
A Family of Exact Goodness-of-Fit Tests for High-Dimensional Discrete Distributions

Feras A. Saad
MIT

Cameron E. Freer
MIT

Nathanael L. Ackerman Vikash K. Mansinghka
Harvard University MIT

Abstract

The objective of goodness-of-fit testing is to assess whether a dataset of observations is likely to have been drawn from a candidate probability distribution. This paper presents a rank-based family of goodness-of-fit tests that is specialized to discrete distributions on high-dimensional domains. The test is readily implemented using a simulation-based, linear-time procedure. The testing procedure can be customized by the practitioner using knowledge of the underlying data domain. Unlike most existing test statistics, the proposed test statistic is distribution-free and its exact (non-asymptotic) sampling distribution is known in closed form. We establish consistency of the test against all alternatives by showing that the test statistic is distributed as a discrete uniform if and only if the samples were drawn from the candidate distribution. We illustrate its efficacy for assessing the sample quality of approximate sampling algorithms over combinatorially large spaces with intractable probabilities, including random partitions in Dirichlet process mixture models and random lattices in Ising models.

1 Introduction

We address the problem of testing whether a dataset of observed samples was drawn from a candidate probability distribution. This problem, known as goodness-of-fit testing, is of fundamental interest and has applications in a variety of fields including Bayesian statistics [10; 31], high-energy physics [34], astronomy [22], genetic association studies [17], and psychometrics [3].

Rank-based methods are a popular approach for assessing goodness-of-fit and have received great attention in the nonparametric statistics literature [15]. However, the majority of existing rank-based tests operate under the assumption of continuous distributions [16, VI.8] and analogous methods for discrete distributions that are theoretically rigorous, customizable using domain knowledge, and practical to implement in a variety of settings remain much less explored.

This paper presents a new connection between rank-based tests and discrete distributions on high-dimensional data structures. By algorithmically specifying an ordering on the data domain, the practitioner can quantitatively assess how typical the observed samples are with respect to resampled data from the candidate distribution. This ordering is leveraged by the test to effectively surface distributional differences.

More specifically, we propose to test whether observations $\{y_1, \dots, y_n\}$, taking values in a countable set \mathcal{T} , were drawn from a given discrete distribution \mathbf{p} on the basis of the rank of each y_i with respect to m i.i.d. samples $\{x_1, \dots, x_m\}$ from \mathbf{p} . If y_i was drawn from \mathbf{p} then we expect its rank to be uniformly distributed over $\{0, 1, \dots, m\}$. When the ranks show a deviation from uniformity, it is unlikely that the y_i were drawn from \mathbf{p} . A key step is to use continuous random variables to break any ties when computing the ranks. We call this statistic the Stochastic Rank Statistic (SRS), which has several desirable properties for goodness-of-fit testing:

1. The SRS is distribution-free: its sampling distribution under the null does not depend on \mathbf{p} . There is no need to construct ad-hoc tables or use Monte Carlo simulation to estimate rejection regions.
2. The exact (non-asymptotic) sampling distribution of the SRS is a discrete uniform. This exactness obviates the need to apply asymptotic approximations in small-sample and sparse regimes.
3. The test is consistent against all alternatives. We show that the SRS is distributed as a discrete uniform if and only if $\{y_1, \dots, y_n\} \sim^{\text{iid}} \mathbf{p}$.

4. The test gives the practitioner flexibility in deciding the set of properties on which the observations be checked to agree with samples from \mathbf{p} . This flexibility arises from the design of the ordering on the domain that is used to compute the ranks.
5. The test is readily implemented using a procedure that is linear-time in the number of observations. The test is simulation-based and does not require explicitly computing $\mathbf{p}(x)$, which is especially useful for distributions with intractable probabilities.

While the test is consistent for any ordering (\mathcal{T}, \prec) over the domain that is used to compute the SRS, the power of the test depends heavily on the choice of \prec . We show how to construct orderings in a variety of domains by (i) defining procedures that traverse and compare discrete data structures; (ii) composing probe statistics that summarize key numerical characteristics; and (iii) using randomization to generate arbitrary orderings.

The remainder of the paper is organized as follows. Section 2 reviews the goodness-of-fit problem and discusses related work. Section 3 presents the proposed test and several theoretical properties. Section 4 gives conceptual examples for distributions over integers, binary strings, and partitions. Section 5 applies the method to (i) compare approximate Bayesian inference algorithms over mixture assignments in a Dirichlet process mixture model and (ii) assess the sample quality of random lattices from approximate samplers for the Ising model.

2 The Goodness-of-Fit Problem

Problem 2.1. Let \mathbf{p} be a candidate discrete distribution over a finite or countably infinite domain \mathcal{T} . Given observations $\{y_1, \dots, y_n\}$ drawn i.i.d. from an unknown distribution \mathbf{q} over \mathcal{T} , is there sufficient evidence to reject the hypothesis $\mathbf{p} = \mathbf{q}$?

In the parlance of statistical testing, we have the following null and alternative hypotheses:

$$\mathbf{H}_0 := [\mathbf{p} = \mathbf{q}] \quad \mathbf{H}_1 := [\mathbf{p} \neq \mathbf{q}].$$

A statistical test $\phi_n: \mathcal{T}^n \rightarrow \{\text{reject}, \text{not reject}\}$ says, for each size n dataset, whether to reject or not reject the null hypothesis \mathbf{H}_0 . We define the significance level

$$\alpha := \Pr \{ \phi_n(Y_{1:n}) = \text{reject} \mid \mathbf{H}_0 \} \quad (1)$$

to be the probability of incorrectly declaring **reject**. For a given level α , the performance of the test ϕ_n is characterized by its power

$$\beta := \Pr \{ \phi_n(Y_{1:n}) = \text{reject} \mid \mathbf{H}_1 \}, \quad (2)$$

which is the probability of correctly declaring **reject**.

Classical goodness-of-fit tests for nominal (unordered) data include the multinomial test [14]; Pearson chi-square test [23]; likelihood-ratio test [33]; nominal Kolmogorov–Smirnov test [13; 24]; and power-divergence statistics [26]. For ordinal data, goodness-of-fit test statistics include the ordinal Watson, Cramér–von Mises, and Anderson–Darling [7] tests as well as the ordinal Kolmogorov–Smirnov [4; 8]. These approaches typically suffer from statistical issues in large domains. They assume that $\mathbf{p}(x)$ is easy to evaluate (which is rarely possible in modern machine-learning applications such as graphical models) and/or require that each discrete outcome $x \in \mathcal{T}$ has a non-negligible expectation $n\mathbf{p}(x)$ [20; 28] (which requires a large number of observations n even when \mathbf{p} and \mathbf{q} are noticeably far from one another). In addition, the rejection regions of these statistics are either distribution-dependent (which requires reestimating the region for each new candidate distribution \mathbf{p}) or asymptotically distribution-free (which is inexact for finite-sample data and imposes additional statistical assumptions on \mathbf{p} and \mathbf{q}). The Mann–Whitney U [19], which is also a rank-based test that bears some similarity to the SRS, is only consistent under median shift, whereas the proposed method is consistent under general distributional inequality.

Recent work in the theoretical computer science literature has established computational and sample complexity bounds for testing approximate equality of discrete distributions [5]. These methods have been primarily studied from a theoretical perspective and have not been shown to yield practical goodness-of-fit tests in practice, nor have they attained widespread adoption in the applied statistics community. For instance, the test in [1] is based on a variant of Pearson chi-square. It requires enumerating over the domain \mathcal{T} and representing $\mathbf{p}(x)$ explicitly. The test in [32] requires specifying and solving a complex linear program. While these algorithms may obtain asymptotically sample-optimal limits, they are designed to detect differences between \mathbf{p} and \mathbf{q} in a way that is robust to highly adversarial settings. These tests do not account for any structure in the domain \mathcal{T} that can be leveraged by the practitioner to effectively surface distributional differences.

Permutation and bootstrap resampling of test statistics are another family of tests for goodness-of-fit [11]. Theoretically rigorous and consistent tests can be obtained using kernel methods, including the maximum mean discrepancy [12] and discrete Stein discrepancy [35]. Since the null distribution is unknown, rejection regions are estimated by bootstrap resampling, which may be inexact due to discreteness of the data. Instead of bootstrapping, the SRS can be used to obtain an exact, distribution-free test by defining an ordering using the kernel. This connection is left for future work.

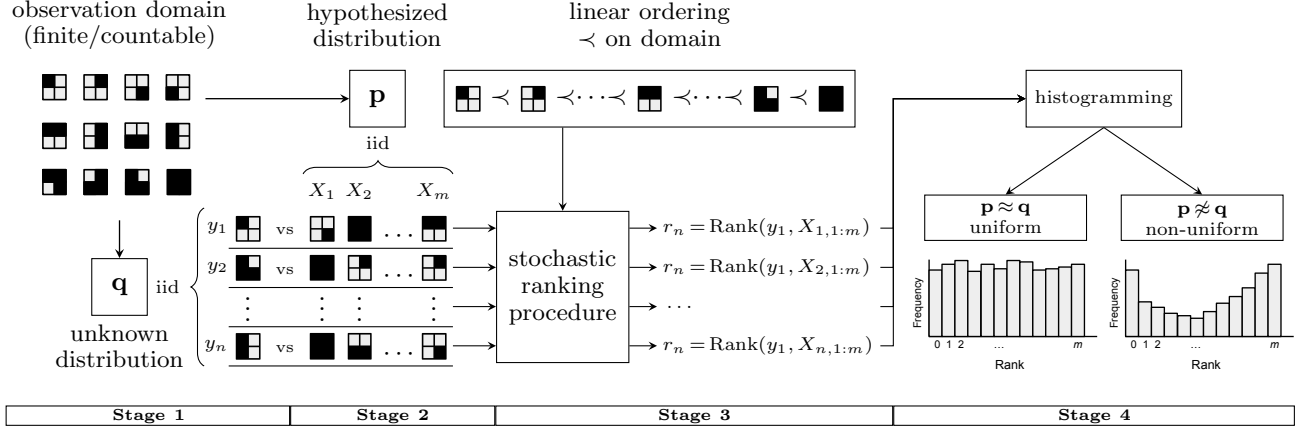


Figure 1: Overview of the proposed goodness-of-fit test for discrete distributions. **Stage 1:** Observations $\{y_1, \dots, y_n\}$ are assumed to be drawn i.i.d. from an unknown discrete distribution \mathbf{q} over a finite or countable observation domain \mathcal{T} (shown in the top-left corner). **Stage 2:** For each y_i , m samples $\{X_{i1}, \dots, X_{im}\}$ are simulated i.i.d. from the candidate distribution \mathbf{p} over \mathcal{T} . **Stage 3:** Given a total order \prec on \mathcal{T} and the observed and simulated data, a stochastic ranking procedure returns the rank r_i of each y_i within $\{X_{i1}, \dots, X_{im}\}$, using uniform random numbers to ensure the ranks are unique. **Stage 4:** The histogram of the ranks $\{r_1, \dots, r_n\}$ is analyzed for uniformity over $\{0, 1, \dots, m\}$.

3 A Family of Exact and Distribution-Free GOF Tests

In this section we describe our proposed method for addressing the goodness-of-fit problem. The proposed procedure combines (i) the intuition from existing methods for ordinal data [7] that the deviation between the expected CDF and empirical CDF of the sample serves as a good signal for goodness-of-fit, with (ii) the flexibility of probe statistics in Monte Carlo-based resampling tests [11] to define, using an ordering \prec on \mathcal{T} , characteristics of the distribution that are of interest to the experimenter. Figure 1 shows the step-by-step workflow of the proposed test and Algorithm 1 formally describes the testing procedure.

Algorithm 1 Exact GOF Test using SRS

Input: $\left\{ \begin{array}{l} \text{simulator for candidate dist. } \mathbf{p} \text{ over } \mathcal{T}; \\ \text{i.i.d. samples } \{y_1, y_2, \dots, y_n\} \text{ from dist. } \mathbf{q}; \\ \text{strict total order } \prec \text{ on } \mathcal{T}, \text{ of any order type;} \\ \text{number } m \geq 1 \text{ of datasets to resimulate;} \\ \text{significance level } \alpha \text{ of hypothesis test;} \end{array} \right.$

Output: Decision to reject the null hypothesis $H_0: \mathbf{p} = \mathbf{q}$ versus alternative hypothesis $H_1: \mathbf{p} \neq \mathbf{q}$ at level α .

- 1: **for** $i = 1, 2, \dots, n$ **do**
- 2: $X_1, X_2, \dots, X_m \sim^{\text{iid}} \mathbf{p}$
- 3: $U_0, U_1, \dots, U_m \sim^{\text{iid}} \text{Uniform}(0, 1)$
- 4: $r_i \leftarrow \sum_{k=1}^m \mathbb{I}[X_k \prec y_i] + \mathbb{I}[X_k = y_i, U_k < U_0]$
- 5: Use a standard hypothesis test to compute p -value of $\{r_1, \dots, r_n\}$ under a discrete uniform on $\{0, \dots, m\}$.
- 6: **return** reject if $p \leq \alpha$, else not reject.

The proposed method addresses shortcomings of existing statistics in sparse regimes. It does not require the ability to compute $\mathbf{p}(x)$ and it is not based on comparing the expected frequency of each $x \in \mathcal{T}$ (which is

often vanishingly small) with its observed frequency. Furthermore, the stochastic rank statistics r_i have an exact and distribution-free sampling distribution. The following theorem establishes that the r_i are uniformly distributed if and only if $\mathbf{p} = \mathbf{q}$. (Proofs are in the Appendix.)

Theorem 3.1. *Let \mathcal{T} be a finite or countably infinite set, let \prec be a strict total order on \mathcal{T} , let \mathbf{p} and \mathbf{q} be two probability distributions on \mathcal{T} , and let m be a positive integer. Consider the following random variables:*

$$X_0 \sim \mathbf{q} \quad (3)$$

$$X_1, X_2, \dots, X_m \sim^{\text{iid}} \mathbf{p} \quad (4)$$

$$U_0, U_1, U_2, \dots, U_m \sim^{\text{iid}} \text{Uniform}(0, 1) \quad (5)$$

$$R = \sum_{j=1}^m \mathbb{I}[X_j \prec X_0] + \mathbb{I}[X_j = X_0, U_j < U_0]. \quad (6)$$

Then $\mathbf{p} = \mathbf{q}$ if and only if for all $m \geq 1$, the rank R is distributed as a discrete uniform random variable on the set of integers $[m+1] := \{0, 1, \dots, m\}$.

Note that the r_i in line 4 of Algorithm 1 are n i.i.d. samples of the random variable R in Eq. (6), which is the rank of $X_0 \sim \mathbf{q}$ within a size m sample $X_{1:m} \sim^{\text{iid}} \mathbf{p}$. For Theorem 3.1, it is essential that ties are broken by pairing each X_i with a uniform random variable U_i , as opposed to, e.g., breaking each tie independently with probability $1/2$, as demonstrated by the next example.

Example 3.2. Let \mathcal{T} contain a single element. Then all the X_i (for $0 \leq i \leq m$) are equal almost surely. Break each tie between X_0 and X_j by flipping a fair coin. Then R is binomially distributed with m trials and weight $1/2$, not uniformly distributed over $[m+1]$.

We now establish theoretical properties of R which form the basis of the goodness-of-fit test in Algorithm 1.

First note that in the case where all the X_i are almost surely distinct, the forward direction of Theorem 3.1, which establishes that if $\mathbf{p} = \mathbf{q}$ then the rank R is uniform for all $m \geq 1$, is easy to show and is known in the statistical literature [2]. However no existing results make the connection between rank statistics and discrete random variables over countable domains with ties broken stochastically. Nor do they establish that $\mathbf{p} = \mathbf{q}$ is a *necessary* condition for uniformity of R (across all m beyond some integer) and can therefore be used as the basis of a consistent goodness-of-fit test. We now state an immediate consequence of Theorem 3.1.

Corollary 3.3. *If $\mathbf{p} \neq \mathbf{q}$, then there is some $M \geq 1$ such that R is not uniformly distributed on $[M + 1]$.*

The next theorem significantly strengthens Corollary 3.3 by showing that if $\mathbf{p} \neq \mathbf{q}$, the rank statistic is non-uniform for *all but finitely many* m .

Theorem 3.4. *Let $\mathbf{p} \neq \mathbf{q}$ and M be defined as in Corollary 3.3. Then for all $m \geq M$, the rank R is not uniformly distributed on $[m + 1]$.*

In fact, unless \mathbf{p} and \mathbf{q} satisfy an adversarial symmetry relationship under the selected ordering \prec , the rank is non-uniform for *all* $m \geq 1$.

Corollary 3.5. *Let \triangleleft denote the lexicographic order on $\mathcal{T} \times [0, 1]$ induced by (\mathcal{T}, \prec) and $([0, 1], <)$. Suppose $\Pr\{(X, U_1) \triangleleft (Y, U_0)\} \neq 1/2$ for $Y \sim \mathbf{q}$, $X \sim \mathbf{p}$, and $U_0, U_1 \sim \text{iid Uniform}(0, 1)$. Then for all $m \geq 1$, the rank R is not uniformly distributed on $[m + 1]$.*

The next theorem establishes the existence of an ordering on \mathcal{T} satisfying the hypothesis of Corollary 3.5.

Theorem 3.6. *If $\mathbf{p} \neq \mathbf{q}$, then there is an ordering \prec^* whose associated rank statistic R is non-uniform for $m = 1$ (and hence by Theorem 3.4 for all $m \geq 1$).*

Intuitively, \prec^* sets elements $x \in \mathcal{T}$ which have a high probability under \mathbf{q} to be “small” in the linear order, and elements $x \in \mathcal{T}$ which have a high probability under \mathbf{p} to be “large” in the linear order. More precisely, \prec^* maximizes the sup-norm distance between the induced cumulative distribution functions $\tilde{\mathbf{p}}$ and $\tilde{\mathbf{q}}$ of \mathbf{p} and \mathbf{q} , respectively (Figure 3). Under a slight variant of this ordering, for finite \mathcal{T} , the next theorem establishes the sample complexity required to obtain exponentially high power in terms of the statistical distance $L_\infty(\mathbf{p}, \mathbf{q}) = \sup_{x \in \mathcal{T}} |\mathbf{p}(x) - \mathbf{q}(x)|$ between \mathbf{p} and \mathbf{q} .

Theorem 3.7. *Given significance level $\alpha = 2\Phi(-c)$ for $c > 0$, there is an ordering for which the proposed test with $m = 1$ achieves power $\beta \geq 1 - \Phi(-c)$ using*

$$n \approx 4c^2 / L_\infty(\mathbf{p}, \mathbf{q})^4 \quad (7)$$

samples from \mathbf{q} , where Φ is the cumulative distribution function of a standard normal.

This key result is independent of the domain size and establishes a lower bound for any \prec because it is based on the optimal ordering \prec^* . The next theorem derives the exact sampling distribution for any pair of distributions (\mathbf{p}, \mathbf{q}) , which is useful for simulation studies (e.g., Figure 3) that characterize the power of the SRS.

Theorem 3.8. *The distribution of R is given by*

$$\Pr\{R = r\} = \sum_{x \in \mathcal{T}} H(x, m, r) \mathbf{q}(x) \quad (8)$$

for $0 \leq r \leq m$, where $H(x, m, r) :=$

$$\left\{ \begin{array}{ll} \sum_{e=0}^m \left\{ \left[\sum_{j=0}^e \binom{m-e}{r-j} \left[\frac{\tilde{\mathbf{p}}(x)}{1 - \mathbf{p}(x)} \right]^{r-j} \right. \right. \\ \left. \left. \left[1 - \frac{\tilde{\mathbf{p}}(x)}{1 - \mathbf{p}(x)} \right]^{(m-e)-(r-j)} \left(\frac{1}{e+1} \right) \right] \right. \\ \left. \binom{m}{e} [\mathbf{p}(x)]^m [1 - \mathbf{p}(x)]^{e-m} \right\} & \text{if } 0 < \mathbf{p}(x) < 1, \\ \binom{r}{m} [\tilde{\mathbf{p}}(x)]^r [1 - \tilde{\mathbf{p}}(x)]^{m-r} & \text{if } \mathbf{p}(x) = 0, \\ \frac{1}{m+1} & \text{if } \mathbf{p}(x) = 1, \end{array} \right.$$

and $\tilde{\mathbf{p}}(x) := \sum_{x' \prec x} \mathbf{p}(x')$ is the CDF of \mathbf{p} .

4 Examples

We now apply the proposed test to a countable domain and two high-dimensional finite domains, illustrating a power comparison and how distributional differences can be detected when the number of observations is much smaller than the domain size. We use Pearson chi-square to assess uniformity of the SRS for Algorithm 1 (see [29] for alternative ways to test for a uniform null).

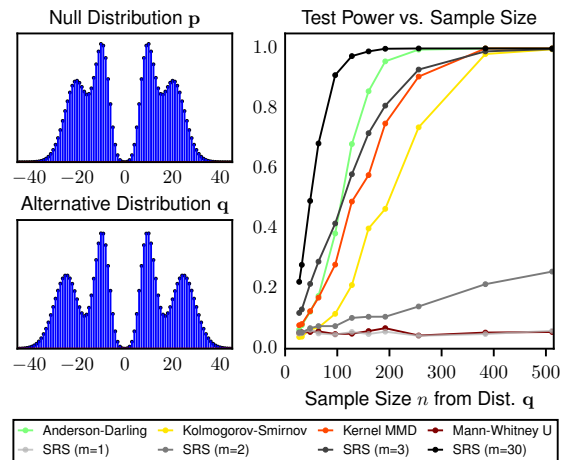


Figure 2: The left panel shows a pair (\mathbf{p}, \mathbf{q}) of reflected, bimodal Poisson distributions with slight location shift. The right plot compares the power of testing $\mathbf{p} = \mathbf{q}$ using the SRS (for various choices of m) to several baseline methods.

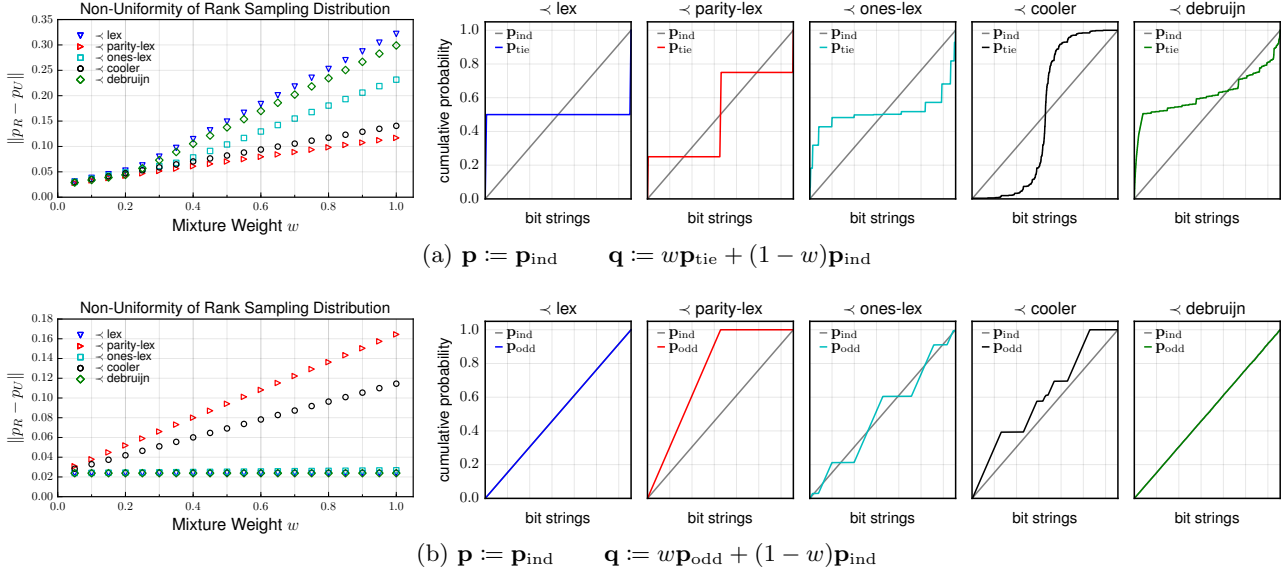


Figure 3: In each of the two panels (a) and (b), the left plot shows the sup-norm distance between the sampling distribution of the rank statistic and the discrete uniform (using Eq. 8 in Theorem 3.8), for a uniform null $\mathbf{p} := \mathbf{p}_{\text{ind}}$ on $\{0, 1\}^{16}$ against alternative distributions of the form $\mathbf{q} := w\mathbf{p}_{\text{alt}} + (1-w)\mathbf{p}_{\text{ind}}$, for increasing mixture weight $0 \leq w \leq 1$ and six different orderings on the binary strings. The right plot compares the cumulative distribution function of the null distribution (diagonal line in gray) with the cumulative distribution functions of the alternative distribution (when $w = 1$) as obtained by sorting the binary strings according to each ordering. Orderings which induce a greater distance between the cumulative distribution functions of the null and alternative distributions result in more power to detect the alternative.

4.1 Bimodal, Symmetric Poisson

We first investigate the performance of the SRS for testing a pair of symmetric, multi-modal distributions over the integers with location shift. In particular, for $x \in \mathbb{Z}$, define distribution $\mathbf{f}(x; \lambda_1, \lambda_2) := \frac{1}{2} (\frac{1}{2} \text{Poisson}(|x|; \lambda_1) + \frac{1}{2} \text{Poisson}(|x|; \lambda_2))$. Note \mathbf{f} is a mixture of Poisson distributions with rates λ_1 and λ_2 , reflected symmetrically about $x = 0$. We set $\mathbf{p}(x) := \mathbf{f}(x; 10, 20)$ and $\mathbf{q} := \mathbf{f}(x; 10, 25)$ so that \mathbf{q} is location-shifted in two of the four modes (Figure 2, left panel).

The right plot of Figure 2 compares the power for various sample sizes n from \mathbf{q} according to the SRS ($m = 1, 2, 3, 30$, shown in increasing shades of gray) and several baselines (shown in color). The baselines (AD, MMD, KS, and Mann–Whitney U) are used to assess goodness-of-fit by performing a two-sample test on n samples from \mathbf{q} with samples drawn i.i.d. from \mathbf{p} . The power (at level $\alpha = 0.05$) is estimated as the fraction of correct answers over 1024 independent trials. The Mann–Whitney U, which is also based on rank statistics with a correction for ties, has no power for all n as it can only detect median shift, as does the SRS with $m = 1$ (see Corollary 3.5). The SRS becomes non-uniform for $m = 2$ although this choice results in low power. The SRS with $m = 3$ has comparable power to the AD and MMD tests. The SRS with $m = 30$ is the most powerful, although it requires more computational effort and samples from \mathbf{p} (Algorithm 1 scales as $O(mn)$).

4.2 Binary strings

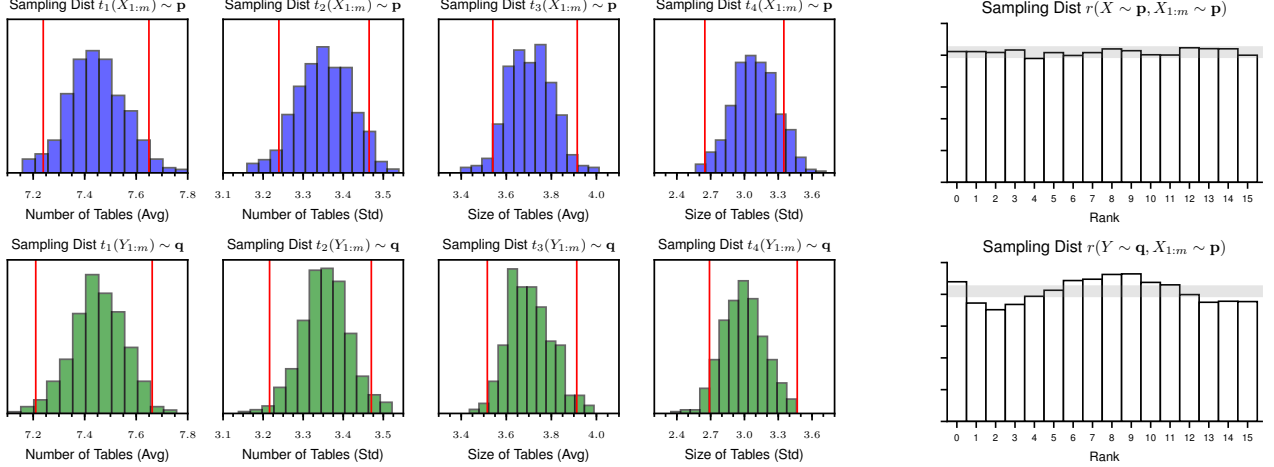
Let $\mathcal{T} := \{0, 1\}^k$ be the set of all length k binary strings. Define the following distributions to be uniform over all strings $x = (x_1, \dots, x_k) \in \{0, 1\}^k$ which satisfy the given predicates:

$$\begin{aligned} \mathbf{p}_{\text{ind}} &: \text{uniform on all strings,} \\ \mathbf{p}_{\text{odd}} &: \sum_{i=1}^k x_i \equiv 1 \pmod{2}, \\ \mathbf{p}_{\text{tie}} &: x_1 = x_2 = \dots = x_{k/2}. \end{aligned}$$

Each of these distributions assigns marginal probability $1/2$ to each bit x_i (for $1 \leq i \leq k$), so all deviations from the uniform distribution \mathbf{p}_{ind} are captured by higher-order relationships. The five orderings used for comparing binary strings are

- \prec_{lex} : Lexicographic (dictionary) ordering,
- \prec_{par} : Parity of ones, ties broken using \prec_{lex} ,
- \prec_{one} : Number of ones, ties broken using \prec_{lex} ,
- \prec_{coo} : Cooler ordering (randomly generated) [30],
- \prec_{dbj} : De Bruijn sequence ordering.

We set the null distribution $\mathbf{p} := \mathbf{p}_{\text{ind}}$ and construct alternative distributions $\mathbf{q} := w\mathbf{p}_c + (1-w)\mathbf{p}_{\text{ind}}$ as mixtures of \mathbf{p}_{ind} with the other two distributions, where $w \in [0, 1]$ and $c \in \{\text{odd}, \text{not}\}$. We take bit strings of length $k = 16$ with $n = 256$ observations so that $|\mathcal{T}| = 65,536$ and 0.4% of the domain size is observed.



(a) Sampling distribution of four different probe statistics $\{t_1, t_2, t_3, t_4\}$ of a dataset of partitions, as sampled from \mathbf{p} (Eq. (9); blue) and from \mathbf{q} (Eq. (10); green) estimated by Monte Carlo simulation. Vertical red lines indicate 2.5% and 97.5% quantiles. Even though $\mathbf{p} \neq \mathbf{q}$, the distributions of these statistics are aligned in such a way that a statistic $t_j(Y_{1:m}) \sim \mathbf{q}$ is unlikely to appear as an extreme value in the sampling distribution of the corresponding statistic $t_j(X_{1:m}) \sim \mathbf{p}$, which leads to under-powered resampling-based tests.

(b) Monte Carlo simulation of the rank statistic illustrates its significant uniform distribution under the null hypothesis (top) and significant non-uniform distribution under the alternative hypothesis (bottom).

Figure 4: Comparison of the sampling distribution of (a) various bootstrapped probe statistics [11] with (b) the stochastic rank statistic, for goodness-of-fit testing the Chinese restaurant processes on $N = 20$ customers. Discussion in main text.

Figure 3 shows how the non-uniformity of the SRS (computed using Theorem 3.8) varies for each of the two alternatives and five orderings ($m = 6$). Each ordering induces a different CDF over $\{0, 1\}^k$ for the alternative distribution, shown in the right panel for $w = 1$. Orderings with a greater maximum vertical distance between the null and alternative CDF attain greater rank non-uniformity. No single ordering is more powerful than all others in both test cases. However, in each case, some ordering detects the difference even at low weights w , despite the sparse observation set.

The alternative $\mathbf{q} = \mathbf{p}_{\text{odd}}$ in Figure 3b is especially challenging: in a sample, all substrings (not necessarily contiguous) of a given length $j < k$ are equally likely. Even though the SRS is non-uniform for all orderings, the powers vary significantly. For example, comparing strings using \prec_{lex} does not effectively distinguish between \mathbf{p}_{ind} and \mathbf{p}_{odd} , as strings with an odd number of ones are lexicographically evenly interspersed within the set of all strings. The parity ordering (which is optimal for this alternative) and the randomly generated cooler ordering have increasing power as w increases.

4.3 Partition testing

We next apply the SRS to test distributions on the space of partitions of the set $\{1, 2, \dots, N\}$. Let Π_N denote the set of all such partitions. We define a distribution on Π_N using the two-parameter Chinese Restaurant Process (CRP) [6, Section 5.1]. Letting $(x|y)_N := (x)(x+y) \cdots (x+(N-1)y)$, the probability

of a partition $\pi := \{\pi_1, \dots, \pi_k\} \in \Pi_N$ with k tables (blocks) is given by

$$\text{CRP}(\pi; a, b) := \begin{cases} \frac{(b|a)_k}{(b|1)_N} \prod_{i=1}^k (1-a)_{c_k-1} & (\text{if } a > 0) \\ \frac{b^k}{(b|1)_N} \prod_{i=1}^k (c_k - 1)! & (\text{if } a = 0), \end{cases}$$

where c_i is the number of customers (integers) at table π_i ($1 \leq i \leq k$). Simulating a CRP proceeds by sequentially assigning customers to tables [6, Def. 7]. Even though we can compute the probability of any partition, the cardinality of Π_N grows exponentially in N (e.g., $|\Pi_{20}| \approx 5.17 \times 10^{13}$). The expected frequency of any partition is essentially zero for sample size $n \ll |\Pi_N|$, so Pearson chi-square or likelihood-ratio tests on the raw data are inappropriate. Algorithm 2 defines a total order on the partition domain Π_N .

Algorithm 2 Total order \prec on the set of partitions Π_N

Input: $\left\{ \begin{array}{l} \text{Partition } \pi := \{\pi_1, \pi_2, \dots, \pi_k\} \in \Pi_N \text{ with } k \text{ blocks.} \\ \text{Partition } \nu := \{\nu_1, \nu_2, \dots, \nu_l\} \in \Pi_N \text{ with } l \text{ blocks.} \end{array} \right.$

Output: LT if $\pi \prec \nu$; GT if $\pi \succ \nu$; EQ if $\pi = \nu$.

- 1: if $k < l$ then return LT $\triangleright \nu$ has more blocks
- 2: if $k > l$ then return GT $\triangleright \pi$ has more blocks
- 3: $\tilde{\pi} \leftarrow$ blocks of π sorted by value of least element in the block
- 4: $\tilde{\nu} \leftarrow$ blocks of ν sorted by value of least element in the block
- 5: for $b = 1, 2, \dots, l$ do
- 6: if $|\tilde{\pi}_b| < |\tilde{\nu}_b|$ then return LT $\triangleright \tilde{\nu}_b$ has more elements
- 7: if $|\tilde{\pi}_b| > |\tilde{\nu}_b|$ then return GT $\triangleright \tilde{\pi}_b$ has more elements
- 8: $\pi'_b \leftarrow$ values in $\tilde{\pi}_b$ sorted in ascending order
- 9: $\nu'_b \leftarrow$ values in $\tilde{\nu}_b$ sorted in ascending order
- 10: for $i = 1, 2, \dots, |\pi'_b|$ do
- 11: if $\pi'_{b,i} < \nu'_{b,i}$ then return LT $\triangleright \pi'_b$ has smallest element
- 12: if $\pi'_{b,i} > \nu'_{b,i}$ then return GT $\triangleright \nu'_b$ has smallest element
- 13: return EQ

We consider the following pair of distributions:

$$\mathbf{p} := \text{CRP}(0.26, 0.76)/2 + \text{CRP}(0.19, 5.1)/2 \quad (9)$$

$$\mathbf{q} := \text{CRP}(0.52, 0.52). \quad (10)$$

These distributions are designed to ensure that partitions from \mathbf{p} and \mathbf{q} have similar distributions on the number and sizes of tables. Figure 4a shows a comparison of using Monte Carlo simulation of various bootstrapped probe statistics for assessing goodness-of-fit versus using the SRS with the ordering in Algorithm 2.

In Figure 4a, each probe statistic takes a size m dataset $X_{1:m}$ (where each X_i is a partition) and produces a numerical summary such as the average of the number of tables in each sample. A resampling test [11] that uses these probe statistics will report (with high probability) that an observed statistic $t(Y_{1:m}) \sim \mathbf{q}$ drawn from the alternative distribution is a non-extreme value in the null distribution $t(X_{1:m}) \sim \mathbf{p}$ (as indicated by alignment of their quantiles, shown in red) and will therefore have insufficient evidence to reject $\mathbf{p} = \mathbf{q}$.

On the other hand, Figure 4b shows that when ranked using the ordering obtained from Algorithm 2 (which is based on a multivariate combination of the univariate probe statistics in Figure 4a specified procedurally), a partition $Y \sim \mathbf{q}$ is more likely to lie in the center of a dataset $X_{1:m} \sim^{\text{iid}} \mathbf{p}$, as illustrated by the non-uniform rank distribution under the alternative hypothesis (the gray band shows 99% variation for a uniform histogram). By comparing the top and bottom panels of Figure 4b, the SRS shows that partitions from \mathbf{q} have a poor fit with respect to partitions from \mathbf{p} , despite their agreement on multiple univariate summary statistics shown in Figure 4a.

5 Applications

We next apply the proposed test to assess the sample quality of random data structures obtained from approximate sampling algorithms over combinatorially large domains with intractable probabilities.

5.1 Dirichlet process mixture models

The recent paper [31] describes simulation-based calibration (SBC), a procedure for validating samples from algorithms that can generate posterior samples for a hierarchical Bayesian model. More specifically, for a prior $\pi(z)$ over the parameters z and likelihood function $\pi(x|z)$ over data x , integrating the posterior over the joint distribution returns the prior distribution:

$$\pi(z) = \int [\pi(z|x')\pi(x'|z')] \pi(z') dz'. \quad (11)$$

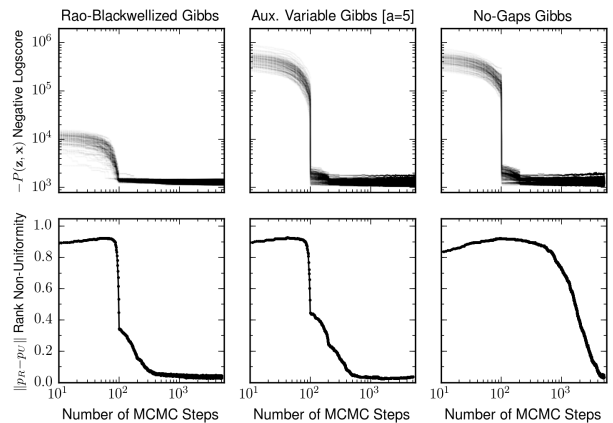


Figure 5: The uniformity of the SRS (bottom row) captures convergence behavior of MCMC sampling algorithms for Dirichlet process mixture models that are not captured by standard diagnostics such as the logscore (top row).

Eq. (11) indicates that by simulating n datasets $\{x_1, \dots, x_n\}$ i.i.d. from the marginal distribution, samples $\{\hat{z}_1, \dots, \hat{z}_n\}$ (where $z_i \approx \pi(z|x_i)$) from an approximate posterior should be i.i.d. samples from the prior $\pi(z)$. An approximate sampler can be thus be diagnosed by performing a goodness-of-fit test to check whether $\hat{z}_{1:n}$ are distributed according to π . Ranks of univariate marginals of a continuous parameter vector $z \in \mathbb{R}^d$ are used in [31]. We extend SBC to handle discrete latent variables z taking values in a large domain.

We sampled $n = 1000$ datasets $\{x_1, \dots, x_n\}$ independently from a Dirichlet process mixture model. Each dataset x_i has $k = 100$ observations and each observation is five-dimensional (i.e., $x_i \in \mathbb{R}^{k \times 5}$) with a Gaussian likelihood. From SBC, samples $\hat{z}_{1:n}$ (where $z_i \in \Pi_k$ and $|\Pi_k| \approx 10^{115}$) of the mixture assignment vector should be distributed according to the CRP prior $\pi(z)$. The top row of Figure 5 shows trace plots of the logscore (unnormalized posterior) of approximate samples from Rao-Blackwellized Gibbs, Auxiliary Variable Gibbs, and No-Gaps Gibbs samplers (Algorithms 3, 8, and 4 in [21]). Each line corresponds to an independent run of MCMC inference. The bottom row shows the evolution of the uniformity of the SRS using $m = 64$ and the ordering on partitions from Algorithm 2.

While logscores typically stabilize after 100 MCMC steps (one epoch through all observations in a dataset) and suggest little difference across the three samplers, the SRS shows that Rao-Blackwellized Gibbs is slightly more efficient than Auxiliary Variable Gibbs and that the sample quality from No-Gaps Gibbs is inferior to those from the other two algorithms up until roughly 5,000 steps. These results are consistent with the observation from [21] that No-Gaps has inefficient mixing (it excessively rejects proposals on singleton clusters).

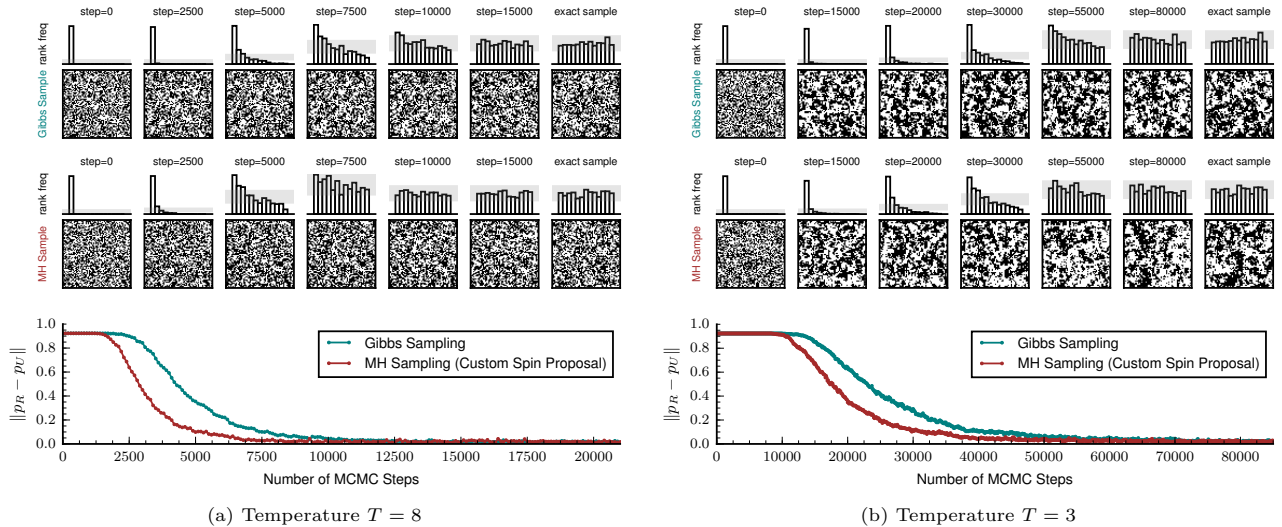


Figure 6: Assessing the goodness-of-fit of approximate samples of a 64×64 Ising model for Gibbs sampling and Metropolis–Hastings sampling (with the custom spin proposal from [18]) at two temperatures using the SRS. In both cases, the SRS converges to its uniform distribution more rapidly for samples obtained from MH than for those from Gibbs sampling.

5.2 Ising models

In this application we use the SRS to assess the sample quality of approximate Ising model simulations. For a ferromagnetic $k \times k$ lattice with temperature T , the probability of a spin configuration $x \in \{-1, +1\}^{k \times k}$ is

$$P(x) \propto \exp\left(-1/T \sum_{i,j} x_i x_j\right). \quad (12)$$

While Eq. (12) is intractable to compute for any x due to the unknown normalization constant, coupling-from-the-past [25] is a popular MCMC technique which can tractably obtain exact samples from the Ising model. For a 64×64 Ising model (domain size $2^{64 \times 64}$), we obtained 650 exact samples using coupling-from-the-past, and used these “ground-truth” samples to assess the goodness-of-fit of approximate samples obtained via Gibbs sampling and Metropolis–Hastings sampling (with a custom spin proposal [18, Section 31.1]).

For each temperature $T = 3$ and $T = 8$, we obtained 7,800 approximate samples using MH and Gibbs. The first two rows of Figure 6 each show the evolution of one particular sample (Gibbs, top; MH, bottom). Two exact samples are shown in the final column of each panel. All approximate and exact samples are independent of one another, obtained by running parallel Markov chains. The SRS of the exact samples with respect to the approximate samples was taken at checkpoints of 100 MCMC steps, using $m = 12$ and an ordering based on the Hamiltonian energy, spin magnetization, and connected components. SRS histograms (and 99% variation bands) evolving at various steps are shown above the Ising model renderings.

The SRS is non-uniform (including in regimes where the difference between approximate and exact samples is too fine-grained to be detected visually) at early steps and more uniform at higher steps. The plots show that MH is a more efficient sampler than Gibbs at moderate temperatures, as its sample quality improves more rapidly. This characteristic was conjectured in [18], which noted that the MH sampler “has roughly double the the probability of accepting energetically unfavourable moves, so may be a more efficient sampler [than Gibbs]”. In addition, the plots suggest that the samples become close to exact (in terms of their joint energy, magnetization, and connected components characteristics) after 20,000 steps for $T = 8$ and 100,000 steps for $T = 3$, even though obtaining exact samples using coupling-from-the-past requires between 500,000 and 1,000,000 MCMC steps for both temperatures.

6 Conclusion

This paper has presented a flexible, simple-to-implement, and consistent goodness-of-fit test for discrete distributions. The test statistic is based on the ranks of observed samples with respect to new samples from the candidate distribution. The key insight is to compute the ranks using an ordering on the domain that is able to detect differences in properties of interest in high dimensions. Unlike most existing statistics, the SRS is distribution-free and has a simple exact sampling distribution. Empirical studies indicate that the SRS is a valuable addition to the practitioner’s toolbox for assessing sample quality in regimes which are not easily handled by existing methods.

Acknowledgments

The authors thank the anonymous referees for their helpful feedback. This research was supported by the DARPA SD2 program (contract FA8750-17-C-0239); the Ethics and Governance of Artificial Intelligence Initiative of the MIT Media Lab and Harvard's Berkman Klein Center; the Systems That Learn Initiative of MIT CSAIL; and an anonymous philanthropic gift.

References

- [1] Jayadev Acharya, Constantinos Daskalakis, and Gautam Kamath. Optimal testing for properties of distributions. In *Advances in Neural Information Processing Systems 28 (NIPS)*, pages 3591–3599. Curran Associates, 2015.
- [2] Mohammad Ahsanullah, Valery B. Nevzorov, and Mohammad Shakil. *An Introduction to Order Statistics*. Atlantis Studies in Probability and Statistics. Atlantis Press, 2013.
- [3] Erling B. Andersen. A goodness of fit test for the Rasch model. *Psychometrika*, 38(1):123–140, 1973.
- [4] Taylor B. Arnold and John W. Emerson. Non-parametric goodness-of-fit tests for discrete null distributions. *The R Journal*, 3(2), 2011.
- [5] Tuğkan Batu, Lance Fortnow, Ronitt Rubinfeld, Warren D. Smith, and Patrick White. Testing that distributions are close. In *Proceedings of the 41st Annual Symposium on Foundations of Computer Science (FOCS)*, pages 259–269. IEEE, 2000.
- [6] Wray Buntine and Marcus Hutter. A Bayesian view of the Poisson–Dirichlet process. *arXiv preprint*, (arXiv:1007.0296), 2010.
- [7] V. Choulakian, R. A. Lockhart, and M. A. Stephens. Cramér–von Mises statistics for discrete distributions. *Canadian Journal of Statistics*, 22(1):125–137, 1994.
- [8] W. J. Conover. A Kolmogorov goodness-of-fit test for discontinuous distributions. *Journal of the American Statistical Association*, 67(339):591–596, 1972.
- [9] J. Dehardt. Generalizations of the Glivenko–Cantelli theorem. *The Annals of Mathematical Statistics*, 42(6):2050–2055, 1971.
- [10] Andrew Gelman, Xiao-Li Meng, and Hal Stern. Posterior predictive assessment of model fitness via realized discrepancies. *Statistica Sinica*, 6:733–807, 1996.
- [11] Phillip I. Good. *Permutation, Parametric, and Bootstrap Tests of Hypotheses*. Springer Series in Statistics. Springer, 2004.
- [12] Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(138):723–773, 2012.
- [13] Wassily Hoeffding. Asymptotically optimal tests for multinomial distributions. *The Annals of Mathematical Statistics*, 36(2):369–401, 1965.
- [14] Susan Dadakis Horn. Goodness-of-fit tests for discrete data: A review and an application to a health impairment scale. *Biometrics*, 33(1):237–247, 1977.
- [15] Erich L. Lehmann and Howard J. D’Abrera. *Non-parametrics: Statistical Methods Based on Ranks*. Holden-Day Series in Probability and Statistics. Holden-Day, 1975.
- [16] Erich L. Lehmann and Joseph P. Romano. *Testing Statistical Hypotheses*. Springer Texts in Statistics. Springer, 3rd edition, 2005.
- [17] Cathryn M. Lewis and Jo Knight. Introduction to genetic association studies. In Ammar Al-Chalabi and Laura Almasy, editors, *Genetics of Complex Human Diseases: A Laboratory Manual*. Cold Spring Harbor Laboratory Press, 2009.
- [18] David J. C. MacKay. *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, 2003.
- [19] Henry B. Mann and Donald R. Whitney. On a test of whether one of two random variables is stochastically larger than the other. *The Annals of Mathematical Statistics*, 18(1):50–60, 1947.
- [20] Alberto Maydeu-Olivares and Carlos Garcia-Forero. Goodness-of-fit testing. In Penelope Peterson, Eva Baker, and Barry McGaw, editors, *International Encyclopedia of Education*, volume 7, pages 190–196. Elsevier, 2010.
- [21] Radford M. Neal. Markov chain sampling methods for Dirichlet process mixture models. *Journal of Computational and Graphical Statistics*, 9(2):249–265, 2000.
- [22] J. A. Peacock. Two-dimensional goodness-of-fit testing in astronomy. *Monthly Notices of the Royal Astronomical Society*, 202(3):615–627, 1983.
- [23] Karl Pearson. On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling. *Philosophical Magazine*, 5:157–175, 1900.
- [24] Anthony N. Pettitt and Michael A. Stephens. The Kolmogorov–Smirnov goodness-of-fit statistic with discrete and grouped data. *Technometrics*, 19(2): 205–210, 1977.

- [25] James G. Propp and David B. Wilson. Exact sampling with coupled Markov chains and applications to statistical mechanics. *Random Structures & Algorithms*, 9(2):223–252, 1996.
- [26] Timothy R. C. Read and Noel A. C. Cressie. *Goodness-of-Fit Statistics for Discrete Multivariate Data*. Springer Series in Statistics. Springer, 1988.
- [27] Walter Rudin. *Principles of Mathematical Analysis*. International Series in Pure and Applied Mathematics. McGraw-Hill, 1976.
- [28] D. S. Starnes, D. Yates, and D. S. Moore. *The Practice of Statistics*. W. H. Freeman and Company, 2010.
- [29] Michael Steele and Janet Chaseling. Powers of discrete goodness-of-fit test statistics for a uniform null against a selection of alternative distributions. *Communications in Statistics—Simulation and Computation*, 35(4):1067–1075, 2006.
- [30] Brett Stevens and Aaron Williams. The coolest order of binary strings. In *Proceedings of the 6th International Conference on Fun with Algorithms (FUN)*, pages 322–333. Springer, 2012.
- [31] Sean Talts, Michael Betancourt, Daniel Simpson, Aki Vehtari, and Andrew Gelman. Validating Bayesian inference algorithms with simulation-based calibration. *arXiv preprint*, (arXiv:1804.06788), 2018.
- [32] Gregory Valiant and Paul Valiant. Estimating the unseen: An $n/\log(n)$ -sample estimator for entropy and support size, shown optimal via new CLTs. In *Proceedings of the 43rd ACM Symposium on Theory of Computing (STOC)*, pages 685–694. ACM, 2011.
- [33] D. A. Williams. Improved likelihood ratio tests for complete contingency tables. *Biometrika*, 63(1):33–37, 1976.
- [34] Michael Williams. How good are your fits? Unbinned multivariate goodness-of-fit tests in high energy physics. *Journal of Instrumentation*, 5(09):P09004, 2010.
- [35] Jiasen Yang, Qiang Liu, Vinayak Rao, and Jennifer Neville. Goodness-of-fit testing for discrete distributions via Stein discrepancy. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, volume 80 of *Proceedings of Machine Learning Research*, pages 5561–5570. PMLR, 2018.