\Omega^*$ . Define the unnormalized importance weights

$$f(\mathbf{z}) = \frac{\kappa P^*(\mathbf{z})}{Q^*(\mathbf{z})} = \frac{\prod_{ij} w_{ij}^{z_{ij}} \mathbb{1}\{\mathbf{z} \in \Omega^*\}}{Q^*(\mathbf{z})} \quad (\mathbf{z} \in \{0, 1\}^{m \times n}),$$

which we can efficiently evaluate for any  $\mathbf{z}$  as described above. In the formula for  $f(\mathbf{z})$  it is important that we use  $\mathbf{w}$ , not  $\bar{\mathbf{w}}$ , particularly if we are approximating  $\kappa$ . We can

approximate  $\kappa$  and  $\mu = \mathbb{E}(h(\mathbf{Z}))$  via importance sampling in the usual way, namely,

$$\hat{\mu}_T = \frac{\sum_{t=1}^T f(\mathbf{Z}_t)h(\mathbf{Z}_t)}{\sum_{t=1}^T f(\mathbf{Z}_t)} \rightarrow \mu, \quad \hat{\kappa}_T = \frac{1}{T} \sum_{t=1}^T f(\mathbf{Z}_t) \rightarrow \kappa \quad (T \rightarrow \infty). \quad (17)$$

Besides being consistent,  $\hat{\kappa}_T$  and  $\hat{\kappa}_T \hat{\mu}_T$  are also unbiased approximations of  $\kappa$  and  $\kappa\mu$ , respectively. See Liu (2001) for details about importance sampling. See Harrison (2012) for modifications when (17) is used to approximate a p-value.

In this context, importance sampling algorithms are usually evaluated empirically by diagnostics related to the variability of the importance weights. The less variable the importance weights, the better the algorithm is judged to be performing. For the numerical illustrations below, we report

$$\hat{c}v_T^2 = \frac{1}{(T-1)\hat{\kappa}_T^2} \sum_{t=1}^T (f(\mathbf{Z}_t) - \hat{\kappa}_T)^2, \quad \hat{\Delta}_T = \frac{\max_{t=1,\dots,T} f(\mathbf{Z}_t)}{\min_{t=1,\dots,T} f(\mathbf{Z}_t)} - 1.$$

As  $T \rightarrow \infty$ , the approximate squared coefficient of variation,  $\hat{c}v_T^2$ , converges to the true squared coefficient of variation,  $cv^2 = \mathbb{V}(f(\mathbf{Z}_1))/\mathbb{E}(f(\mathbf{Z}_1))^2 = \mathbb{V}(P^*(\mathbf{Z}_1)/Q^*(\mathbf{Z}_1))$ , where  $\mathbb{V}$  denotes variance.  $\tilde{T} = T/(1 + cv^2)$  has been suggested as a rough diagnostic for effective sample size, meaning that a sample size of  $T$  from  $Q^*$  behaves roughly like a sample size of  $\tilde{T}$  from  $P^*$  for the purposes of Monte Carlo approximating  $\mu$  for well-behaved functions  $h$  (Kong et al., 1994; Liu, 2001). For many but not all examples we find  $\hat{c}v_T^2 < 1$ , suggesting that  $Q^*$  is appropriate for efficient importance sampling. The relative range of importance weights reported by  $\hat{\Delta}_T$  is an especially stringent diagnostic. For nearly constant margins and  $P^*$  close to uniform, we often find  $\hat{\Delta}_T \approx 0$ , suggesting that  $Q^*$  is an excellent approximation of  $P^*$ ; see Table 1.

## 6 Numerical illustrations

We experiment with four different classes of weights based on a canonical matrix  $\mathbf{y}$  whose entries are independently sampled from the uniform(0, 1) distribution: (I)  $w_{ij} = 1$ , which is the uniform distribution over  $\Omega^*$ , (II)  $w_{ij} = y_{ij} + 1$ , (III)  $w_{ij} = y_{ij}$ , and (IV)  $w_{ij} = -\mathbb{1}(y_{ij} < 0.99) \log(y_{ij})$ , for all  $i, j$ . The specific entries of  $\mathbf{y}$  for different sized matrices are in the supplementary material. The resulting  $P^*$  is increasingly non-uniform in each of the latter three cases and has 1% structural zeros in case (IV). Recall that each  $\mathbf{w}$  corresponds to a family of weights of the form  $\boldsymbol{\alpha}\boldsymbol{\beta}^T \circ \mathbf{w}$  that give the same  $P^*$  and  $Q^*$ ; see Section 4.3. In all cases we report results with  $T = 1000$ .

We begin with  $500 \times 500$   $r_1$ -regular matrices, i.e.,  $r_i = c_j = r_1$  for all  $i, j = 1, \dots, 500$ . Results are summarized in Table 1 for  $r_1 = 1, 2, 4, 8, \dots, 256$ . The diagnostics are striking, especially in the uniform case, for which the importance weights are essentially constant. Performance degrades slightly as  $P^*$  becomes strongly non-uniform, but in all cases the estimated  $cv^2$  is less than one. Low variability in importance weights corresponds to high precision in estimates of  $\kappa$ . For example, in the uniform case, where  $\kappa = |\Omega^*|$ , for  $r_1 = 256$  we obtain  $\hat{\kappa} = (1.478301 \pm 0.000044) \times 10^{73781}$ , where the errors are approximate standard errors estimated from the same importance samples, and for  $r_1 = 2$  we obtain  $\hat{\kappa}_T = (2.27653 \pm 0.00017) \times 10^{2266}$ , the latter of which is close to the true value of  $\kappa = 2.27658 \dots \times 10^{2266}$ ; see

supplementary material. To our knowledge the exact value of  $\kappa$  in these examples can only be efficiently computed for the special case of the uniform distribution over either 1-regular or 2-regular matrices (Anand et al., 1966). Sampling from the uniform distribution over 1-regular matrices is trivial,  $\kappa = m!$ , and there is no need to use our algorithm, although it is comforting that  $Q^* = P^*$  in this case.

We remark that the distributions corresponding to different weight classes in Table 1 are almost singular with respect to each other. For example, in the 1-regular case, if we use the  $Q^*$  for weight class I as a proposal distribution for the  $P^*$  corresponding to one of the other weight classes, then we obtain, for weight class II,  $\hat{\Delta}_T = 2 \times 10^{13}$  and  $\hat{c}\hat{v}_T^2 = 8 \times 10^2$ , for weight class III,  $\hat{\Delta}_T = 2 \times 10^{67}$  and  $\hat{c}\hat{v}_T^2 = 1 \times 10^3$ , and for weight class IV, only 6 of the 1000 observations from  $Q^*$  were even in the support of  $P^*$ , owing to the structural zeros. Results are similar for other combinations and become even more extreme as  $r_1$  increases.

Table 1:  $500 \times 500$   $r_1$ -regular matrices

$r_1$	uniform		$\mathbf{w}$ class II		$\mathbf{w}$ class III		$\mathbf{w}$ class IV	
	$\hat{\Delta}_T$	$\hat{c}\hat{v}_T^2$	$\hat{\Delta}_T$	$\hat{c}\hat{v}_T^2$	$\hat{\Delta}_T$	$\hat{c}\hat{v}_T^2$	$\hat{\Delta}_T$	$\hat{c}\hat{v}_T^2$
1	0	0	$2 \times 10^{-1}$	$5 \times 10^{-4}$	$4 \times 10^0$	$4 \times 10^{-2}$	$5 \times 10^1$	$3 \times 10^{-1}$
2	$4 \times 10^{-2}$	$5 \times 10^{-6}$	$2 \times 10^{-1}$	$4 \times 10^{-4}$	$6 \times 10^0$	$4 \times 10^{-2}$	$8 \times 10^1$	$2 \times 10^{-1}$
4	$1 \times 10^{-2}$	$1 \times 10^{-6}$	$1 \times 10^{-1}$	$4 \times 10^{-4}$	$5 \times 10^0$	$3 \times 10^{-2}$	$2 \times 10^2$	$2 \times 10^{-1}$
8	$2 \times 10^{-2}$	$1 \times 10^{-6}$	$2 \times 10^{-1}$	$3 \times 10^{-4}$	$3 \times 10^0$	$3 \times 10^{-2}$	$4 \times 10^1$	$2 \times 10^{-1}$
16	$1 \times 10^{-2}$	$1 \times 10^{-6}$	$2 \times 10^{-1}$	$3 \times 10^{-4}$	$3 \times 10^0$	$3 \times 10^{-2}$	$4 \times 10^1$	$1 \times 10^{-1}$
32	$8 \times 10^{-3}$	$8 \times 10^{-7}$	$1 \times 10^{-1}$	$2 \times 10^{-4}$	$2 \times 10^0$	$2 \times 10^{-2}$	$1 \times 10^1$	$1 \times 10^{-1}$
64	$9 \times 10^{-3}$	$9 \times 10^{-7}$	$1 \times 10^{-1}$	$2 \times 10^{-4}$	$3 \times 10^0$	$2 \times 10^{-2}$	$2 \times 10^1$	$9 \times 10^{-2}$
128	$1 \times 10^{-2}$	$9 \times 10^{-7}$	$1 \times 10^{-1}$	$9 \times 10^{-5}$	$1 \times 10^0$	$1 \times 10^{-2}$	$5 \times 10^0$	$5 \times 10^{-2}$
256	$9 \times 10^{-3}$	$9 \times 10^{-7}$	$5 \times 10^{-2}$	$5 \times 10^{-5}$	$1 \times 10^0$	$1 \times 10^{-2}$	$9 \times 10^0$	$7 \times 10^{-2}$

For the special case of 1-regular matrices, corresponding to the first row in Table 1,  $\kappa$  is the permanent of  $\mathbf{w}$  and various generalizations of the permanent correspond to expectations under  $P^*$ . The current state-of-the-art algorithm for approximating permanents and  $\alpha$ -permanents of general matrices, see (3) above, seems to be the importance sampling algorithm of Kou & McCullagh (2009), which has about the same computational complexity as our algorithm. For the case  $\alpha = 1$ , their algorithm is nearly identical to ours, the main differences being the choice of column order and our use of  $\tilde{\mathbf{w}}$ , which seems to give our algorithm slightly better performance. Their algorithm is generally preferable for  $\alpha \neq 1$ , since it is tailored to the specific choice of  $\alpha$ , although in many cases performance is comparable. The supplementary materials have numerical comparisons for each of the  $\mathbf{w}$  used in Table 1 and for all of the examples in Kou & McCullagh (2009), which include cases with  $\alpha = 1$  and  $\alpha = 1/2$ . It is interesting that in many cases our generic approach is competitive with specialized software.

In Table 2 we repeat the simulations of Table 1 for  $50 \times 100$  irregular matrices with margins  $\mathbf{r} = k\tilde{\mathbf{r}}$  and  $\mathbf{c} = k\tilde{\mathbf{c}}$ , for the cases  $k = 1, \dots, 4$ , where  $\tilde{\mathbf{r}}^T = (24^1, 22^2, 17^4, 13^3, 12^2, 11^3, 10^2, 9^3, 8^6, 7^1, 6^4, 5^4, 4^5, 3^6, 2^4)$  and  $\tilde{\mathbf{c}} = (12^2, 10^2, 9^5, 8^4, 7^6, 6^{11}, 5^{10}, 4^{18}, 3^9, 2^{13}, 1^{20})$  using  $i^j$  to denote  $j$  copies of  $i$ . Performance degrades in the irregular case as the matrices become more dense. In most, but not all cases, the diagnostics suggest  $Q^*$  could be used for efficient importance sampling.

For the special case of the uniform distribution, corresponding to the far left category of weights in Tables 1 and 2, the sequential importance sampling algorithm of Chen et al. (2005), as implemented in the publicly available R package `networksis` (Admiraal & Handcock,

Table 2:  $50 \times 100$  irregular matrices with  $\mathbf{r} = k\tilde{\mathbf{r}}$ ,  $\mathbf{c} = k\tilde{\mathbf{c}}$

$k$	uniform		$\mathbf{w}$ class II		$\mathbf{w}$ class III		$\mathbf{w}$ class IV	
	$\hat{\Delta}_T$	$\hat{c}\hat{V}_T^2$	$\hat{\Delta}_T$	$\hat{c}\hat{V}_T^2$	$\hat{\Delta}_T$	$\hat{c}\hat{V}_T^2$	$\hat{\Delta}_T$	$\hat{c}\hat{V}_T^2$
1	$4 \times 10^{-1}$	$1 \times 10^{-3}$	$3 \times 10^0$	$5 \times 10^{-2}$	$8 \times 10^1$	$5 \times 10^{-1}$	$5 \times 10^3$	$3 \times 10^0$
2	$3 \times 10^0$	$3 \times 10^{-2}$	$7 \times 10^0$	$1 \times 10^{-1}$	$7 \times 10^2$	$2 \times 10^0$	$6 \times 10^4$	$7 \times 10^0$
3	$2 \times 10^2$	$7 \times 10^{-1}$	$2 \times 10^2$	$6 \times 10^{-1}$	$2 \times 10^4$	$6 \times 10^0$	$3 \times 10^6$	$4 \times 10^1$
4	$3 \times 10^6$	$3 \times 10^1$	$3 \times 10^6$	$2 \times 10^1$	$4 \times 10^9$	$2 \times 10^2$	$2 \times 10^{13}$	$8 \times 10^2$

2008), appears to be the current state-of-the-art algorithm for practical Monte Carlo approximation. Our algorithm is a substantial improvement, especially for dense or irregular margins. Using `networksis` gives  $\hat{\Delta}_T = (1 \times 10^1, 1 \times 10^2, 2 \times 10^5, 1 \times 10^{11})$  and  $\hat{c}\hat{V}_T^2 = (2 \times 10^{-1}, 6 \times 10^{-1}, 1 \times 10^1, 4 \times 10^2)$  for the first pair of columns in Table 2. The `networksis` implementation is several orders of magnitude slower than our implementation, and is too slow for most of the examples in Table 1.

## Supplementary Material

Supplementary material includes (i) a more detailed description of the column-wise factorization described in Section 3.1, (ii) details about the solution to (16) and other preprocessing of the weights and margins, (iii) alternative combinatorial approximations for sparse matrices with irregular margins, (iv) extensions to Theorem 2 for the case of structural zeros with at most one structural zero in each row and column, including the case of structural zeros along the diagonal, (v) more principled treatments of structural zeros in the approximations to  $U$  and  $V$ , (vi) details for the numerical simulations, (vii) additional numerical illustrations, including examples using real data, and (viii) a Matlab implementation of the algorithm.

## A Column-wise factorization

Define the set of binary matrices with margins  $\mathbf{r} \in \mathbb{N}^{m \times 1}$  and  $\mathbf{c} \in \mathbb{N}^{1 \times n}$  to be

$$\Omega_{m,n}^*(\mathbf{r}, \mathbf{c}) = \{\mathbf{z} \in \{0, 1\}^{m \times n} : \mathbf{R}(\mathbf{z}) = \mathbf{r}, \mathbf{C}(\mathbf{z}) = \mathbf{c}\},$$

and let

$$N_{m,n}(\mathbf{r}, \mathbf{c}) = |\Omega_{m,n}^*(\mathbf{r}, \mathbf{c})|$$

denote the number of such matrices, where  $\mathbb{N} = \{0, 1, \dots\}$  denotes the nonnegative integers. For an  $m \times n$  matrix  $\mathbf{w} \in [0, \infty)^{m \times n}$  define the function

$$P_{m,n}^*(\mathbf{z} \mid \mathbf{r}, \mathbf{c}, \mathbf{w}) = \frac{\mathbb{1}\{\mathbf{z} \in \Omega_{m,n}^*(\mathbf{r}, \mathbf{c})\}}{\kappa_{m,n}(\mathbf{r}, \mathbf{c}, \mathbf{w})} \prod_{i=1}^m \prod_{j=1}^n w_{ij}^{z_{ij}} \quad (\mathbf{z} \in \{0, 1\}^{m \times n})$$

with normalization constant

$$\kappa_{m,n}(\mathbf{r}, \mathbf{c}, \mathbf{w}) = \sum_{\mathbf{z} \in \Omega_{m,n}^*(\mathbf{r}, \mathbf{c})} \prod_{i=1}^m \prod_{j=1}^n w_{ij}^{z_{ij}}$$

using the convention that  $0/0 = 0$ . If  $\kappa_{m,n}(\mathbf{r}, \mathbf{c}, \mathbf{w}) > 0$ , then  $P_{m,n}^*$  is a probability mass function and we use  $\mathbf{Z}$  to denote a random binary matrix with this distribution. We use  $\mathbf{Z}^1$  to denote the first column of  $\mathbf{Z}$ , which has probability mass function

$$P_{m,n}(\mathbf{x} \mid \mathbf{r}, \mathbf{c}, \mathbf{w}) = \mathbb{P}(\mathbf{Z}^1 = \mathbf{x}) = \sum_{\mathbf{z} \in \Omega_{m,n}^*(\mathbf{r}, \mathbf{c})} P_{m,n}^*(\mathbf{z} \mid \mathbf{r}, \mathbf{c}, \mathbf{w}) \mathbb{1}\{\mathbf{z}^1 = \mathbf{x}\} \quad (\mathbf{x} \in \{0, 1\}^{m \times 1}),$$

the support of which is a subset of

$$\Omega_{m,n}(\mathbf{r}, \mathbf{c}) = \{\mathbf{x} \in \{0, 1\}^{m \times 1} : \mathbf{x} = \mathbf{z}^1, \mathbf{z} \in \Omega_{m,n}^*(\mathbf{r}, \mathbf{c})\}.$$

If the entries of  $\mathbf{w}$  are strictly positive, then the support is all of  $\Omega_{m,n}(\mathbf{r}, \mathbf{c})$ .

The algorithmic challenge of sampling the entire matrix  $\mathbf{Z}$  reduces to the challenge of sampling from the first column  $\mathbf{Z}^1$ , because once we have the first column, then we can update the margins and proceed sequentially, treating successive columns like the first. Indeed, it is straightforward to verify that

$$\mathbb{P}(\mathbf{Z}^j = \mathbf{x} \mid \mathbf{Z}^{1:j-1} = \mathbf{y}) = P_{m,n-j+1}(\mathbf{x} \mid \mathbf{r} - \mathbf{R}(\mathbf{y}), \mathbf{c}^{j:n}, \mathbf{w}^{j:n}),$$

so that the generic decomposition  $\mathbb{P}(\mathbf{Z}) = \mathbb{P}(\mathbf{Z}^1) \prod_{j=2}^n \mathbb{P}(\mathbf{Z}^j \mid \mathbf{Z}^{1:j-1})$  gives

$$P_{m,n}^*(\mathbf{z} \mid \mathbf{r}, \mathbf{c}, \mathbf{w}) = P_{m,n}(\mathbf{z}^1 \mid \mathbf{r}, \mathbf{c}, \mathbf{w}) \prod_{j=2}^n P_{m,n-j+1}(\mathbf{z}^j \mid \mathbf{r} - \mathbf{R}(\mathbf{z}^{1:j-1}), \mathbf{c}^{j:n}, \mathbf{w}^{j:n}).$$

(In the main text, we primarily focus on sampling the first column  $\mathbf{Z}_1$ , we suppress  $m, n, \mathbf{r}, \mathbf{c}, \mathbf{w}$  in the notation as much as possible, and we assume that  $\kappa > 0$ .) To summarize, our target distribution is the binary random vector  $\mathbf{Z}_1 \in \{0, 1\}^{m \times 1}$  with probability mass function

$$P(\mathbf{x}) \propto \sum_{\mathbf{z} \in \{0,1\}^{m \times n}} \mathbb{1}\{\mathbf{R}(\mathbf{z}) = \mathbf{r}, \mathbf{C}(\mathbf{z}) = \mathbf{c}, \mathbf{z}^1 = \mathbf{x}\} \prod_{i=1}^m \prod_{j=1}^n w_{ij}^{z_{ij}}.$$

## B Preprocessing the weights and margins

The preprocessing that affects the definition of  $Q^*$  consists of transforming  $\mathbf{w}$  into  $\bar{\mathbf{w}}$  and choosing an ordering of the columns. Other elements of the preprocessing are merely for computational efficiency. All of the preprocessing of  $\mathbf{w}$ , but not reordering the columns, can be skipped when it is known that  $P^*$  is uniform over  $\Omega^*$ , e.g., when  $\bar{\mathbf{w}} \equiv 1$ .

### B.1 Computing $\bar{\mathbf{w}}$ , the solution to equation (16) in the main text

Fix  $\mathbf{w} \in [0, \infty)^{m \times n}$ . Define

$$n_i = \sum_{j=1}^n \mathbb{1}\{w_{ij} > 0\}, \quad m_j = \sum_{i=1}^m \mathbb{1}\{w_{ij} > 0\} \quad (i = 1, \dots, m; j = 1, \dots, n).$$

We are looking for the  $\bar{\mathbf{w}} \in [0, \infty)^{m \times n}$  with the following properties:

$$\bar{w}_{ij} = \alpha_i \beta_j w_{ij}, \quad \sum_{j=1}^n \bar{w}_{ij} = n_i, \quad \sum_{i=1}^m \bar{w}_{ij} = m_j \quad (\alpha_i, \beta_j > 0; i = 1, \dots, m; j = 1, \dots, n).$$

Initializing  $\bar{\mathbf{w}}^{(0)} = \mathbf{w}$  and  $t = 1$ , we iterate the following fixed point equations until convergence:

$$\begin{aligned}\bar{w}_{ij}^{(2t-1)} &= \frac{n_i \bar{w}_{ij}^{(2t-2)}}{\sum_{\ell=1}^n \bar{w}_{i\ell}^{(2t-2)}} \quad (i = 1, \dots, m; j = 1, \dots, n) \\ \bar{w}_{ij}^{(2t)} &= \frac{m_j \bar{w}_{ij}^{(2t-1)}}{\sum_{\ell=1}^m \bar{w}_{\ell j}^{(2t-1)}} \quad (i = 1, \dots, m; j = 1, \dots, n),\end{aligned}$$

where superscripts are indices, and where we take  $0/0 = 0$ . If we iterate this for  $T$  steps, then we use  $\bar{\mathbf{w}} = \bar{\mathbf{w}}^{(2T)}$ . In our experience, a small  $T$  is usually adequate to reach convergence. Since each  $\bar{\mathbf{w}}^{(T)} \in \Lambda(\mathbf{w})$ , iterating to convergence is not important for validity of the algorithm. Rothblum & Schneider (1989) prove existence and uniqueness of  $\bar{\mathbf{w}}$ . They also show that the solution can also be found using a convex programming algorithm, but we have not experimented with this approach.

The computational cost of computing  $\bar{\mathbf{w}}$  takes at least  $O(mn)$  operations, but we do not have theoretical bounds on the computational complexity. In our experience, it can be treated as a negligible preprocessing step. Although this choice of  $\bar{\mathbf{w}}$  outperforms many alternatives, we have found no theoretical justification for its use. It is closely related to Sinkhorn balancing of  $\mathbf{w}$  (Sinkhorn, 1964, 1967), which has appeared in the literature in both algorithmic and theoretical treatments of permanents (e.g., Ando, 1989; Beichl & Sullivan, 1999), and it has the nice property that  $\bar{\mathbf{w}} \equiv 1$  whenever  $\mathbf{w} = \boldsymbol{\alpha}\boldsymbol{\beta}^T$ . In any case,  $P_{m,n}^*(\cdot \mid \mathbf{r}, \mathbf{c}, \bar{\mathbf{w}}) = P_{m,n}^*(\cdot \mid \mathbf{r}, \mathbf{c}, \mathbf{w})$ , so switching from  $\mathbf{w}$  to  $\bar{\mathbf{w}}$  does not change the target distribution.

## B.2 Choosing a column ordering

Our algorithm is not invariant to the ordering of the columns, nor to the pattern of zeros in  $\bar{\mathbf{w}}$ . We use the following heuristic ordering of the columns. First, if  $\bar{\mathbf{w}}$  is banded, then we leave the columns in their original order. The special case of banded weights arises frequently in some applications and we find that the banded ordering works best for accommodating so many zero weights. In other cases, we reorder the columns first by decreasing column sum and then by decreasing variance of the entries of  $\bar{\mathbf{w}}$  within each column. These preprocessing steps, and the accompanying postprocessing steps of returning the columns to their original orders, all require negligible additional computation. In practice, if one is interested in a specific matrix for which these heuristics do not work well, then it can often be advantageous to experiment with different column orders or perhaps swapping the roles of rows and columns. For the description of the algorithm, when referring to the  $j$ th column, we mean the  $j$ th column after any reordering of the columns.

## B.3 Precomputing the constants $v$ for all columns

Define the symmetric polynomials

$$G_n(\mathbf{y}, k) = \sum_{\mathbf{b} \in \{0,1\}^n} \mathbb{1}\{\sum_{j=1}^n b_j = k\} \prod_{j=1}^n y_j^{b_j} \quad (\mathbf{y} \in \mathbb{R}^n, k = 0, \dots, n), \quad (18)$$

and let  $\bar{\mathbf{w}}_i^{j:k} = (\bar{w}_{ij}, \dots, \bar{w}_{ik})$  denote the  $i$ th row of  $\bar{\mathbf{w}}^{j:k}$ . Before sampling we also precompute and store

$$G(i, j, k) = G_{n-j+1}(\bar{\mathbf{w}}_i^{j:n}, k) \quad (19)$$

for all  $i = 1, \dots, m$ ,  $j = 1, \dots, n$ , and  $k = 0, \dots, \min(r_i, n - j + 1)$ . The entire collection can be computed in  $O(nd)$  operations (where  $d = \sum_i r_i = \sum_j c_j$ ), by initializing with  $G(i, j, 0) = 1$ ,  $G(i, n, 1) = \bar{w}_{in}$ , and  $G(i, j, k) = 0$  for all  $i, j$  and  $k > n - j + 1$ , and then using the recursive formula

$$G(i, j, k) = G(i, j + 1, k) + \bar{w}_{ij}G(i, j + 1, k - 1).$$

In particular, in equation (15) in the main text we see that for the first column

$$v_i = \frac{\bar{w}_{i1} \binom{n-1}{r_i-1}^{-1} G(i, 2, r_i - 1)}{\binom{n-1}{r_i}^{-1} G(i, 2, r_i)}.$$

(Recall that we use  $\bar{\mathbf{w}}$ , not  $\mathbf{w}$ , in our implementation of the algorithm.) For the  $j$ th column ( $1 < j < n$ ) we will have

$$v_i = \frac{\bar{w}_{ij} \binom{n-j}{r_i - R_i(\mathbf{z}^{1:j-1}) - 1}^{-1} G(i, j + 1, r_i - R_i(\mathbf{z}^{1:j-1}) - 1)}{\binom{n-j}{r_i - R_i(\mathbf{z}^{1:j-1})}^{-1} G(i, j + 1, r_i - R_i(\mathbf{z}^{1:j-1}))},$$

where  $\mathbf{z}^{1:j-1}$  are the previously sampled columns so that  $\mathbf{r} - \mathbf{R}(\mathbf{z}^{1:j-1})$  are the updated row sums when preparing to sample the  $j$ th column.

## C Alternative combinatorial approximations

For each positive integer  $\ell$  and any nonnegative integer  $a$  we define

$$[a]_\ell = a(a - 1) \cdots (a - \ell + 1),$$

and for a  $k$ -vector  $\mathbf{t}$  of nonnegative integers we define

$$[\mathbf{t}]_\ell = \sum_{i=1}^k [t_i]_\ell.$$

In Section 4.2 of the main text we used a combinatorial approximation due to Canfield et al. (2008), however, other approximations can also be used and may give better performance for some problems. For instance, Greenhill et al. (2006, Theorem 1.3) provide an alternative combinatorial approximation for  $N_{m,n}(\mathbf{r}, \mathbf{c})$  that is accurate, asymptotically, for sparse matrices, except perhaps when the margins are extremely variable:

$$\tilde{N}_{m,n}(\mathbf{r}, \mathbf{c}) = \frac{[\mathbf{c}]_1!}{\prod_{i=1}^m r_i! \prod_{j=1}^n c_j!} \exp(-\alpha_1(\mathbf{c})[\mathbf{r}]_2 - \alpha_2(\mathbf{c})[\mathbf{r}]_3 - \alpha_3(\mathbf{c})[\mathbf{r}]_2^2),$$

$$\alpha_1(\mathbf{c}) = \frac{[\mathbf{c}]_2}{2[\mathbf{c}]_1^2} + \frac{[\mathbf{c}]_2}{2[\mathbf{c}]_1^3} + \frac{[\mathbf{c}]_2^2}{4[\mathbf{c}]_1^4}, \quad \alpha_2(\mathbf{c}) = -\frac{[\mathbf{c}]_3}{3[\mathbf{c}]_1^3} + \frac{[\mathbf{c}]_2^2}{2[\mathbf{c}]_1^4}, \quad \alpha_3(\mathbf{c}) = \frac{[\mathbf{c}]_2}{4[\mathbf{c}]_1^4} + \frac{[\mathbf{c}]_3}{2[\mathbf{c}]_1^4} - \frac{[\mathbf{c}]_2^2}{2[\mathbf{c}]_1^5},$$

where we take  $0/0 = 0$ . Following Section 4.2 in the main text, a straightforward calculation gives

$$u_i = r_i \exp[(r_i - 1)(2\alpha_1(\mathbf{c}^{2:n}) + 3\alpha_2(\mathbf{c}^{2:n})(r_i - 2) + 4\alpha_3(\mathbf{c}^{2:n})([r]_2 - r_i + 1))].$$

This combinatorial approximation is an improvement of an approximation in O’Neil (1969), which was mentioned in Chen et al. (2005) and studied by Blanchet (2009). Both are exactly uniform for the pathological cases in Bezáková et al. (2006); see Supplementary Section E.5.

## D Structural zeros and ones

### D.1 Remarks

The algorithm can be improved to better accommodate structural zeros. We avoided this in the main text to simplify the exposition, but the complexity of the algorithm does not change significantly. The numerical experiments in the main text do not use these improvements, even though some of the examples have structural zeros.

We use the term structural zeros to denote positions  $(i, j)$  such that  $w_{ij} = 0$ , which allows the investigator to explicitly force the binary matrix to zeros at those positions. It is possible that the row and column sums also force some entries to be zero, but we are not referring to those types of implicit structural zeros.

Sometimes it is desirable to force an entry to be one. These structural ones can be accommodated using structural zeros. In a preprocessing step, we replace structural ones with structural zeros and decrement the row and column sums appropriately. Then we sample as usual. In a postprocessing step, we reinsert the structural ones. Henceforth, we only discuss structural zeros.

### D.2 Extensions to Theorem 2 in the main text for zero diagonal

Here we report an extension of Theorem 2 in the main text to the case where  $w$  has at most one zero entry in each row and column. This includes the special case of a zero diagonal, which arises frequently when the binary matrices of interest are adjacency matrices of directed graphs. Unlike the main text, the order of the columns is important for the validity of the algorithm. The columns must be reordered during preprocessing so that  $c_1 \geq \dots \geq c_n$ .

**Theorem 3.** (Chen, 2006, 2007) Assume that  $c_1 \geq \dots \geq c_n$ , fix  $\mathbf{w} \in [0, \infty)^{m \times n}$ , define  $a_{ij} = \mathbf{1}\{w_{ij} > 0\}$  for each  $i, j$ , assume  $R_i(\mathbf{a}) \geq n - 1$  and  $C_j(\mathbf{a}) \geq m - 1$  for each  $i, j$ , and assume that  $\kappa_{m,n}(\mathbf{r}, \mathbf{c}, \mathbf{w}) > 0$ . Choose  $\pi$  so that  $r_{\pi_1} \geq \dots \geq r_{\pi_m}$  and so that whenever  $r_{\pi_i} = r_{\pi_{i+1}}$  we also have  $y_{\pi_i} \leq y_{\pi_{i+1}}$ , where

$$y_i = \begin{cases} \text{the unique } j \text{ such that } w_{ij} = 0 & \text{if there exists such a } j; \\ n + 1 & \text{otherwise.} \end{cases}$$

For each  $i = 1, \dots, m$ , define

$$\mathcal{A}_i = \begin{cases} \{0\} & (a_{\pi_i 1} r_{\pi_i} = 0); \\ \{0, 1\} & (0 < a_{\pi_i 1} r_{\pi_i} < R_{\pi_i}(a)); \\ \{1\} & (a_{\pi_i 1} r_{\pi_i} = R_{\pi_i}(a)), \end{cases} \quad \mathcal{B}_i = \begin{cases} \{\max(0, b_i), \dots, c_1\} & (i < m); \\ \{c_1\} & (i = m), \end{cases}$$

for

$$b_i = \left( \sum_{\ell=1}^i r_{\pi_\ell} \right) - \min_{j=1, \dots, n} \left\{ \sum_{k=j+1}^n c_k + \sum_{\ell=1}^i \sum_{k=2}^j a_{\pi_\ell k} \right\}.$$

Define  $\tilde{\Omega}$  according to (5) in the main text. Then  $\tilde{\Omega}$  is the support of  $P_{m,n}(\cdot \mid \mathbf{r}, \mathbf{c}, \mathbf{w})$ .

### D.3 Alternative treatments of structural zeros

Here we redefine  $\mathcal{A}$ ,  $U$ , and  $V$  from the main text to account for structural zeros differently. We define  $\mathcal{A}$  according to Supplementary Theorem 3 above, which allows trivial cases to be handled by  $\tilde{\Omega}$ .

Let  $\mathbf{Y}$  be a random matrix chosen uniformly over the support of  $P^*$  and define

$$U(\mathbf{x}) = \mathbb{P}(\mathbf{Y}^1 = \mathbf{x}) = \frac{N_{m,n-1}(\mathbf{r} - \mathbf{x}, \mathbf{c}^{2:n})}{N_{m,n}(\mathbf{r}, \mathbf{c})}, \quad V(\mathbf{x}) = \mathbb{E} \left( \prod_{ij} w_{ij}^{Y_{ij}} \mid \mathbf{Y}^1 = \mathbf{x} \right),$$

so that, for any  $\mathbf{x}$ ,  $P(\mathbf{x}) \propto U(\mathbf{x})V(\mathbf{x})\mathbb{1}\{\mathbf{x} \in \Omega\}$ . In the main text, these definition were the same except that  $\mathbf{Y}$  was chosen uniformly over  $\Omega^*$ . If  $\mathbf{w}$  forces structural zeros, then the support of  $P^*$  may be smaller than  $\Omega^*$ . We proceed exactly as in the main text to develop approximations of the new  $U$  and  $V$ .

For the new  $U$ , we follow section 4.2 and note that  $N_{m,n}(\mathbf{r}, \mathbf{c})$  needs to be replaced by the size of the support of  $P_{m,n}^*(\cdot \mid \mathbf{r}, \mathbf{c}, \mathbf{w})$ , say  $N_{m,n}(\mathbf{r}, \mathbf{c}, \mathbf{w})$ , and, consequently,  $\tilde{N}$  needs to be replaced by a combinatorial approximation of the size of the support of  $P^*$ . Greenhill & McKay (2009) provide the modified asymptotic enumeration results corresponding to those that led to equation (13) in the main text. Define  $a_{ij} = \mathbb{1}\{w_{ij} > 0\}$  for all  $i, j$ . Note that  $N_{m,n}(\mathbf{r}, \mathbf{c}, \mathbf{w}) = N_{m,n}(\mathbf{r}, \mathbf{c}, \mathbf{a})$ . For an approximation of  $\tilde{N}_{m,n}(\mathbf{r}, \mathbf{c}, \mathbf{w})$ , Greenhill & McKay (2009, Theorem 2.1) suggest

$$\begin{aligned} \tilde{N}_{m,n}(\mathbf{r}, \mathbf{c}, \mathbf{a}) &= \left( \frac{\sum_{\ell=1}^m \sum_{k=1}^n a_{\ell k}}{\sum_{k=1}^n c_k} \right)^{-1} \prod_{i=1}^m \binom{R_i(\mathbf{a})}{r_i} \prod_{j=1}^n \binom{C_j(\mathbf{a})}{c_j} \\ &\quad \times \exp \left[ -\frac{1}{2} (1 - \mu_{m,n}(\mathbf{r}, \mathbf{c})) (1 - \nu_{m,n}(\mathbf{c})) - \delta_{m,n}(\mathbf{r}, \mathbf{c}, \mathbf{a}) \right], \\ \delta_{m,n}(\mathbf{r}, \mathbf{c}, \mathbf{a}) &= \eta_{m,n}(\mathbf{c}) \sum_{i=1}^m \sum_{j=1}^n \left( (1 - a_{ij}) \left( r_i - \frac{R_i(\mathbf{a})}{mn} \sum_{k=1}^n c_k \right) \left( c_j - \frac{C_j(\mathbf{a})}{mn} \sum_{k=1}^n c_k \right) \right), \end{aligned}$$

which reduces to the formula in the main text when  $\mathbf{a} \equiv 1$ . The functions  $\mu, \nu, \eta$  are defined in the main text. This approximation leads to

$$\begin{aligned} u_i &= \frac{r_i}{R_i(\mathbf{a}^{2:n}) - r_i + 1} \exp \left( \eta_{m,n-1}(\mathbf{c}^{2:n}) \left[ (1 - \nu_{m,n-1}(\mathbf{c}^{2:n})) \left( \frac{1}{2} - r_i + \frac{1}{m} \sum_{k=2}^n c_k \right) \right. \right. \\ &\quad \left. \left. + \sum_{j=2}^n (1 - a_{ij}) \left( c_j - \frac{C_j(\mathbf{a})}{m(n-1)} \sum_{k=2}^n c_k \right) \right] \right), \end{aligned}$$

where, as before, we set  $u_i = 1$  whenever  $r_i = 0$  or  $r_i = R_i(\mathbf{a}^{2:n}) + 1$ .

For the new  $V$ , we follow section 4.3 in the main text, but define  $\mathbf{B}$  to be a matrix of independent Bernoulli random variables where  $B_{ij}$  is Bernoulli( $a_{ij}/2$ ) for  $a_{ij} = \mathbb{1}\{w_{ij} > 0\}$ .

Following equations (14) and (15) from the main text, the first change comes after the second equality in (15), giving

$$\begin{aligned}
v_i &= \frac{w_{i1} \binom{R_i(\mathbf{a}^{2:n})}{r_i-1}^{-1} \sum_{\mathbf{b} \in \{0,1\}^{n-1}} \mathbb{1}\{\sum_{j=1}^{n-1} b_j = r_i - 1\} \prod_{j=2}^n w_{ij}^{b_j-1}}{\binom{R_i(\mathbf{a}^{2:n})}{r_i}^{-1} \sum_{\mathbf{b} \in \{0,1\}^{n-1}} \mathbb{1}\{\sum_{j=1}^{n-1} b_j = r_i\} \prod_{j=2}^n w_{ij}^{b_j-1}} \\
&= \frac{w_{i1} (R_i(\mathbf{a}^{2:n}) - r_i + 1) G_{n-1}(\mathbf{w}_i^{2:n}, r_i - 1)}{r_i G_{n-1}(\mathbf{w}_i^{2:n}, r_i)} = \frac{w_{i1} (R_i(\mathbf{a}^{2:n}) - r_i + 1) G(i, 2, r_i - 1)}{r_i G(i, 2, r_i)},
\end{aligned}$$

where we have made use of the fact that terms inside the summations in the first expression are zero whenever  $\mathbf{b}$  has an entry of one in a place where there is a structural zero. The functions  $G_n$  and  $G$  are defined above in supplementary section B.3. For zeros in the numerator or denominator of the expression for  $v_i$  we set  $v_i = 1$  and allow  $\mathcal{A}_i$  to deterministically choose the appropriate value of the  $i$ th entry. As in the main text, we suggest replacing  $\mathbf{w}$  with  $\bar{\mathbf{w}}$  throughout.

## E Numerical illustrations

### E.1 Pseudorandom number generator

The pseudorandom number generator used by the importance sampling algorithm for the numerical illustrations is the default pseudorandom number generator in Matlab version 7.14, which is the Mersenne twister algorithm `mt19937ar` (c.f. Matsumoto & Nishimura, 1998) described at

<http://www.math.sci.hiroshima-u.ac.jp/~m-mat/MT/emt.html>.

### E.2 Canonical weight matrices

The weight matrices  $\mathbf{w}$  used in Section 6 of the main text are built from a canonical matrix  $\mathbf{y}$ . The  $m \times n$  canonical matrix  $\mathbf{y}$  is constructed as follows:

$$y_{ij} = \frac{R((j-1)m+i)}{2^{31}-1} \quad (i = 1, \dots, m; j = 1, \dots, n),$$

where  $R(0) = 1$  and

$$R(k) = 7^5 R(k-1) \bmod (2^{31}-1) \quad (k = 1, \dots, mn).$$

The sequence  $R(1), R(2), \dots$  is a simple, well-known multiplicative congruential pseudo-random number generator, known as MINSTD, for the discrete uniform distribution over  $\{1, \dots, 2^{31}-2\}$  (Park & Miller, 1988). It was the default pseudorandom number generator in Matlab for many years and is fine for our purpose of creating a matrix  $\mathbf{y}$  with independent uniform(0,1) entries whose values are easy to communicate to others.

### E.3 The number of $n \times n$ two-regular binary matrices

Anand et al. (1966, Eq. (27)) give a simple recursive formula for the number of  $n \times n$  two-regular binary matrices, say  $H_n$ . Initialize  $H_1 = 0$ ,  $H_2 = 1$ ,  $H_3 = 6$ , and then

$$H_k = \frac{1}{2} k(k-1)^2 ((2k-3)H_{k-2} + (k-2)^2 H_{k-3}) \quad (k = 4, 5, \dots).$$

The exact value of  $H_{500}$  can be found in the appendix of this supplement. As noted in Section 6 of the main text, our algorithm provides an extremely accurate approximation.

Chen et al. (2005) used their importance sampling algorithm to approximate  $H_{100}$  as  $(2.96 \pm 0.03) \times 10^{314}$  based on a sample size of 100. For comparison, using a sample of size 100 from our algorithm gives an approximation of  $(2.969 \pm 0.001) \times 10^{314}$ , which appears to be almost 1000 times more efficient for the purposes of approximate enumeration. The true value is  $2.9692 \dots \times 10^{314}$ . The full number can be found in the appendix of this supplement. We should also note that the importance sampling approximations are much more accurate than the combinatorial approximations upon which the importance sampling algorithm is based. For instance, using the approximation  $\tilde{N}$  from Section 4.2 of the main text gives  $2.957 \times 10^{314}$ .

#### E.4 Approximating $\alpha$ -permanents

Here we report comparisons between using our algorithm for approximating  $\alpha$ -permanents and using the custom importance sampling algorithm of Kou & McCullagh (2009). We thank Sam Kou for sharing his code with us. The  $\alpha$ -permanent of  $\mathbf{w}$  can be expressed as

$$\text{per}_\alpha(\mathbf{w}) = \kappa \mathbb{E}(\alpha^{\text{cyc}(\mathbf{Z})}),$$

where  $\mathbf{Z}$  has distribution  $P^*$  with the same  $\mathbf{w}$  and all row and column sums equal to one; see equation (3) in the main text. We approximate it using the consistent, unbiased approximation

$$\hat{\text{per}}_{\alpha,T}(\mathbf{w}) = \hat{\kappa}_T \hat{\mu}_T = \frac{1}{T} \sum_{t=1}^T f(\mathbf{Z}_t) h(\mathbf{Z}_t)$$

for  $h(\mathbf{z}) = \alpha^{\text{cyc}(\mathbf{z})}$ ; see Section 5.2 and equation (17) in the main text.

The Kou & McCullagh algorithm does not attempt to generate  $\mathbf{Z}$  from a distribution that is close to  $P^*$ , like ours does, but rather from a distribution proportional to  $h(\mathbf{z})P^*(\mathbf{z})$ . In the case where  $\alpha = 1$  so that  $h \equiv 1$ , the two approaches agree and the empirical results are quite similar. But when  $\alpha \neq 1$ , their algorithm is generally better, because it is tailored for the choice of  $\alpha$ . Nevertheless, our algorithm might be useful in cases where  $\text{per}_\alpha(\mathbf{w})$  is needed for many  $\alpha$  simultaneously, or in cases where  $\alpha$  is very close to one.

Supplementary Table 3 reports  $\hat{\text{per}}_{\alpha,T}(\mathbf{w})$  along with approximate standard errors defined as  $\hat{\sigma}_T/\sqrt{T}$ , where

$$\hat{\sigma}_T^2 = \frac{1}{T-1} \sum_{t=1}^T (f(\mathbf{Z}_t)h(\mathbf{Z}_t) - \hat{\kappa}_T \hat{\mu}_T)^2.$$

It also reports an approximate relative standard error defined as

$$\hat{\text{rel}}_T = \frac{\hat{\sigma}_T/\sqrt{T}}{\hat{\kappa}_T \hat{\mu}_T} \times 100\%.$$

We use  $T = 1000$  for the examples with  $n = 500$  to match Table 1 in the main text. The other  $\mathbf{w}$  are taken from Kou & McCullagh (2009) and we use  $T = 20000$  to facilitate comparison with their results. In some cases the true value of  $\text{per}_\alpha(\mathbf{w})$  is known and this is shown in the final column of the table; see the supplementary appendix. Except for the  $n = 500$  examples and the results from our algorithm, the entries of Supplementary Table 3 come directly from Table 1 in Kou & McCullagh (2009).

Table 3: Approximating  $\alpha$ -permanents

parameters			our algorithm		Kou & McCullagh		true value
$\mathbf{w}$	$n$	$\alpha$	$\hat{\text{pê}}_{\alpha,T}(\mathbf{w})$	$\hat{\text{rê}}_T\%$	$\hat{\text{pê}}_{\alpha,T}(\mathbf{w})$	$\hat{\text{rê}}_T\%$	$\text{per}_{\alpha}(\mathbf{w})$
I	500	1	$(1.220 \pm 0.000) \times 10^{1134}$	0.00	$(1.220 \pm 0.000) \times 10^{1134}$	0.00	$1.220 \times 10^{1134}$
II	500	1	$(1.437 \pm 0.001) \times 10^{1222}$	0.08	$(1.441 \pm 0.001) \times 10^{1222}$	0.08	?
III	500	1	$(3.998 \pm 0.028) \times 10^{983}$	0.69	$(3.975 \pm 0.033) \times 10^{983}$	0.82	?
IV	500	1	$(3.523 \pm 0.056) \times 10^{1133}$	1.60	$(3.546 \pm 0.066) \times 10^{1133}$	1.85	?
I	500	1/2	$(2.963 \pm 0.167) \times 10^{1132}$	5.62	$(3.078 \pm 0.000) \times 10^{1132}$	0.00	$3.078 \times 10^{1132}$
II	500	1/2	$(3.296 \pm 0.174) \times 10^{1220}$	5.28	$(3.662 \pm 0.012) \times 10^{1220}$	0.32	?
III	500	1/2	$(8.889 \pm 0.459) \times 10^{981}$	5.16	$(1.021 \pm 0.012) \times 10^{982}$	1.21	?
IV	500	1/2	$(9.759 \pm 0.650) \times 10^{1131}$	6.66	$(9.228 \pm 0.228) \times 10^{1131}$	2.47	?
$A_1$	20	1	$(9.800 \pm 0.008) \times 10^{32}$	0.08	$(9.787 \pm 0.014) \times 10^{32}$	0.14	$9.784 \times 10^{32}$
$A_2$	20	1	$(3.513 \pm 0.004) \times 10^{32}$	0.10	$(3.506 \pm 0.007) \times 10^{32}$	0.21	$3.514 \times 10^{32}$
$A_3$	15	1/2	$(1.456 \pm 0.009) \times 10^{22}$	0.60	$(1.437 \pm 0.003) \times 10^{22}$	0.22	$1.439 \times 10^{22}$
$A_4$	15	1/2	$(7.049 \pm 0.044) \times 10^{21}$	0.63	$(7.043 \pm 0.022) \times 10^{21}$	0.32	$7.034 \times 10^{21}$
$A_5$	20	1	$(3.290 \pm 0.003) \times 10^{49}$	0.09	$(3.294 \pm 0.012) \times 10^{49}$	0.38	$3.290 \times 10^{49}$
$A_6$	20	1	$(5.928 \pm 0.024) \times 10^{40}$	0.40	$(5.782 \pm 0.103) \times 10^{40}$	1.72	$5.946 \times 10^{40}$
$A_7$	15	1/2	$(2.069 \pm 0.013) \times 10^{31}$	0.63	$(2.092 \pm 0.008) \times 10^{31}$	0.37	$2.095 \times 10^{31}$
$A_8$	15	1/2	$(1.579 \pm 0.020) \times 10^{25}$	1.27	$(1.549 \pm 0.027) \times 10^{25}$	1.68	$1.579 \times 10^{25}$
$K(x)_9$	9	1/2	$(4.524 \pm 0.036) \times 10^0$	0.79	$(4.504 \pm 0.020) \times 10^0$	0.43	$4.505 \times 10^0$
$K(x)_{11}$	11	1/2	$(1.634 \pm 0.014) \times 10^2$	0.86	$(1.622 \pm 0.009) \times 10^2$	0.56	$1.623 \times 10^2$
$K(x)_{13}$	13	1/2	$(5.815 \pm 0.050) \times 10^3$	0.86	$(5.844 \pm 0.026) \times 10^3$	0.45	$5.816 \times 10^3$
$K(x)_{15}$	15	1/2	$(2.134 \pm 0.019) \times 10^5$	0.89	$(2.117 \pm 0.011) \times 10^5$	0.53	$2.114 \times 10^5$
$K(x)_{100}^{Tr}$	100	1/2	$(1.876 \pm 0.118) \times 10^{-16}$	6.28	$(1.928 \pm 0.037) \times 10^{-16}$	1.90	$1.911 \times 10^{-16}$

## E.5 Additional numerical illustrations for the uniform distribution

Our original interest in these problems was motivated by the uniform distribution over  $\Omega^*$  and we have a variety of simulations investigating this special case. This section is largely reproduced from one of our 2009 unpublished preprints, [arXiv:0906.1004v1](#), which focused on comparing different combinatorial approximations and was the basis for our emphasis on the Canfield et al. (2008) approximation in the main text. The simulations from this section were carried out in 2009 on a MacBook laptop with 2 GB of RAM and a 2.16 GHz dual core processor using Matlab. Everything in this section refers to the uniform distribution with  $\mathbf{w} \equiv 1$ .

Supplementary Table 4 details the speed of the algorithm on  $1000 \times 1000$  binary matrices with all row and column sums identical. These run-times are merely meant to provide a feel for how the algorithm behaves — no attempt was made to control the other processes operating simultaneously on the laptop. Presumably a careful C or assembly language implementation would run much faster. The observed runtime scales closely with the computational complexity of  $O(md)$ . So, for example,  $100 \times 100$   $r_1$ -regular matrices can be sampled about 100 times faster than  $1000 \times 1000$   $r_1$ -regular matrices, and  $10 \times 10$  matrices can be sampled about 10000 times faster.

Table 4: Sampling time per  $1000 \times 1000$   $r_1$ -regular matrix

$r_1$	2	4	8	16	32	64	128	256	512
time (s)	1.2	1.6	2.4	4.0	6.7	12.4	24.4	39.2	46.6

Supplementary Table 5 reports diagnostics on these examples using  $T = 1000$ . We

note that the true number of  $1000 \times 1000$  two-regular matrices is  $1.75147 \dots \times 10^{5133}$ ; see Supplementary Section E.3. The approximation from the first row of Supplementary Table 5 is quite accurate.

Table 5: Performance for the uniform distribution over  $1000 \times 1000$   $r_1$ -regular binary matrices

$r_1$	$\hat{\Delta}_T$	$\hat{cV}_T^2$	$\hat{\kappa}_T$
2	0.049	$4.2 \times 10^{-6}$	$(1.75148 \pm 0.00011) \times 10^{5133}$
4	0.075	$6.4 \times 10^{-6}$	$(7.64296 \pm 0.00061) \times 10^{9910}$
8	0.041	$2.1 \times 10^{-6}$	$(1.01879 \pm 0.00005) \times 10^{18531}$
16	0.008	$3.9 \times 10^{-7}$	$(2.31580 \pm 0.00005) \times 10^{33629}$
32	0.005	$2.3 \times 10^{-7}$	$(6.50167 \pm 0.00010) \times 10^{59218}$
64	0.004	$2.2 \times 10^{-7}$	$(1.22048 \pm 0.00002) \times 10^{100716}$
128	0.004	$1.8 \times 10^{-7}$	$(9.38861 \pm 0.00013) \times 10^{163302}$
256	0.004	$2.2 \times 10^{-7}$	$(6.70630 \pm 0.00010) \times 10^{243964}$
512	0.004	$2.2 \times 10^{-7}$	$(5.02208 \pm 0.00007) \times 10^{297711}$

Bezáková et al. (2006) investigates the performance of the Chen et al. (2005) algorithm on pathological margins with very large  $r_1$  and  $c_1$ , but with all other row and column sums exactly 1. They prove that the Chen et al. (2005) proposal distribution is extremely far from uniform for such margins, too far for importance sampling to be practical. It seems likely that our  $Q^*$  suffers from the same problem, because of the similarities between the combinatorial approximations in each approach. The empirical performance of the  $Q^*$  from the main text is quite bad in these cases; see below. On the other hand, it is straightforward to show that using the combinatorial approximations in Supplementary Section C gives  $Q^* = P^*$  for these types of margins.

Following Bezáková et al. (2006), we experiment with the margins  $\mathbf{r}^T = (240, 1, \dots, 1)$  and  $\mathbf{c} = (179, 1, \dots, 1)$  for a  $240 \times 301$  matrix. By conditioning on the entry in the first row and the first column and then using symmetry, one can see that

$$N_{m,n}(\mathbf{r}, \mathbf{c}) = \binom{300}{240} \binom{239}{179} 60! + \binom{300}{239} \binom{239}{178} 61! = 9.6843 \dots \times 10^{205}.$$

Generating a single observation takes about 0.077 s. Using  $T = 10^5$  gives  $\hat{\Delta}_T = 4.1 \times 10^{11}$ ,  $\hat{cV}_T^2 = 1.7 \times 10^3$ , and  $\hat{\kappa}_T = (2.2 \pm 0.3) \times 10^{205}$ , which is quite bad and highly misleading: approximate 95% confidence intervals created by doubling the standard errors would not come close to covering the true value of  $\kappa$ . Alternatively, using the algorithm with  $u_i$  from Supplementary Section C gives  $\hat{\Delta}_T = \hat{cV}_T^2 = 0$  and  $\hat{\kappa}_T = \kappa = 9.6843 \dots \times 10^{205}$ , since  $Q^* = P^*$  in this case. In most practical examples, however, the algorithm presented in the main text is superior.

Finally, consider Darwin’s finch data (c.f. Chen et al., 2005) which is a  $13 \times 17$  occurrence matrix with  $\mathbf{r}^T = (14, 13, 14, 10, 12, 2, 10, 1, 10, 11, 6, 2, 17)$  and  $\mathbf{c} = (4, 4, 11, 10, 10, 8, 9, 10, 8, 9, 3, 10, 4, 7, 9, 3, 3)$ . A single sample takes about 0.001 s. With  $T = 10^6$ , we find  $\hat{\Delta}_T = 2.8 \times 10^3$  and  $\hat{cV}_T^2 = 0.44$  with  $\hat{\kappa}_T = (6.722 \pm 0.004) \times 10^{16}$ . Chen et al. (2005) report the true value of  $\kappa = 67149106137567626$ , and they also report a  $\hat{cV}_T^2$  of “around one” for their algorithm on this problem. Generally speaking, these importance sampling algorithms tend to be less uniform for small irregular problems like this one, than for the larger and/or more regular examples above.

The previous experiments are based primarily on the internal diagnostics of samples from the proposal distribution  $Q^*$ . Other than the asymptotic analysis in Blanchet (2009) concerning approximate enumeration using a variation of the algorithm of Chen et al. (2005), there are no external checks on the uniformity of  $Q^*$ . Using a complicated, high dimensional proposal distribution without external checks can be dangerous. Indeed, consider the following worst-case scenario. Suppose that  $\Omega^* = E \cup E^c$ , where  $E$  is much smaller than  $E^c$ , and suppose that  $Q^*$  is uniform over each of  $E$  and  $E^c$ , but far from uniform over  $\Omega^*$ , namely,

$$Q^*(\mathbf{z}) = \frac{1 - \epsilon}{|E|} \mathbb{1}\{\mathbf{z} \in E\} + \frac{\epsilon}{|E^c|} \mathbb{1}\{\mathbf{z} \in E^c\} \quad (|E|/|E^c| \ll \epsilon \ll 1).$$

If  $\epsilon$  is extremely tiny, say  $\epsilon = 10^{-100}$ , then Monte Carlo samples from  $Q^*$  will, practically speaking, always lie in  $E$ , which itself is a tiny fraction of  $\Omega^*$ . Furthermore, all internal diagnostics will report that  $Q^*$  is exactly uniform, since it is uniform over  $E$ . But, of course, statistical inferences based on samples from  $Q^*$  will tend to be completely wrong. This section describes two types of experiments designed to provide external checks on the uniformity of  $Q^*$ .

For the first set of experiments we generate a binary matrix  $\mathbf{Z}$  from the uniform distribution over all binary matrices with row sums  $\mathbf{r}$ . This is easy to do by independently and uniformly choosing each row of  $\mathbf{Z}$  from one of the  $\binom{n}{r_i}$  possible configurations. Since the conditional distribution of  $\mathbf{Z}$  given its columns sums  $\mathbf{C}$  is uniform over  $\Omega_{m,n}^*(\mathbf{r}, \mathbf{C})$ , we can view  $\mathbf{Z}$  as a single observation from the uniform distribution over  $\Omega_{m,n}^*(\mathbf{r}, \mathbf{C})$ . Of course, there is no practical way to uniformly and independently generate another such  $\mathbf{Z}$  with the same  $\mathbf{C}$ . Notice that the importance weight  $f(\mathbf{Z})$  gives external information about the uniformity of  $Q^*$  for these margins, since it gives the value of  $1/Q^*$  at a uniformly chosen location in  $\Omega_{m,n}^*(\mathbf{r}, \mathbf{C})$ . Indeed, in the pathological thought experiment described above,  $\mathbf{Z}$  would almost certainly be in  $E^c$  and  $f(\mathbf{Z})$  would be substantially larger than any of the importance weights. Alternatively, if  $Q^*$  is nearly uniform, then  $f(\mathbf{Z})$  should be indistinguishable from the other importance weights. In summary, we can compare  $Q^*$  to  $P^*$  by comparing the importance weights to  $f(\mathbf{Z})$ . This observation can also be used to give valid Monte Carlo p-values with importance sampling, even if the importance sampling distribution is far from the target distribution (Harrison, 2012).

Each experiment of this type proceeds identically. We fix  $m$ ,  $n$ , and  $\mathbf{r}$ . Then we generate  $L$  iid observations, say  $\mathbf{Z}_0^{(1)}, \dots, \mathbf{Z}_0^{(L)}$ , from the uniform distribution over all  $m \times n$  binary matrices with row sums  $\mathbf{r}$ . The column sums of these matrices are  $\mathbf{C}^{(1)}, \dots, \mathbf{C}^{(L)}$ . Then, for each  $\ell = 1, \dots, L$ , we generate  $T$  iid observations, say  $\mathbf{Z}_1^{(\ell)}, \dots, \mathbf{Z}_T^{(\ell)}$ , from the proposal distribution  $Q^*$  over  $\Omega_{m,n}^*(\mathbf{r}, \mathbf{C}^{(\ell)})$ . We compute the ratio of maximum to minimum importance weights *including the original observation* for each  $\ell$ , namely,

$$\hat{\Delta}_T^{(\ell)} = \frac{\max_{t=0, \dots, T} f(\mathbf{Z}_t^{(\ell)})}{\min_{t=0, \dots, T} f(\mathbf{Z}_t^{(\ell)})} - 1,$$

and we report the final summary  $\hat{\Delta}_T^{\max} = \max_{\ell=1, \dots, L} \hat{\Delta}_T^{(\ell)}$ . If  $\hat{\Delta}_T^{\max}$  is close to zero, then this provides evidence that  $Q^*$  is approximately uniform over a large part of each  $\Omega_{m,n}^*(\mathbf{r}, \mathbf{C}^{(\ell)})$ .

We begin with  $1000 \times 1000$  matrices with regular row sums  $r_1 = \dots = r_m$ , but the column sums will not be regular, since they are generated randomly. We use  $L = 10$  and

$T = 10$  for the cases  $r_1 = 2, 8, 32$ , finding  $\hat{\Delta}_T^{\max} = 0.0002, 0.0023, 0.0051$ , respectively. For another example, take the row sums for the irregular  $50 \times 100$  case that was used for Table 2 in the main text and take  $k = 1$ , i.e.,  $\mathbf{r} = \tilde{\mathbf{r}}$ . We use  $L = 100$  and  $T = 1000$  and find that  $\hat{\Delta}_T^{\max} = 1.264$ . These preliminary experiments are encouraging, and suggest that  $Q^*$  is indeed a good approximation of uniform  $P^*$  in many cases.

For the second type of experiment, we try to design an extreme  $\mathbf{z}' \in \Omega^*$  and compare the importance weights to  $f(\mathbf{z}')$ . Again, if  $Q^*$  is approximately uniform over all of  $\Omega^*$  then  $f(\mathbf{z}')$  should be indistinguishable from the other importance weights. For these experiments we report

$$\hat{\Delta}'_T = \frac{\max \{f(\mathbf{z}'), f(\mathbf{Z}_1), \dots, f(\mathbf{Z}_T)\}}{\min \{f(\mathbf{z}'), f(\mathbf{Z}_1), \dots, f(\mathbf{Z}_T)\}} - 1,$$

which should be close to zero if the region in  $\Omega^*$  where  $Q^*$  is approximately uniform includes  $\mathbf{z}'$ .

Consider the regular case where  $m = n = 1000$  and  $r_1 = r_i = c_j$  for all  $i, j$ . Suppose that  $r_1$  evenly divides 1000 and let  $\mathbf{z}'$  be comprised only of disjoint  $r_1 \times r_1$  blocks of ones. In particular, take  $\mathbf{z}'_{ij} = 1$  for  $(k-1)r_1 + 1 \leq i, j \leq kr_1$  and for  $k = 1, \dots, 1000/r_1$ . For the cases  $r_1 = 2, 4, 8$  we compute  $f(\mathbf{z}')$  and compare it to the data that generated the corresponding parts of table 5, obtaining  $\hat{\Delta}'_T = 0.741, 26.24, 6.25 \times 10^4$ , respectively. Clearly,  $Q^*$  is not a uniformly accurate approximation of  $P^*$  over all of  $\Omega^*$  and is unlikely to be useful as a proposal distribution for rejection sampling to get exact samples from  $P^*$ . Nevertheless,  $Q^*$  seems to be extremely well-suited as a proposal distribution for importance sampling. For another example, consider the irregular  $50 \times 100$  case that was used for Table 2 in the main text and take  $k = 1$ , i.e.,  $\mathbf{r} = \tilde{\mathbf{r}}$  and  $\mathbf{c} = \tilde{\mathbf{c}}$ . We construct a pathological  $\mathbf{z}'$  as follows. Place  $c_1$  ones in the last  $c_1$  rows, corresponding to the smallest row sums, of the first column. Place  $c_2$  ones in the last available  $c_2$  rows of the second column, where a row is available if placing a one in that row will not exceed the row sum for that row. Continue in this manner until all the columns are assigned or until a column cannot be assigned successfully. In general, this procedure is not guaranteed to terminate successfully, but it does for this choice of margins. The resulting  $\mathbf{z}'$  is unusual because rows and columns with large sums tend to have zeros at the intersecting entry. Using the data from the corresponding part of Table 2 in the main text gives  $\hat{\Delta}'_T = 14.37$ .

## Supplementary Appendix

### The number of $100 \times 100$ two-regular binary matrices

2969298425 4879211020 5463258948 9046531125 6932010720 0899043082 6661472985 5602957737  
5386603250 7914169840 3947972542 0803105057 9494091210 8196163985 3132939771 8223074880  
1582489734 4113002630 0345104451 5505567811 8301236764 6670284335 5753266570 2919415207  
2361422613 1731302283 4023510256 2089359423 4174989926 4000000000 0000000000 00000

## The number of $500 \times 500$ two-regular binary matrices

2276586004 3872645654 7163822917 6140246378 6529219189 6007058852 1885701633 9308224336  
7024859918 5873168947 8428993358 7710991052 6831024823 1020957186 1359882527 3634597638  
7751901014 9459428637 5300752209 6236400145 2272455600 2450498447 6886449802 2657577100  
8803085437 1426603063 9060350752 5676829379 2441654640 4384402364 9178512515 5701834312  
5382285704 7911170936 9213162976 1124369611 0263144354 2492660647 6317501009 4702298551  
3783877264 5366936440 0850289755 0247749665 4582496735 4778933695 9359401807 4728987947  
4052084791 8351006525 6516882276 6819426986 4276522770 8754690714 8153703130 7689579335  
5313886817 9879619523 6757312609 9563935644 5860973860 5720751902 8525628015 1655464790  
3607836217 2202522127 9381851238 5339132917 8663772909 4697618230 9562268584 1389355037  
4200343275 4426328049 4429348983 4734923700 0635594018 1200043308 9996436581 2082429967  
1420144526 3238392163 0625410465 1147246306 0267066287 2838455102 1984436331 4795820153  
4878729606 4682614593 4828351763 2549610945 2823414530 6966187549 3636469942 1582542169  
0511243887 9654470644 8952801709 4100687806 1803581920 0502635810 6084543151 8196763100  
9226192052 8186323173 8128828855 7307283447 5486503911 0996089630 7969624574 8668199425  
1430690842 9240854111 3288457886 5068062328 1130147009 2410850737 0194640624 5215023611  
0105313331 5631006370 7547904555 8541951209 3762970404 4299114208 6898539174 1261578007  
5271576323 7806458898 5197173413 2333790169 8450503603 6175432120 4646913329 9283772618  
0789892314 7885014128 9831206980 1470933069 2885920165 3886059912 3547627990 2473766270  
0084914243 1261925800 3966112818 5515090740 2869173796 5265773700 6653705150 0776999823  
6682749949 6649629337 6729065663 7740220752 0069908832 1026134189 8109544591 4141299020  
9944691129 8101632276 5735759559 3131678694 4342947732 7389063830 1146871076 6098180223  
5086650691 0193318778 3650834389 5788540935 3233656140 3425148468 8948999361 5539721393  
2767810044 6245991329 5809908199 9005968612 6446584189 0334076925 7082772956 3377889631  
0446650398 8183375905 5124117054 7434261832 8900372657 5745038153 2952534928 4112016395  
9467531245 7165626500 2517876951 1088955612 4288697963 9375087520 6487400471 4382991165  
8206541306 8546637026 9648941941 8803223917 8589969888 6361729999 1147924387 2385375087  
0828596942 2197021633 2700563010 0820849326 1167561772 1388697124 8640000000 0000000000  
0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000  
0000000000 0000000000 00000000

## Exactly computing the $\alpha$ -permanent of a constant matrix

If  $\pi = (\pi_1, \dots, \pi_n)$  is a permutation chosen uniformly at random and  $C$  is the number of disjoint cycles in  $\pi$ , then  $C$  has the same distribution as  $B_1 + \dots + B_n$ , where each  $B_i$  is independent Bernoulli( $1/i$ ) (Durrett, 2010, Lemma 2.2.5). If  $\mathbf{w}$  is an  $n \times n$  constant matrix with common entry  $b$ , then

$$\text{per}_\alpha(\mathbf{w}) = n!b^n \mathbb{E}(\alpha^C) = n!b^n \prod_{i=1}^n \mathbb{E}(\alpha^{B_i}) = n!b^n \prod_{i=1}^n \left( \frac{\alpha}{i} + \left(1 - \frac{1}{i}\right) \right) = n!b^n \prod_{i=1}^n \frac{i + \alpha - 1}{i}.$$

We used this formula with  $n = 500$  and  $b = 1$  to get the true value of  $\text{per}_\alpha(\mathbf{w})$  for  $\mathbf{w}$  in class I in Supplementary Table 3.

## Matlab implementation of the algorithm

This is a place-holder for cleaner, shorter code that will be inserted prior to publication. Software will also be available on the author's website.

```
function [logQ,logP,alist] = BernoulliMarginsRnd(SampN,rN,cN,wN,pflag,wflag,cflag,bIN)
%function [logQ,logP,alist] = BernoulliMarginsRnd(N,r,c,w,pflag,wflag,cflag,Binput)
%
% Approximate sampling from independent Bernoulli random variables B(i,j)
% arranged as an m x n matrix B given the m-vector of row sums r and the
```

```

% n-vector of column sums c, i.e., given that sum(B,2)=r and sum(B,1)=c.
%
% An error is generated if no binary matrix agrees with r and c.
%
% B(i,j) is Bernoulli(p(i,j)) where p(i,j)=w(i,j)/(1+w(i,j)), i.e.,
% w(i,j)=p(i,j)/(1-p(i,j)). [The case p(i,j)=1 must be handled by the user
% in a preprocessing step, by converting to p(i,j)=0 and decrementing the
% row and column sums appropriately.]
%
% Use w=[] for w identically 1, i.e., approximate uniform sampling over
% binary matrices with margins r and c.
%
% N is the sample size. Because of pre-processing, it is more efficient
% per matrix to use larger sample sizes.
%
% alist stores the locations of the ones in the samples.
% If d = sum(r) = sum(c), then alist is 2 x d x N.
%
% The 1-entries of the kth matrix are stored as alist(:,k). The
% (row,column) indices are (alist(1,t,k),alist(2,t,k)) for t=1:d.
%
% If B is the kth matrix, then B can be created from alist via:
%
% B = false(m,n); for t = 1:size(alist,2), B(alist(1,t,k),alist(2,t,k)) = true; end
%
% logQ(k)=log(probability that algorithm generates B)
% logP(k)=log(prod(w(B)))
%
% If the algorithm is used for importance sampling, then the kth
% unnormalized importance weight is exp(logP(k)-logQ(k)).
%
% NOTE for w(i,j)=0:
%
% If the entries of w are not strictly positive, then the algorithm can
% sometimes generate matrices with logP(k)=-inf. In these cases, some of
% the entries of alist(:,k) may be zero and logQ(k) corresponds to the
% probability of generating that particular alist(:,k).
%
% OPTIONS:
%
% pflag: 'canfield' or '' (default, works best in most cases)
%         'greenhill' (perhaps useful for sparse and highly irregular margins)
% pflag controls which combinatorial approximations are used
%
% wflag: 'sinkhorn' or '' (default)
% wflag controls the initial balancing of w; it is passed to canonical.m
%
% cflag: 'descend' or '' (default)
%         'none' (sample columns in original order)
% cflag controls the order in which the columns are sampled
%
% Binput is a m x n binary matrix. If it is provided, then the algorithm
% computes the probability of generating this matrix.

if nargin < 8 || isempty(bin)
    doIN = false;
else
    doIN = true;
end
if nargin < 7 || isempty(cflag)
    cflag = 'descend';
end
if nargin < 6 || isempty(wflag)
    wflag = 'sinkhorn';
end
if nargin < 5 || isempty(pflag)
    pflag = 'canfield';
end
if nargin < 4
    wN = [];
end

doW = true;
if isempty(wN), doW = false; end

doA = true;
if nargout < 2, doA = false; end

if ~isscalar(SampN) || SampN < 1 || SampN ~= round(SampN), error('SampN must be a positive integer'), end

ptype = 0;
switch lower(pflag)
    case 'canfield'
        ptype = 1;
    case 'greenhill'
        ptype = 2;
    otherwise
        error('unknown pflag')
end

```

```

%-----%
%----- START: PREPROCESSING -----%
%-----%

% sizing
mT = numel(rN);
nT = numel(cN);

% sort the marginals (descending)
rT = rN(:);
[rsort,rndxT] = sort(rT,'descend');

if doW
    % balance the weights
    [r,wopt] = canonical(wN,wflag);
    % reorder the columns
    switch lower(cflag)
        case 'none'
            cndx = 1:nT;
        case 'descend'
            [r,cndx] = sortrows(-[cN(:) var(wopt,0,1)']);
        otherwise
            error('unknown cflag')
    end
    csort = cN(cndx);
    wopt = wopt(:,cndx);
    % precompute log weights
    logw = log(wN);

    % -----
    % precompute G

    logwopt = log(wopt);

    rmax = max(rT);
    G = -inf(rmax+1,mT,nT-1);
    G(1, :, :) = 0;
    G(2, :, nT-1) = logwopt(:,nT);

    for i = 1:mT
        ri = rT(i);
        for j = nT-1:-1:2
            wij = logwopt(i,j);
            for k = 2:ri+1
                b = G(k-1,i,j)+wij;
                a = G(k,i,j);
                if a > -inf || b > -inf
                    if a > b
                        G(k,i,j-1) = a + log(1+exp(b-a));
                    else
                        G(k,i,j-1) = b + log(1+exp(a-b));
                    end
                end
            end
        end
    end

    for j = 1:nT-1
        for k = 1:rmax
            Gknum = G(k,i,j);
            Gkden = G(k+1,i,j);
            if isinf(Gkden)
                G(k,i,j) = -1;
            else
                G(k,i,j) = wopt(i,j)*exp(Gknum-Gkden)*((nT-j-k+1)/k);
            end
        end
        if isinf(Gkden)
            G(rmax+1,i,j) = -1;
        end
    end
end
% -----
else
    switch lower(cflag)
        case 'none'
            cndx = 1: numel(cN);
        case 'descend'
            [csort,cndx] = sort(cN(:),'descend');
        otherwise
            error('unknown cflag')
    end
end

% generate the inverse index for the row orders to facilitate fast
% sorting during the updating
irndxT = (1:mT)'; irndxT(rndxT) = irndxT;

% basic input checking
if rsort(1) > nT || rsort(mT) < 0 || csort(1) > mT || csort(nT) < 0 || any(rsort ~= round(rsort)) || any(csort ~= round(csort))
    error('marginal entries invalid')
end

```

```

end

% compute the conjugate of c
cconjT = conjugate_local(csort,mT);

% get the running total of number of ones to assign
countT = sum(rsort);

% get the running total of sum of c squared
ccount2T = sum(csort.^2);
% get the running total of (2 times the) column marginals choose 2
ccount2cT = sum(csort.*(csort-1));
% get the running total of (6 times the) column marginals choose 3
ccount3cT = sum(csort.*(csort-1).*(csort-2));

% get the running total of sum of r squared
rcount2T = sum(rsort.^2);
% get the running total of (2 times the) row marginals choose 2
rcount2cT = sum(rsort.*(rsort-1));
% get the running total of (6 times the) row marginals choose 3
rcount3cT = sum(rsort.*(rsort-1).*(rsort-2));

% check for compatible marginals
if countT ~= sum(csort) || any(cumsum(rsort) > cumsum(cconjT)), error('marginal sums invalid'), end

% initialize the memory
logQ = zeros(SampN,1);
logP = zeros(SampN,1);
if doA, AN = SampN; else AN = 1; end
alist = zeros(2,countT,AN);
% initialize the memory
M = csort(1)+3; % index 1 corresponds to -1; index 2 corresponds to 0, index 3 corresponds to 1, ..., index M corresponds to c(1)+1
S = zeros(M,mT);
SS = zeros(M,1);

eps0 = eps(0); % used to prevent divide by zero

%-----%
%----- END: PREPROCESSING -----%
%-----%

% loop over the number of samples
for SampLoop = 1:SampN

    %----- INITIALIZATION -----%
    if doA, ALoop = SampLoop; else ALoop = 1; end

    % copy in initialization
    r = rT;
    rndx = rndxT;
    irndx = irndxT;

    cconj = cconjT;
    count = countT;
    ccount2 = ccount2T;
    ccount2c = ccount2cT;
    ccount3c = ccount3cT;
    rcount2 = rcount2T;
    rcount2c = rcount2cT;
    rcount3c = rcount3cT;
    m = mT;
    n = nT;

    % initialize
    place = 0; % most recent assigned column in alist
    logq = 0; % running log probability
    logp = 0;

    %-----%
    %----- START: COLUMN-WISE SAMPLING -----%
    %-----%

    %----- loop over columns -----%
    for c1 = 1:nT

        %-----%
        %----- START: SAMPLE THE NEXT "COLUMN" -----%
        %-----%

        % remember the starting point for this columns
        placestart = place + 1;

        %-----%
        % sample a col
        %-----%

        label = cndx(c1); % current column label

        colval = csort(c1); % current column value

```

```

if colval == 0 || count == 0, break, end

% update the conjugate
for i = 1:colval
    cconj(i) = cconj(i)-1;
end
% update the number of columns remaining
n = n - 1;

%----- DP initialization -----

smin = colval;
smax = colval;
cumsums = count;
% update the count
count = count - colval;
% update running total of sum of c squared
ccount2 = ccount2 - colval^2;
% update the remaining (two times the) sum of column sums choose 2
ccount2c = ccount2c - colval*(colval-1);
% update the remaining (six times the) sum of column sums choose 3
ccount3c = ccount3c - colval*(colval-1)*(colval-2);

cumconj = count;

SS(colval+3) = 0;
SS(colval+2) = 1;
SS(colval+1) = 0;

% get the constants for computing the probabilities
% it is faster to compute them all, than to check pflag
d = ccount2c/count^2;
if (count == 0) || (m*n == count)
    weightA = 0;
else
    weightA = m*n/(count*(m*n-count));
    weightA = weightA*(1-weightA*(ccount2-count^2/n))/2;
end

d2 = ccount2c/(2*count^2+eps0) + ccount2c/(2*count^3+eps0) + ccount2c^2/(4*count^4+eps0);
d3 = -ccount3c/(3*count^3+eps0) + ccount2c^2/(2*count^4+eps0);
d22 = ccount2c/(4*count^4+eps0) + ccount3c/(2*count^4+eps0) - ccount2c^2/(2*count^5+eps0);

%----- dynamic programming -----
SSS = 0;
% loop over (remaining and sorted descending) rows in reverse
for i = m:-1:1

    % get the value of this row and use it to compute the
    % probability of a 1 for this row/column pair
    rlabel = rndx(i);
    val = r(rlabel);
    if ptype == 1
        % canfield
        p = val*exp(weightA*(1-2*(val-count/m)));
        p = p./(m+1-val*p);
        q = 1-p;
    elseif ptype == 2
        % greenhill
        q = 1/(1+val*exp((2*d2+3*d3*(val-2)+4*d22*(rccount2c-val+1))*(val-1)));
        p = 1-q;
    else
        % never get here
        p = 0; q = 0; % helps compiler
    end

    % incorporate weights
    if doW && n > 0 && val > 0
        Gk = G(val,rlabel,c1);
        if Gk < 0
            q = 0;
        else
            p = p*Gk;
        end
    end

    % update the feasibility constraints
    cumsums = cumsums - val;
    cumconj = cumconj - cconj(i);

    sminold = smin;
    smaxold = smax;

    % incorporate the feasibility constraints into bounds on the
    % running column sum
    smin = max(0,max(cumsums-cumconj,sminold-1));
    smax = min(smaxold,i-1);

    % DP iteration
    SSS = 0;

```

```

SS(smin+1) = 0; % no need to set S(1:smin) = 0, since it is not accessed
for j = smin+2:smax+2
    a = SS(j)*q;
    b = SS(j+1)*p;
    apb = a + b;
    SSS = SSS + apb;
    SS(j) = apb;
    S(j,i) = b/(apb+eps0);
end
SS(smax+3) = 0; % no need to set S(smax+4:end) = 0, since it is not accessed

% check for impossible
if SSS <= 0, break, end

% normalize to prevent overflow/underflow
for j = smin+2:smax+2
    SS(j) = SS(j) / SSS;
end

end

% check for impossible
if SSS <= 0, logp = -inf; break, end

%----- sampling -----
j = 2; % running total (offset to match indexing offset)
jmax = colval + 2;
if j < jmax % skip assigning anything when colval == 0
    if doIN
        for i = 1:m
            % get the transition probability of generating a one
            p = S(j,i);
            % get the current row
            rlabel = rndx(i);
            if bIN(rlabel,label)

                % if we have a generated a one, then decrement the current
                % row total
                val = r(rlabel);
                r(rlabel) = val-1;

                rcount2 = rcount2 - 2*val + 1;
                rcount2c = rcount2c - 2*val + 2;
                rcount3c = rcount3c - 3*(val-1)*(val-2);

                % record the entry and update the log probability
                place = place + 1;
                logq = logq + log(p);
                if doW, logp = logp + logw(rlabel,label); end
                alist(1,place,ALoop) = rlabel;
                alist(2,place,ALoop) = label;
                j = j + 1;
                % the next test is not necessary, but seems more efficient
                % since all the remaining p's must be 0
                if j == jmax, break, end
            else
                logq = logq + log(1-p);
            end
        end
    else
        for i = 1:m
            % get the transition probability of generating a one
            p = S(j,i);
            if rand <= p

                % if we have a generated a one, then decrement the current row total
                rlabel = rndx(i);
                val = r(rlabel);
                r(rlabel) = val-1;

                rcount2 = rcount2 - 2*val + 1;
                rcount2c = rcount2c - 2*val + 2;
                rcount3c = rcount3c - 3*(val-1)*(val-2);

                % record the entry and update the log probability
                place = place + 1;
                logq = logq + log(p);
                if doW, logp = logp + logw(rlabel,label); end
                alist(1,place,ALoop) = rlabel;
                alist(2,place,ALoop) = label;
                j = j + 1;
                % the next test is not necessary, but seems more efficient
                % since all the remaining p's must be 0
                if j == jmax, break, end
            else
                logq = logq + log(1-p);
            end
        end
    end
end
end

```

```

end

%-----%
%----- END: SAMPLE THE NEXT "COLUMN" -----%
%-----%

if count == 0, break, end

%-----%
% everything is updated except the sorting
%-----%

%-----%
%----- START: RESORT THE NEW ROW SUMS -----%
%-----%

% re-sort the assigned rows

% this code block takes each row that was assigned to the list
% and either leaves it in place or swaps it with the last row
% that matches its value; this leaves the rows sorted (descending)
% since each row was decremented by only 1

% looping in reverse ensures that least rows are swapped first
for j = place:-1:placestart
    % get the row label and its new value (old value -1)
    k = alist(1,j,ALoop);
    val = r(k);
    % find its entry in the sorting index
    irndxk = irndx(k);
    % look to see if the list is still sorted
    irndxk1 = irndxk + 1;
    if irndxk1 > m || r(rndx(irndxk1)) <= val
        % no need to re-sort
        continue;
    end
    % find the first place where k can be inserted
    irndxk1 = irndxk1 + 1;
    while irndxk1 <= m && r(rndx(irndxk1)) > val
        irndxk1 = irndxk1 + 1;
    end
    irndxk1 = irndxk1 - 1;
    % now swap irndxk and irndxk1
    rndxk1 = rndx(irndxk1);
    rndx(irndxk) = rndxk1;
    rndx(irndxk1) = k;
    irndx(k) = irndxk1;
    irndx(rndxk1) = irndxk;
end

%-----%
%----- END: RESORT THE NEW ROW SUMS -----%
%-----%

% r(rndx(rndx1:rndxm)) is sorted descending and has exactly those
% unassigned rows
% rndx(rndx1:rndxm) still gives the labels of those rows
% rndx(irndx(k)) = k
%
% c(c1+1:cn) is sorted descending and has exactly those unassigned columns
% cndx(c1+1:cn) still gives the labels of those columns
%
% m, n, count, ccount2, ccount2c are valid for the remaining rows, cols

end

logQ(SampLoop) = logq;
logP(SampLoop) = logp;

end

%-----%
%-----%
%-----%
%----- END OF MAIN FUNCTION -----%
%-----%
%-----%

% helper function (just to keep everything together... not for efficiency,
% since it is only called once)

function cc = conjugate_local(c,n)
% function cc = conjugate(c,n)
%
% let c(:) be nonnegative integers
% cc(k) = sum(c >= k) for k = 1:n

cc = zeros(n,1);

```

```

%c = min(c,n);

for j = 1:numel(c)
    k = c(j);
    if k >= n
        cc(n) = cc(n) + 1;
    elseif k >= 1
        cc(k) = cc(k) + 1;
    end
end

s = cc(n);
for j = n-1:-1:1
    s = s + cc(j);
    cc(j) = s;
end

%-----
function [a,b,abw,k] = canonical(w,flag,tol,maxiter,r,c)

[m,n] = size(w);

if nargin <6 || isempty(c)
    c = ones(1,n);
elseif size(c,1) ~= 1
    c = c(:).';
end
if nargin <5 || isempty(r)
    r = ones(m,1);
elseif size(r,2) ~= 1
    r = r(:);
end
if nargin <4 || isempty(maxiter)
    maxiter = 10^5;
end
if nargin <3 || isempty(tol)
    tol = 1e-8;
end
if nargin <2 || isempty(flag)
    flag = 'sinkhorn';
end

switch lower(flag)

    case 'sinkhorn'

        M = sum(w>0,1); N = sum(w>0,2);
        a = N./sum(w,2); a = a/mean(a);
        b = M./sum(bsxfun(@times,a,w),1);

        if tol >= 0, a0 = a; b0 = b; end

            k = 0;
            tolcheck = inf;
            while k < maxiter && tolcheck > tol
                k = k + 1;

                a = N./sum(bsxfun(@times,b,w),2); a = a/mean(a);
                b = M./sum(bsxfun(@times,a,w),1);

                if tol >= 0
                    tolcheck = sum(abs(a-a0))+sum(abs(b-b0));
                a0 = a; b0 = b;
            end

            case 'sinkhorn-col'

                w = fliplr(w);

                M = sum(w>0,1); N = cumsum(w>0,2);
                aa = N./cumsum(w,2);
                b = M./sum(w.*aa,1); b = b / mean(b);
                a = aa(:,n);

                if tol >= 0, a0 = a; b0 = b; end

                    k = 0;
                    tolcheck = inf;
                    while k < maxiter && tolcheck > tol
                        k = k + 1;

                        aa = N./cumsum(bsxfun(@times,b,w),2);
                        b = M./sum(w.*aa,1); b = b / mean(b);
                        a = aa(:,n);

                        if tol >= 0
                            tolcheck = sum(abs(a-a0))+sum(abs(b-b0));
                        a0 = a; b0 = b;
                    end
                end
            end
        end
    end
end

```

```

end
end

w = fliplr(w);
b = fliplr(b);

case 'log'

w0 = w > 0;
M = sum(w0,1); N = sum(w0,2);
logw = log(w+~w0);
a = exp(-sum(logw,2) ./ N);
b = exp(-sum(logw,1) ./ M);

case 'entropy'

w1 = (w > 0) ./ max(w,eps(0));
a = sqrt(sum(w1,2) ./ sum(w,2)); a = a/mean(a);
b = sqrt(sum(bsxfun(@rdivide,w1,a),1) ./ sum(bsxfun(@times,a,w),1));

if tol >= 0, a0 = a; b0 = b; end

    k = 0;
    tolcheck = inf;
    while k < maxiter && tolcheck > tol
        k = k + 1;

a = sqrt(sum(bsxfun(@rdivide,w1,b),2) ./ sum(bsxfun(@times,b,w),2)); a = a/mean(a);
b = sqrt(sum(bsxfun(@rdivide,w1,a),1) ./ sum(bsxfun(@times,a,w),1));

if tol >= 0
    tolcheck = sum(abs(a-a0))+sum(abs(b-b0));
a0 = a; b0 = b;
end
    end

case 'l2'

w2 = w.^2;

a = sum(w,2) ./ sum(w2,2); a = a/mean(a);
b = sum(bsxfun(@times,a,w),1) ./ sum(bsxfun(@times,a.^2,w2),1);

if tol >= 0, a0 = a; b0 = b; end

    k = 0;
    tolcheck = inf;
    while k < maxiter && tolcheck > tol
        k = k + 1;

a = sum(bsxfun(@times,b,w),2) ./ sum(bsxfun(@times,b.^2,w2),2); a = a/mean(a);
b = sum(bsxfun(@times,a,w),1) ./ sum(bsxfun(@times,a.^2,w2),1);

if tol >= 0
    tolcheck = sum(abs(a-a0))+sum(abs(b-b0));
a0 = a; b0 = b;
end
    end

    case 'l2p'

        w2 = w.^2;

        c = (1+w).^3;
a = sum(w ./ c,2) ./ sum(w2 ./ c,2); a = a/mean(a);
c = (1+bsxfun(@times,a,w)).^3;
b = sum(bsxfun(@times,a,w) ./ c,1) ./ sum(bsxfun(@times,a.^2,w2) ./ c,1);

if tol >= 0, a0 = a; b0 = b; end

        k = 0;
        tolcheck = inf;
        while k < maxiter && tolcheck > tol
            k = k + 1;

            c = (1+a*b.*w).^3;
a = sum(bsxfun(@times,b,w) ./ c,2) ./ sum(bsxfun(@times,b.^2,w2) ./ c,2); a = a/mean(a);
c = (1+a*b.*w).^3;
b = sum(bsxfun(@times,a,w) ./ c,1) ./ sum(bsxfun(@times,a.^2,w2) ./ c,1);

if tol >= 0
    tolcheck = sum(abs(a-a0))+sum(abs(b-b0));
a0 = a; b0 = b;
end
        end

    case 'ratio'

        wz = w > 0;
        w(~wz) = eps(0);

```

```

    a = sqrt(sum(wz./w,2)./sum(w,2)); a = a/mean(a);
    b = sqrt(sum(wz./(bsxfun(@times,a,w),1)./sum(bsxfun(@times,a,w),1)));

    if tol >= 0, a0 = a; b0 = b; end

    k = 0;
    tolcheck = inf;
    while k < maxiter && tolcheck > tol
        k = k + 1;

    a = sqrt(sum(wz./(bsxfun(@times,b,w),2)./sum(bsxfun(@times,b,w),2)); a = a/mean(a);
        b = sqrt(sum(wz./(bsxfun(@times,a,w),1)./sum(bsxfun(@times,a,w),1)));

    if tol >= 0
        tolcheck = sum(abs(a-a0))+sum(abs(b-b0));
    a0 = a; b0 = b;
    end
    end

    case 'barvinok'

    s = log(r/n);
    t = log(c/m);

    M = w.*(exp(s)*exp(t));
    M = M ./ (1+M);

    sMr = sum(M,2)-r;
    sMc = sum(M,1)-c;

    tolcheck = sum(abs(sMr))+sum(abs(sMc));

    alpha = .01;

    while tolcheck > tol

    s = s - alpha*sMr;
    t = t - alpha*sMc;

    M = w.*(exp(s)*exp(t));
    M = M ./ (1+M);

    sMr = sum(M,2)-r;
    sMc = sum(M,1)-c;

    tolcheck = sum(abs(sMr))+sum(abs(sMc));
    end

    a = exp(s);
    b = exp(t);

    otherwise

        error('unknown flag')
    end

    if nargout > 2, abw = a*b.*w; end

```

## References

- Admiraal, R., & Handcock, M. S. (2008). networksis: a package to simulate bipartite graphs with fixed marginals through sequential importance sampling. *J. Statist. Software*, 24(8), 1–21.
- Anand, H., Dumir, V. C., & Gupta, H. (1966). A combinatorial distribution problem. *Duke Math. J.*, 33(4), 757–769.
- Ando, T. (1989). Majorization, doubly stochastic matrices, and comparison of eigenvalues. *Linear Algebra Appl.*, 118, 163–248.
- Bapat, R. B., & Beg, M. I. (1989). Order statistics for nonidentically distributed variables and permanents. *Sankhyā Ser. A*, 51, 79–93.
- Barvinok, A. (2010a). Matrices with prescribed row and column sums. *Linear Algebra Appl.*, 436, 820–834.

- Barvinok, A. (2010b). On the number of matrices and a random matrix with prescribed row and column sums and 0-1 entries. *Adv. Math.*, 224(1), 316–339.
- Bayati, M., Kim, J. H., & Saberi, A. (2010). A sequential algorithm for generating random graphs. *Algorithmica*, 58(4), 860–910.
- Beichl, I., & Sullivan, F. (1999). Approximating the permanent via importance sampling with application to the dimer covering problem. *J. Comput. Phys.*, 149(1), 128–147.
- Békéssy, A., Bekessy, P., & Komlós, J. (1972). Asymptotic enumeration of regular matrices. *Stud. Sci. Math. Hungar.*, 7, 343–353.
- Besag, J., & Clifford, P. (1989). Generalized Monte Carlo significance tests. *Biometrika*, 76(4), 633–642.
- Bezáková, I., Bhatnagar, N., & Vigoda, E. (2007). Sampling binary contingency tables with a greedy start. *Random Struct. Algor.*, 30, 168–205.
- Bezáková, I., Sinclair, A., Štefankovič, D., & Vigoda, E. (2006). Negative examples for sequential importance sampling of binary contingency tables. In Y. Azar & T. Erlebach (Eds.), *Algorithms – ESA 2006* (Vol. 4168, p. 136–147). Berlin/Heidelberg: Springer.
- Blanchet, J. H. (2009). Efficient importance sampling for binary contingency tables. *Ann. Appl. Probab.*, 19(3), 949–982.
- Brazzale, A. R. (2005). *hoa*: An R package bundle for higher order likelihood inference. *Rnews*, 5, 20–27. (ISSN 609-3631)
- Brazzale, A. R., & Davison, A. C. (2008). Accurate parametric inference for small samples. *Statist. Sci.*, 23(4), 465–484.
- Canfield, E. R., Greenhill, C., & McKay, B. D. (2008). Asymptotic enumeration of dense 0–1 matrices with specified line sums. *J. Comb. Theory A*, 115(1), 32–66.
- Chen, Y. (2006). Simple existence conditions for zero-one matrices with at most one structural zero in each row and column. *Discrete Math.*, 306(22), 2870–2877.
- Chen, Y. (2007). Conditional inference on tables with structural zeros. *J. Comput. Graph. Stat.*, 16(2), 445–467.
- Chen, Y., Diaconis, P., Holmes, S. P., & Liu, J. S. (2005). Sequential Monte Carlo methods for statistical analysis of tables. *J. Am. Statist. Assoc.*, 100(469), 109–120.
- Chen, Y., & Small, D. (2005). Exact tests for the Rasch model via sequential importance sampling. *Psychometrika*, 70(1), 11–30.
- Connor, E. F., & Simberloff, D. (1979). The assembly of species communities: chance or competition? *Ecology*, 60, 1132–1140.
- Cox, D. R. (1958). The regression analysis of binary sequences. *J. R. Statist. Soc. B*, 20, 215–242.
- Cytel. (2010). *LogXact 9*. Cambridge, MA: Cytel Inc.

- Diaconis, P., & Evans, S. N. (2000). Immanants and finite point processes. *J. Comb. Theory A*, *91*(1-2), 305–321.
- Durrett, R. (2010). *Probability: theory and examples* (4th ed.). New York: Cambridge Univ. Pr.
- Fienberg, S., Meyer, M., & Wasserman, S. (1985). Statistical analysis of multiple sociometric relations. *Journal of the American Statistical Association*, *80*(389), 51–67.
- Frey, B. J. (1998). *Graphical models for machine learning and digital communication*. Cambridge, MA: MIT Press.
- Gale, D. (1957). A theorem on flows in networks. *Pac. J. Math.*, *7*, 1073–1082.
- Goldenberg, A., Zheng, A., Fienberg, S., & Airolidi, E. (2010). A survey of statistical network models. *Foundations and Trends in Machine Learning*, *2*(2), 129–233.
- Gotelli, N. J. (2000). Null model analysis of species co-occurrence patterns. *Ecology*, *81*(9), 2606–2621.
- Greenhill, C., & McKay, B. D. (2009). Random dense bipartite graphs and directed graphs with specified degrees. *Random Struct. Algor.*, *35*(2), 222–249.
- Greenhill, C., McKay, B. D., & Wang, X. (2006). Asymptotic enumeration of sparse 0–1 matrices with irregular row and column sums. *J. Comb. Theory A*, *113*(2), 291–324.
- Harrison, M. T. (2012). Conservative hypothesis tests and confidence intervals using importance sampling. *Biometrika*, *99*(1), 57–69.
- Harrison, M. T., & Geman, S. (2009). A rate and history-preserving resampling algorithm for neural spike trains. *Neural comput.*, *21*(5), 1244–1258.
- Holland, P. W., & Leinhardt, S. (1981). An exponential family of probability distributions for directed graphs. *J. Am. Statist. Assoc.*, *76*, 33–50.
- Jerrum, M., Sinclair, A., & Vigoda, E. (2004). A polynomial-time approximation algorithm for the permanent of a matrix with nonnegative entries. *J. Assoc. Comp. Mach.*, *51*(4), 671–697.
- Kannan, R., Tetali, P., & Vempala, S. (1999). Simple Markov-chain algorithms for generating bipartite graphs and tournaments. *Random Struct. Algor.*, *14*(4), 293–308.
- Kong, A., Liu, J. S., & Wong, W. H. (1994). Sequential imputations and Bayesian missing data problems. *J. Am. Statist. Assoc.*, *89*, 278–288.
- Kou, S. C., & McCullagh, P. (2009). Approximating the  $\alpha$ -permanent. *Biometrika*, *96*(3), 635–644.
- Littlewood, D. (1950). *The theory of group characters and matrix representations of groups*. Oxford Univ. Press.
- Liu, J. S. (2001). *Monte Carlo strategies in scientific computing*. New York: Springer.

- Macchi, O. (1975). The coincidence approach to stochastic point processes. *Adv. Appl. Probab.*, 7, 83–122.
- Matsumoto, M., & Nishimura, T. (1998). Mersenne twister: a 623-dimensionally equidistributed uniform pseudo-random number generator. *ACM Trans. Model. Comp. Simul.*, 8(1), 3–30.
- McCullagh, P., & Møller, J. (2006). The permanental process. *Adv. Appl. Probab.*, 38, 873–888.
- McKay, B. D. (1984). Asymptotics for 0-1 matrices with prescribed line sums. In D. M. Jackson & S. A. Vanstone (Eds.), *Enumeration and design* (pp. 225–238). Academic Press.
- McKay, B. D., & Wormald, N. C. (1990). Uniform generation of random regular graphs of moderate degree. *Journal of Algorithms*, 11(1), 52–67.
- Mehta, C. R., & Patel, N. R. (1995). Exact logistic regression: theory and examples. *Statist. Med.*, 14(19), 2143–2160.
- O’Neil, P. E. (1969). Asymptotics and random matrices with row-sum and column-sum restrictions. *B. Am. Math. Soc.*, 75, 1276–1282.
- Park, S. K., & Miller, K. W. (1988). Random number generators: good ones are hard to find. *Commun. Assoc. Comp. Mach.*, 31(10), 1192–1201.
- Ponocny, I. (2001). Nonparametric goodness-of-fit tests for the Rasch model. *Psychometrika*, 66(3), 437–459.
- Rao, A., Jana, R., & Bandyopadhyay, S. (1996). A Markov chain Monte Carlo method for generating random (0, 1)-matrices with given marginals. *Sankhyā: The Indian Journal of Statistics, Series A*, 225–242.
- Rasch, G. (1960). *Probabilistic models for some intelligence and attainment tests*. Copenhagen: Danmarks Paedagogiske Institut.
- Rasch, G. (1961). On general laws and the meaning of measurement in psychology. In J. Neyman (Ed.), *Proceedings of the fourth berkeley symposium on mathematical statistics and probability: Probability theory* (Vol. 4, pp. 321–334). Berkeley, CA.
- Rothblum, U. G., & Schneider, H. (1989). Scalings of matrices which have prespecified row sums and column sums via optimization. *Linear Algebra Appl.*, 114, 737–764.
- Ryser, H. J. (1957). Combinatorial properties of matrices of zeros and ones. *Can. J. Math.*, 9, 371–377.
- Shirai, T., & Takahashi, Y. (2003). Random point fields associated with certain Fredholm determinants I: fermion, Poisson and boson point processes. *J. Funct. Anal.*, 205(2), 414–463.
- Sinkhorn, R. (1964). A relationship between arbitrary positive matrices and doubly stochastic matrices. *Ann. Math. Statist.*, 35(2), 876–879.

- Sinkhorn, R. (1967). Diagonal equivalence to matrices with prescribed row and column sums. *Am. Math. Mon.*, *74*(4), 402–405.
- Snijders, T. A. B. (1991). Enumeration and simulation methods for 0–1 matrices with given marginals. *Psychometrika*, *56*(3), 397–417.
- StataCorp. (2009). *Stata statistical software: Release 11*. College Station, TX: StataCorp LP.
- Valiant, L. G. (1979). The complexity of computing the permanent. *Theor. Comput. Sci.*, *8*(2), 189–201.
- Vaughan, R. J., & Venables, W. N. (1972). Permanent expressions for order statistic densities. *J. R. Statist. Soc. B*, *34*, 308–310.
- Vere-Jones, D. (1988). A generalization of permanents and determinants. *Linear Algebra Appl.*, *111*, 119–124.
- Vere-Jones, D. (1997). Alpha-permanents and their applications to multivariate gamma, negative binomial and ordinary binomial distributions. *New Zeal. J. Math.*, *26*, 125–149.
- Verhelst, N. D. (2008). An efficient MCMC algorithm to sample binary matrices with fixed marginals. *Psychometrika*, *73*(4), 705–728.
- Wasserman, S. S. (1977). Random directed graph distributions and the triad census in social networks. *J. Math. Sociol.*, *5*(1), 61–86.
- Zamar, D., McNeney, B., & Graham, J. (2007). elm: Software implementing exact-like inference for logistic regression models. *J. Statist. Software*, *21*, 1–18.