# Generalised Bayesian Inference for Discrete Intractable Likelihood

Takuo Matsubara[1]   Jeremias Knoblauch[2]   François-Xavier Briol[2,4]   Chris. J. Oates[3,4]

[1]The University of Edinburgh, UK
[2]University College London, UK
[3]Newcastle University, UK
[4]The Alan Turing Institute, UK

**Abstract**

Discrete state spaces represent a major computational challenge to statistical inference, since the computation of normalisation constants requires summation over large or possibly infinite sets, which can be impractical. This paper addresses this computational challenge through the development of a novel generalised Bayesian inference procedure suitable for discrete intractable likelihood. Inspired by recent methodological advances for continuous data, the main idea is to update beliefs about model parameters using a discrete Fisher divergence, in lieu of the problematic intractable likelihood. The result is a generalised posterior that can be sampled from using standard computational tools, such as Markov chain Monte Carlo, circumventing the intractable normalising constant. The statistical properties of the generalised posterior are analysed, with sufficient conditions for posterior consistency and asymptotic normality established. In addition, a novel and general approach to calibration of generalised posteriors is proposed. Applications are presented on lattice models for discrete spatial data and on multivariate models for count data, where in each case the methodology facilitates generalised Bayesian inference at low computational cost.

## 1 Introduction

This paper focuses on statistical models for data defined on a discrete set $\mathcal{X}$, whose probability mass function $p_\theta$ involves a parameter $\theta$ to be inferred. In this setting, there is an urgent need for computational methodology applicable to models that are *intractable*, in the specific sense that

$$p_\theta(\boldsymbol{x}) = \frac{\tilde{p}_\theta(\boldsymbol{x})}{Z_\theta}, \qquad Z_\theta := \sum_{\boldsymbol{x} \in \mathcal{X}} \tilde{p}_\theta(\boldsymbol{x}), \tag{1}$$

where the positive function $\tilde{p}_\theta$ is straightforward to evaluate but direct computation of the normalising constant $Z_\theta \in (0, \infty)$ is impractical. This situation is ubiquitous in the discrete data context, since it is often impractical to compute a sum over a large or infinite discrete set. To limit scope, this paper considers generalised Bayesian inference where, to date, several computational approaches have been proposed. These approaches, which are recalled in Section 2, are mainly applicable in settings where it is possible to simulate data $\boldsymbol{x}$, conditional on the parameter $\theta$. However, in several of the most scientifically important instances of (1), exact (or even approximate) simulation from the model is not practical.

Important examples of statistical models exhibiting these computational challenges include lattice models of spatial data [Moores et al., 2020], statistical models for graph-valued data [Lusher et al., 2013], and statistical models for multivariate count data [Inouye et al., 2017]. In each case, the normalising constant involves summation over a set whose cardinality is exponential in the dimension of the lattice, in the size of the nodal set of the graph, or even infinite, rendering direct computation and simulation of data intractable in general.

To circumvent both computation of the normalising constant and simulation from the statistical model, Matsubara et al. [2022] proposed a generalised Bayesian posterior, called *KSD-Bayes*, which is based on a Stein discrepancy. The resulting generalised posterior is consistent and asymptotically normal, and thus shares many of the properties of the standard Bayesian posterior whilst admitting a form which does not require the computation of an intractable normalisation constant. However, a major limitation of KSD-Bayes is the dependence of the generalised posterior on a user-specified symmetric positive definite function, called a kernel, which determines precisely how beliefs are updated. In continuous domains, such as $\mathbb{R}^d$, there are several natural choices of kernel available, and their associated Stein discrepancies have been well-studied [Anastasiou et al., 2023]. However, in discrete domains there are often no natural choices of kernel, or when natural choices exists [such as a heat kernel; Chung and Graham, 1997] they can be computationally impractical.

This paper presents *DFD-Bayes*, the first generalised Bayesian inference method tailored to inference with discrete intractable likelihood. The approach is based on a novel discrete version of the Fisher divergence which, in contrast to KSD-Bayes, does not require a kernel to be specified. Further, the DFD-Bayes posterior has computational complexity $O(nd)$, where $n$ is the number of data and $d$ is the data dimension, which compares favourably to the KSD-Bayes computational complexity of $O(n^2 d)$. The DFD-Bayes methodology is supported by asymptotic guarantees, presented in Section 3, and empirical results, in Section 4, demonstrate state-of-the-art performance in the applications considered. Before setting out the proposed methodology, we first we review related work in Section 2.

## 2    Background

The aim of this section is to briefly review existing Bayesian and generalised Bayesian methodology for intractable statistical models, extending the discussion to include both continuous and discrete data. Frequentist estimation for intractable models is not discussed [we refer the reader to e.g. Hyvärinen, 2005].

**Approximate Likelihood**    Faced with an intractable model, a pragmatic approach is simply to employ standard Bayesian inference with a tractable approximation to the likelihood [e.g. Bhattacharyya and Atchade, 2019]. A classical example of approximate likelihood is the pseudolikelihood of Besag [1974], which replaces the joint probability mass function of the data with a product of conditional probability mass functions, each of which is sufficiently low-dimensional (or otherwise tractable enough) to permit normalising constants to be computed. Generalisations of this approach are sometimes referred to as composite likelihood [Varin et al., 2011]. These approximations are usually model-specific, and analysis of the approximation error may be difficult in general [Lindsay et al., 2011].

**Simulation-Based Methods**  One class of intractable statistical models that has been explored in detail are models for which it is possible to simulate data $\boldsymbol{x}$ conditional on the parameter $\theta$. A well-known approach to inference in this class of models is the exchange algorithm of Møller et al. [2006] and Murray et al. [2006], which constructs a Markov chain on an extended state space for which the standard Bayesian posterior occurs as a marginal. Simulation of the Markov chain requires both exact simulation from the statistical model and evaluation of $\tilde{p}_\theta(\boldsymbol{x})$. Further methodological development has been focused on removing the requirement to evaluate $\tilde{p}_\theta(\boldsymbol{x})$, with approximate Bayesian computation [Marin et al., 2012], Bayesian synthetic likelihood [Price et al., 2018], MMD-Bayes [Cherief-Abdellatif and Alquier, 2020, Pacchiardi and Dutta, 2021] and the posterior boostrap [Dellaporta et al., 2022] emerging as likelihood-free methods, which require only that data can be simulated. Unfortunately, for many statistical models of discrete data, exact simulation [the state-of-the-art being e.g. Propp and Wilson, 1998] from the model is impractical.

**Markov Chain-Based Methods**  Another pragmatic approach is to substitute exact simulations with approximate simulations, such as obtained from a Markov chain. This idea works in specific instances; see the review of Park and Haran [2018]. The main drawback of these approaches, as far as this paper is concerned, is that they require the design of a rapidly mixing Markov chain on a possibly large (or infinite) discrete set. As such, these methods require bespoke implementations for each class of statistical model considered, and for many models of interest appropriate Markov chains have yet to be developed. Thus Markov chain-based methods do not represent a general solution to discrete intractable likelihood.

**Russian Roulette**  The pseudo-marginal approach justifies replacing the intractable likelihood $p_\theta(\boldsymbol{x})$ with a positive unbiased estimator $\hat{p}_\theta(\boldsymbol{x})$ of the likelihood in the context of a Metropolis–Hastings algorithm [Andrieu and Roberts, 2009]. The practical difficulty of this approach is to construct a positive unbiased estimator. Lyne et al. [2015] proposed the Russian roulette estimator for intractable statistical models, a simulation technique from the physics literature which involves random truncation of the sum (or of an integral in the continuous context) defining the normalising constant. The Russian roulette estimator is unbiased but is not guaranteed to be positive, meaning that post hoc re-weighting of the Markov chain sample path is required. The ergodicity of Russian roulette has not, to the best of our knowledge, been theoretically studied. Further, the mixing time of the Markov chain is known to be sensitive to the variance of $\hat{p}_\theta(\boldsymbol{x})$, which can be large for estimators based on random truncation (especially when there is no clear a priori ordering for the summands, which can occur in the discrete context). As such, the pseudo-marginal approach does not at present represent a general computational solution to intractable likelihood.

**Generalised Bayesian Inference**  Motivated by the absence of general computational methodology for intractable likelihood, Matsubara et al. [2022] proposed a solution called KSD-Bayes. The setting for this approach is the nascent field of generalised Bayesian inference. Given a prior $\pi(\theta)$, a dataset $\{\boldsymbol{x}_i\}_{i=1}^n \subset \mathcal{X}$, and a constant $\beta > 0$, generalised Bayesian inference updates beliefs using a loss function $D_n(\theta)$, producing a generalised posterior

$$\pi_n^D(\theta) \propto \pi(\theta) \exp(-\beta D_n(\theta)). \tag{2}$$

The standard posterior is recovered by the negative log-likelihood $D_n(\theta) = -\sum_{i=1}^n \log p_\theta(\boldsymbol{x}_i)$, while several alternative loss functions have been developed to confer robustness in settings where the sta-

tistical model is misspecified (see the survey in Bissiri et al. [2016] for the case of additive loss functions, and Knoblauch et al. [2022] for further generalisation). KSD-Bayes [Matsubara et al., 2022] is distinguished among existing generalised Bayesian inference methods by its applicability to statistical models involving an intractable normalising constant [see also Section 4.2 of Giummolè et al., 2019]. This was achieved by selecting $D_n(\theta)$ to be a Stein discrepancy between the statistical model $p_\theta$ and the empirical distribution of the dataset, which can be computed without the normalising constant. Strikingly, a fully conjugate treatment of the continuous exponential family model, and a straight-forward treatment of the discrete exponential family model using Markov chain Monte Carlo, is possible using KSD-Bayes; this in principle provides a solution to many of the aforementioned instances of intractable likelihood. However, the dependence of KSD-Bayes on a user-specified kernel renders the approach unattractive for discrete domains, where there are often no natural choices of kernel, or where natural choices[1] are computationally impractical. Furthermore, the $O(n^2 d)$ computational cost of KSD-Bayes is super-linear in the size of the dataset.

This paper presents general methodology for inferring the parameters of a intractable discrete statistical model. The main idea is to employ a discrete Fisher divergence as a loss function in a generalised Bayesian inference context. The resulting *DFD-Bayes* method does not require a choice of kernel, enjoys theoretical guarantees, and can be computed at cost $O(nd)$ linear in the size of the dataset. Full details are provided next.

## 3 Methodology

This section presents and analyses DFD-Bayes. First, we present a novel discrete formulation of the Fisher divergence in Section 3.1. DFD-Bayes is introduced in Section 3.2, where posterior consistency and asymptotic normality are established. Section 3.3 presents a novel approach to calibration of generalised posteriors, which may be of independent interest. Limitations of DFD-Bayes are discussed in Section 3.4.

**Notation**   Denote by $\mathcal{X}$ a countable set in which data are contained, and by $\Theta$ the set of permitted values for the parameter $\theta$, where $\Theta$ is a Borel subset of $\mathbb{R}^v$ for some $v \in \mathbb{N}$. Probability distributions on $\mathcal{X}$ are identified with their probability mass functions, with respect to the counting measure on $\mathcal{X}$. The $i$-th coordinate of a function $f : \mathcal{X} \to \mathbb{R}^d$ is denoted by $f_i : \mathcal{X} \to \mathbb{R}$. For a probability distribution $q$ on $\mathcal{X}$ and $d, p \in \mathbb{N}$, denote by $L^p(q, \mathbb{R}^d)$ the Lebesgue space of measurable functions $f : \mathcal{X} \to \mathbb{R}^d$ such that $\sum_{i=1}^d \mathbb{E}_{X \sim q}[|f_i(X)|^p] < \infty$, in which two elements $f, g \in L^p(q, \mathbb{R}^d)$ are identified if they are $p$-almost everywhere equal. The notation $\|\cdot\|$ indicates the Euclidean norm of $\mathbb{R}^m$, and will be applied also to matrices and tensors interpreted, respectively, as elements of $\mathbb{R}^{v \times v}$ and $\mathbb{R}^{v \times v \times v}$. A Dirac measure at $x \in \mathcal{X}$ is denoted by $\delta_x$.

### 3.1   A Discrete Fisher Divergence

The Fisher divergence underpins several frequentist estimators for intractable statistical models, most notably score matching [Hyvärinen, 2005], and has been used in the context of Bayesian model

---

[1]A natural choice is the heat kernel, whose origins lie in spectral graph theory [Chung and Graham, 1997]. However, computation of the heat kernel requires a $O(D^3)$ cost where $D = \text{card}(\mathcal{X})$, which is often impractical. For example, the Ising model on a lattice $\mathcal{X} = \{0, 1\}^d$ has $D = 2^d$, while the Conway–Maxwell–Poisson model of Section 4.1 has $D = \infty$, meaning approximation of the heat kernel would be required.

selection [e.g. Dawid and Musio, 2015]. It is classically defined for continuous domains; for (sufficiently regular) densities $p$ and $q$ on $\mathbb{R}^d$, the Fisher divergence is $\mathrm{FD}(p\|q) = \mathbb{E}_{X\sim q}[\|\nabla \log p(X) - \nabla \log q(X)\|^2]$ where $\nabla$ denotes the gradient operator in $\mathbb{R}^d$. Its main advantage is that it can be computed without knowledge of the normalising constant[2] of $p$ and, furthermore, expectations with respect to $p$ are not required. The Fisher divergence was extended to discrete domains in Lyu [2009], Xu et al. [2022]. However, existing work focuses on domains $\mathcal{X}$ of finite cardinality or one-dimensional models, and a technical contribution of this paper, which may be of independent interest, is to present an extension of Fisher divergence to certain discrete domains which may be a countably infinite set in multiple dimensions. The extended divergence satisfies the requirements of a proper local scoring rule and thus complements existing scoring rule methodology developed in the finite domain context in Dawid et al. [2012].

**Standing Assumption 1.** *Let $\mathcal{X} = S_1 \times \cdots \times S_d$, where for each $i = 1, \ldots, d$ there is an order isomorphism $S_i \cong I_i \subseteq \mathbb{Z}$, and $d \in \mathbb{N}$.*

This setting is general enough to include diverse data types, such as multivariate count data, or network data with a fixed vertex set. For any set $S \cong I \subseteq \mathbb{Z}$, precisely one of the following must hold: (i) no smallest or largest elements of $S$ exist; (ii) both a smallest element, $s_{\min}$, and a largest element, $s_{\max}$, exist; (iii) only $s_{\min}$ exists; (iv) only $s_{\max}$ exists. Without loss of generality, we will identify the case (iv) with (iii) by reversing the ordering of $S$. In addition, it will be useful to extend the domains $S_i$ to include an additional state (not part of the ordering), denoted $\star$, and to this end we let $S_i^\star = S_i \cup \{\star\}$ and $\mathcal{X}^\star = S_1^\star \times \cdots \times S_d^\star$. A function $h : \mathcal{X} \to \mathbb{R}$ extends to a function $h : \mathcal{X}^\star \to \mathbb{R}$ by setting $h(\boldsymbol{x}) = 0$ whenever any of the coordinates of $\boldsymbol{x}$ are equal to $\star$.

**Definition 1.** *Let $S \cong I \subseteq \mathbb{Z}$. For consecutive elements $r < s < t$ in $S$ we let $s^- := r$ and $s^+ := t$. If both $s_{\min}$ and $s_{\max}$ exist, we let $s_{\min}^- := s_{\max}$ and $s_{\max}^+ := s_{\min}$ or, if only $s_{\min}$ exists, we let $s_{\min}^- := \star$ and $\star^+ = s_{\min}$. For $\boldsymbol{x} = (x_1, \ldots, x_d) \in \mathcal{X}$, define $\boldsymbol{x}^{i+} := (x_1, \ldots, x_i^+, \ldots, x_d)$ and $\boldsymbol{x}^{i-} := (x_1, \ldots, x_i^-, \ldots, x_d)$.*

Simply put, this ensures that each element $s$ has both a preceding and proceeding element, so that increments and decrements are well-defined. The above structure can be exploited to define an operator for $\mathcal{X}$ that is analogous to the gradient operators for $\mathbb{R}^d$:

**Definition 2.** *For $h : \mathcal{X} \to \mathbb{R}$, define the backward difference operator by*

$$\nabla^- h(\boldsymbol{x}) := \left[ h(\boldsymbol{x}) - h(\boldsymbol{x}^{1-}), \ \cdots \ , \ h(\boldsymbol{x}) - h(\boldsymbol{x}^{d-}) \right]^\top \in \mathbb{R}^d.$$

Based on Definitions 1 and 2, we can construct a divergence applicable to discrete domains $\mathcal{X}$, which we term a *discrete Fisher divergence*. Recall that values of $f \in L^p(q, \mathbb{R}^d)$ in a measure zero domain of $q$ i.e. $\{\boldsymbol{x} \in \mathcal{X} \mid q(\boldsymbol{x}) = 0\}$ are arbitrary and not involved in the integral with respect to $q$ [Rudin, 1987, Remark 1.37, p.29]. In what follows, it is sufficient for functions $(\nabla^- p)/p, (\nabla^- q)/q \in L^2(q, \mathbb{R}^d)$ to be well-defined in the support of $q$.

**Definition 3.** *Let $p$ and $q$ be probability distributions on $\mathcal{X}$, such that $(\nabla^- p)/p, (\nabla^- q)/q \in L^2(q, \mathbb{R}^d)$. The discrete Fisher divergence is defined as*

$$\mathrm{DFD}(p\|q) := \mathbb{E}_{X\sim q}\left[ \left\| \frac{\nabla^- p(X)}{p(X)} - \frac{\nabla^- q(X)}{q(X)} \right\|^2 \right]. \tag{3}$$

---

[2]The Fisher divergence depends only on $\nabla \log p$, equal to the ratio $(\nabla p)/p$, meaning it is sufficient to know $p$ up to a normalising constant.

The choice of a Euclidean norm in (3) is not critical and other norms could be employed, but for expository purposes the standard Euclidean norm will be used throughout. Proposition 1 justifies the name 'divergence' and offers an alternative, computable formula for (3).

**Proposition 1.** *The discrete Fisher divergence satisfies* $\mathrm{DFD}(p\|q) \geq 0$ *for any* $p, q$, *with equality if and only if* $p = q$. *Furthermore, if* $p(\boldsymbol{x}^{j+}) > 0$ *for all* $\boldsymbol{x}$ *and* $j = 1, \ldots, d$ *in the support of* $q$, *it admits the following alternative formula*

$$\mathrm{DFD}(p\|q) = \mathbb{E}_{X \sim q} \left[ \sum_{j=1}^{d} \left( \frac{p(X^{j-})}{p(X)} \right)^2 - 2 \left( \frac{p(X)}{p(X^{j+})} \right) \right] + C(q), \tag{4}$$

*where the term* $C(q) := \mathbb{E}_{X \sim q}[\sum_{j=1}^{d} 1 + (1 - q(X^{j-})/q(X))^2]$ *is* $p$-*independent.*

The proof is provided in Appendix B.1. Note that $\mathrm{DFD}(p\|q)$ can be computed without the normalising constant of $p$, analogously to $\mathrm{FD}(p\|q)$ in $\mathbb{R}^d$. All models $p_\theta$ used in this paper are positive on $\mathcal{X}$, for which the assumption $p(\boldsymbol{x}^+) > 0$ in Proposition 1 is automatically satisfied. From Proposition 1, the discrete Fisher divergence between a model $p_\theta$ and an empirical distribution $p_n = \frac{1}{n} \sum_{i=1}^{n} \delta_{\boldsymbol{x}_i}$ corresponding to data $\{\boldsymbol{x}_i\}_{i=1}^{n}$, is computed as

$$\mathrm{DFD}(p_\theta\|p_n) \overset{\theta}{=} \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{d} \left( \frac{p_\theta(\boldsymbol{x}_i^{j-})}{p_\theta(\boldsymbol{x}_i)} \right)^2 - 2 \left( \frac{p_\theta(\boldsymbol{x}_i)}{p_\theta(\boldsymbol{x}_i^{j+})} \right) \tag{5}$$

where $\overset{\theta}{=}$ indicates equality up to an additive, $\theta$-independent constant. In contrast to the continuous Fisher divergence, the $\theta$-independent constant $C(p_n) = \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{d} 1 + (1 - p_n(\boldsymbol{x}_i^{j-})/p_n(\boldsymbol{x}_i))^2$ is well-defined for an empirical density $p_n$ in the discrete Fisher divergence.

**Remark 1.** *The computational cost associated with evaluation of* (5) *is* $O(nd)$, *which improves on the* $O(n^2 d)$ *cost of kernel Stein discrepancy. Furthermore, if* $\mathcal{X}$ *is a finite set and count data are provided, indicating the number of times each of the elements of* $\mathcal{X}$ *occurred, then the complexity of* (5) *reduces to* $O(d)$, *independent of the size of the dataset.*

**Remark 2.** *The discrete Fisher divergence can also be interpreted as a Stein discrepancy constructed based on an* $L^2$-*ball Stein set [Barp et al., 2019]. This implies that discrete Fisher divergence is stronger than popular kernel Stein discrepancies; see Appendix D.*

## 3.2 A Generalised Posterior

We are now in a position to present DFD-Bayes.

**Definition 4** (DFD-Bayes). *Given a prior distribution* $\pi$ *on* $\Theta$, *a statistical model* $p_\theta : \mathcal{X} \to (0, \infty)$ *parametrised by* $\theta \in \Theta$, *and a dataset* $\{\boldsymbol{x}_i\}_{i=1}^{n}$, *the DFD-Bayes posterior is*

$$\pi_n^D(\theta) \propto \pi(\theta) \exp\left(-\beta n \, \mathrm{DFD}(p_\theta\|p_n)\right), \tag{6}$$

*where* $\beta \in (0, \infty)$ *is a constant to be specified.*

This is clearly a special case of the generalised posterior in (2) with $D_n(\theta) = n\,\mathrm{DFD}(p_\theta\|p_n)$. The $\theta$-independent constant $C(p_n)$ of $\mathrm{DFD}(p_\theta\|p_n)$ will be cancelled out by normalisation of the DFD-Bayes posterior. It is thus sufficient to use (5) in place of $\mathrm{DFD}(p_\theta\|p_n)$ for computation. The role of $n$ in (6) is to ensure correct scaling of the generalised posterior as $n \to \infty$ limit, while the appropriate choice of $\beta$ is crucial in calibrating the coverage of the generalised posterior at finite $n$, and will be discussed in Section 3.3. Appendix A contains a detailed worked example of the DFD-Bayes posterior and a comparison with other posteriors using simple tractable models. For the moment, two important properties are highlighted:

**Remark 3.** *In contrast to standard posteriors for intractable likelihoods, the DFD-Bayes posterior is directly amenable to standard Markov chain Monte Carlo because* (5) *is independent of the intractable constant, with the cost of evaluating* (6) *as low as $O(d)$ (c.f. Remark 1).*

**Remark 4.** *In contrast to KSD-Bayes, DFD-Bayes is invariant to order-preserving transformations of the data. Note that the discrete Fisher divergence upper bounds the kernel Stein discrepancies; see Appendix D.3.*

The asymptotic behaviour of the standard Bayesian posterior is well-understood, with sufficient conditions for posterior consistency and asymptotic normality providing frequentist justification for Bayesian inference in the large data limit. Our attention now turns to establishing analogous conditions for DFD-Bayes.

**Standing Assumption 2.** *The data $\{\boldsymbol{x}_i\}_{i=1}^n$ consist of independent samples from a probability distribution $p$ on $\mathcal{X}$. The distribution $p$ and the statistical model $p_\theta$ for these data satisfy $(\nabla^- p)/p, (\nabla^- p_\theta)/p_\theta \in L^2(p, \mathbb{R}^d)$, for all $\theta \in \Theta$.*

The setting of independent data is broad enough to contain important examples of discrete intractable likelihood, including the models studied in Section 4. The other assumption simply ensures that $\mathrm{DFD}(p_\theta\|p_n)$ is well-defined, due to Proposition 1. In this setting, a natural first requirement is that the statistical model is identifiable in the large data limit:

**Assumption 1.** *There exists a unique minimiser $\theta_*$ of $\theta \mapsto \mathrm{DFD}(p_\theta\|p)$ and there exists a sequence $\{\theta_n\}_{n=1}^\infty$ such that $\theta_n$ minimises $\theta \mapsto \mathrm{DFD}(p_\theta\|p_n)$ almost surely for all $n$ sufficiently large. Further, there exists a bounded convex open set $U \subseteq \Theta$ such that $\theta_* \in U$ and $\theta_n \in U$ almost surely for all $n$ sufficiently large.*

The existence of $U$ in Assumption 1 essentially implies that for large enough $n$, we can restrict our theoretical analysis to a bounded subset $U \subseteq \Theta$. This is not restrictive: it can be enforced by re-parameterising the model $p_\theta$ so that its new parameter space is bounded and convex.[3] The existence of $\{\theta_n\}_{n=1}^\infty$ and $\theta_*$ is more difficult to assess in practice, since the true data generating distribution $p$ is unknown. That being said, assuming their existence is common in the asymptotic analysis of Bayesian procedures [see e.g. van der Vaart, 1998, Section 10]. It is worth highlighting that Assumption 1 does not require the model family $\{p_\theta \mid \theta \in \Theta\}$ to contain $p$, which is in contrast to the assumptions needed for the classical asymptotic normality result [van der Vaart, 1998, Theorem 10.1]. On the other hand, if the model family $\{p_\theta \mid \theta \in \Theta\}$ contains $p$ uniquely, existence

---

[3]For example, we can re-parameterise any unbounded parameter $\kappa$ through the logistic function and define the invertible transformation $\theta = (1 + e^{-\kappa})^{-1} \in [0, 1]$.

of $\theta_*$ is immediate since the discrete Fisher divergence is a divergence and hence $\mathrm{DFD}(p_\theta\|p) = 0$ if and only if $p_\theta = p$.

Our second main requirement is a technical condition on the derivatives and moments of the model, to ensure that the asymptotic limit is well-defined. It is helpful to introduce the shorthand $r_{j-}(\boldsymbol{x},\theta) := p_\theta(\boldsymbol{x}^{j-})/p_\theta(\boldsymbol{x})$. For a function $g : \Theta \to \mathbb{R}$, let $\nabla_\theta^2 g(\theta) \in \mathbb{R}^{v \times v}$ with entries $\partial_i \partial_j g(\theta)$, and let $\nabla_\theta^3 g(\theta) \in \mathbb{R}^{v \times v \times v}$ with entries $\partial_i \partial_j \partial_k g(\theta)$.

**Assumption 2.** *Assume that $\theta \mapsto p_\theta(\boldsymbol{x})$ is three times continuously differentiable in $U$ for any $\boldsymbol{x} \in \mathcal{X}$, and*

$$\mathbb{E}_{X \sim p}\left[\sup_{\theta \in U} \|\nabla_\theta^s r_{j-}(X^{j+}, \theta)\|\right] < \infty \qquad and \qquad \mathbb{E}_{X \sim p}\left[\sup_{\theta \in U} \|\nabla_\theta^s (r_{j-}(X, \theta)^2)\|\right] < \infty$$

*for all $j = 1, \ldots, d$ and $s = 1, 2, 3$.*

In contrast to Assumption 1, it is easier to verify Assumption 2, as illustrated in Example 1. It considers the exponential family, a large class of models which encompasses the models in our experiments in Section 4. For example, any model on a space $\mathcal{X}$ of finite cardinality is an exponential family model [Amari, 2016, Ch. 2.2.2].

**Example 1** (Exponential Family). *Consider an exponential family model $p_\theta(\boldsymbol{x}) \propto \exp(\eta(\theta) \cdot T(\boldsymbol{x}) + b(\boldsymbol{x}))$, where $\eta : \Theta \to \mathbb{R}^k$, $T : \mathcal{X} \to \mathbb{R}^k$ and $b : \mathcal{X} \to \mathbb{R}$ for some $k \in \mathbb{N}$. For this model, we have $r_{j-}(\boldsymbol{x}, \theta) = \exp(\eta(\theta) \cdot (T(\boldsymbol{x}^{j-}) - T(\boldsymbol{x})) + b(\boldsymbol{x}^{j-}) - b(\boldsymbol{x}))$. Assumption 2 is satisfied if, for $j = 1, \ldots, d$, (i) $\|\eta(\theta)\|$ and $\|\nabla_\theta^s \eta(\theta)\|$ for $s = 1, 2, 3$ are bounded over $\theta \in U$, (ii) $\|T(\boldsymbol{x}^{j-}) - T(\boldsymbol{x})\|$ is bounded over $\boldsymbol{x} \in \mathcal{X}$, and (iii) $\mathbb{E}_{X \sim p}[\exp(b(X^{j-}) - b(X))^2] < \infty$. The requirements (ii) and (iii) are immediate if $\mathcal{X}$ is a finite set.*

The calculations that accompany Example 1 are provided in Appendix E.1. The following theorem establishes that both consistency and asymptotic normality hold. The former implies that our generalised posterior concentrates around the population minimiser $\theta_*$ with probability 1 when $n \to \infty$. The latter establishes that our generalised posterior is normal around $\theta_*$ in the same asymptotic limit.

**Theorem 1.** *Suppose Assumptions 1 and 2 hold. Assume that the prior $\pi$ is positive and continuous at $\theta_*$. Let $B_\epsilon(\theta_*) := \{\theta \in \Theta \mid \|\theta - \theta_*\|_2 < \epsilon\}$. Then for any $\epsilon > 0$,*

$$\int_{B_\epsilon(\theta_*)} \pi_n^D(\theta)\mathrm{d}\theta \xrightarrow{a.s.} 1 \qquad as\ n \to \infty. \tag{7}$$

*Denote by $\widetilde{\pi}_n^D$ a density on $\mathbb{R}^d$ of a random variable $\sqrt{n}(\theta - \theta_n)$ for $\theta \sim \pi_n^D$. If $H_* := \beta \nabla_\theta^2 \mathrm{DFD}(p_\theta\|p)|_{\theta=\theta_*}$ is nonsingular, then*

$$\int_{\mathbb{R}^p} \left|\widetilde{\pi}_n^D(\theta) - \frac{1}{\sqrt{\det(2\pi H_*^{-1})}}\exp\left(-\frac{1}{2}\theta \cdot H_*\theta\right)\right| \mathrm{d}\theta \to 0 \qquad as\ n \to \infty. \tag{8}$$

The proof of Theorem 1 is provided in Appendix B.2. The result was established using similar arguments from early work by Hooker and Vidyashankar [2014], Ghosh and Basu [2016] and extended techniques of Miller [2021], Matsubara et al. [2022].

8

## 3.3 A New Approach to Calibration of Generalised Posteriors

The weight $\beta$ in (2) controls the scale of the generalised posterior, and the selection of an appropriate value for $\beta$ is critical to ensure the generalised posterior is calibrated. The literature on this topic is under-developed, but two existing approaches stand out. The first approach was proposed in the recent review paper of Syring and Martin [2019]. It consists of a new stochastic sequential update algorithm for choosing $\beta$, such that a 95% highest posterior density region coincides with a 95% confidence interval. Unfortunately, this approach leads to a large computational cost and is therefore often impractical. The second approach is due to Lyddon et al. [2019] and consists in setting $\beta$ such that the scale of the posterior's asymptotic covariance matrix coincides with that of a frequentist counterpart with correct coverage. Matsubara et al. [2022] numerically showed that this approach is unstable when $n$ is not large enough or when $\theta$ is high dimensional. In addition, the second approach does not take the prior $\pi$ into account, because it depends only the generalised posterior's asymptotic covariance matrix.

In order to remedy some of these issues, the present paper proposes a novel selection criterion for $\beta$ that can be viewed as a more analytically tractable alternative to Syring and Martin [2019]. This criterion is applicable to generalised posteriors beyond DFD-Bayes and may therefore be of independent interest. Our approach consists of two steps: (i) computing minimisers of $B$ "bootstrapped" losses and (ii) estimating an appropriate value of $\beta$ using the closed-form expression in Theorem 2. In contrast to Syring and Martin [2019], step (ii) is non-iterative and exact. Additionally, computation of each minimiser in step (i) is embarrassingly parallel. Relative to the approach of Lyddon et al. [2019], the advantage of our method is that it does not rely on asymptotic quantities, takes the prior into account, and maintains numerical stability even if the parameter $\theta$ is high-dimensional.

To describe the method we first define the minimiser $\theta_n \in \arg\min_{\theta \in \Theta} D_n(\theta)$, where $D_n$ is a loss function based on a dataset $\{\boldsymbol{x}_i\}_{i=1}^n$. To make the dependence on $\beta$ explicit, we denote the posterior $\pi_n^D$ by $\pi_{n,\beta}^D$. In step (i), bootstrap datasets $\{\boldsymbol{x}_i^{(b)}\}_{i=1}^n$, $b = 1, \ldots, B$, are generated by sampling each $\boldsymbol{x}_i^{(b)}$ uniformly with replacement from the original dataset. Then, for each bootstrap dataset, we compute a minimiser $\theta_n^{(b)} = \arg\min_{\theta \in \Theta} D_n^{(b)}(\theta)$, where the superscript indicates that $D_n^{(b)}$ is based on the $b^{\text{th}}$ bootstrap dataset. This leads to an empirical measure $\delta_\theta^B = \frac{1}{B} \sum_{b=1}^B \delta(\theta_n^{(b)})$ which approximates the sampling distribution of the estimator $\theta_n$. In step (ii), we choose $\beta$ to minimise a statistical divergence between $\pi_{n,\beta}^D$ and $\delta_\theta^B$. However, this is not straight-forward, since the majority of statistical divergences (e.g. Kullback–Liebler divergence) require the normalising constant of $\pi_{n,\beta}^D$ for every $\beta$. Interestingly, this is the same computational challenge posed by intractable likelihood. Our proposal is therefore to employ a divergence that circumvents computational of the normalisation constant; here we minimise the score matching loss in the continuous domain $\Theta$ [Hyvärinen, 2005]:

$$\beta_* \in \arg\min_{\beta > 0} \frac{1}{n} \sum_{b=1}^B \left\| \nabla \log \pi_{n,\beta}^D\big(\theta_n^{(b)}\big) \right\|^2 + 2\operatorname{Tr}\left(\nabla^2 \log \pi_{n,\beta}^D\big(\theta_n^{(b)}\big)\right). \tag{9}$$

This leads to an explicit score-matching estimator for $\beta$, circumventing intractability of (2):

**Theorem 2.** *Consider a generalised posterior $\pi_{n,\beta}^D$ with twice differentiable loss function $D_n :$ $\Theta \to \mathbb{R}$. Suppose that there exists at least one $\theta_n^{(b)}$ s.t. $\nabla_\theta D_n(\theta_n^{(b)}) \neq 0$ and that $\sum_{b=1}^B \nabla_\theta D_n(\theta_n^{(b)}) \cdot$*

$\nabla_\theta \log \pi(\theta_n^{(b)}) + \text{Tr}(\nabla_\theta^2 D_n(\theta_n^{(b)})) > 0$. *Then $\beta_*$ in (9) is unique, with*

$$\beta_* = \frac{\sum_{b=1}^B \nabla_\theta D_n(\theta_n^{(b)}) \cdot \nabla_\theta \log \pi(\theta_n^{(b)}) + \text{Tr}(\nabla_\theta^2 D_n(\theta_n^{(b)}))}{\sum_{b=1}^B \|\nabla_\theta D_n(\theta_n^{(b)})\|^2} > 0. \tag{10}$$

The proof is provided in Appendix B.3. The condition in Theorem 2 directly implies existence and positivity of (10). However, in practice, computing (10) and verifying the existence and positivity directly is strikingly easier than validating the local convexity of $D_n$ and $\log \pi$. Note that (10) is straight-forward to compute whenever the loss $D_n$ is amenable to automatic differentiation. For completeness, we also provide an explicit expression in Appendix E.2 for the case of the DFD-Bayes posterior with an exponential family model.

**Remark 5.** *Step (i) of our algorithm is embarrassingly parallelisable over bootstrap samples. Each component inside the sum in (10) can also be computed in parallel during step (ii). Overall, the total cost can be reduced linearly in the number of available cores $K$, and the cost of step (ii) is $O(p^2 \times C \times B/K)$, where $C$ is the cost of evaluating $D_n(\theta)$ and $\pi(\theta)$ at $\theta$.*

### 3.4 Limitations

There are at least two important limitations of the DFD-Bayes methodology, which will now be discussed. First, DFD-Bayes was not derived as an approximation to standard Bayesian inference, and thus the semantics associated with the generalised posterior should not be confused with the semantics of standard Bayesian inference; see Bissiri et al. [2016], Knoblauch et al. [2022] for a detailed discussion of this point. In particular, we need to calibrate DFD-Bayes through the selection of $\beta$, which is not a feature of standard Bayesian inference under well-specified models. Although we expect our bootstrap approach to outperform existing alternative approaches for small sample size $n$, it is possible that in those cases the bootstrap criterion for selecting $\beta$ in Section 3.3 will fail, and in these circumstances the generalised posterior will fail to be calibrated. Second, the generalised posterior may suffer from similar drawbacks to score-based methods for continuous data, including insensitivity to mixing proportions [Wenliang and Kanagawa, 2021]. Indeed, for a two-component mixture model $p_\theta(\boldsymbol{x}) = (1-\theta)p_1(\boldsymbol{x}) + \theta p_2(\boldsymbol{x})$, we can compute the ratios

$$\rho_j := \frac{p_\theta(\boldsymbol{x}^{j-})}{p_\theta(\boldsymbol{x})} = \frac{(1-\theta)p_1(\boldsymbol{x}^{j-}) + \theta p_2(\boldsymbol{x}^{j-})}{(1-\theta)p_1(\boldsymbol{x}) + \theta p_2(\boldsymbol{x})}$$

on which the discrete Fisher divergence is based. Suppose, informally, that the high probability regions $R_1$ of $p_1$ and $R_2$ of $p_2$ are separated, meaning $p_2 \approx 0$ on $R_1$ and $p_1 \approx 0$ on $R_2$. Then these ratios are approximately independent of $\theta$ on $R_1 \cup R_2$, since $\rho_j \approx p_1(\boldsymbol{x}^{j-})/p_1(\boldsymbol{x})$ for $\boldsymbol{x} \in R_1$ and $\rho_j \approx p_2(\boldsymbol{x}^{j-})/p_2(\boldsymbol{x})$ for $\boldsymbol{x} \in R_2$. It follows that $\text{DFD}(p_\theta \| p_n)$ is approximately independent of $\theta$ whenever the data $\{\boldsymbol{x}\}_{i=1}^n \subseteq R_1 \cup R_2$. See Appendix A.4 for an empirical demonstration using a mixture model of two Poisson distributions. Thus, although DFD-Bayes may be applied to mixture models, supported by the theoretical guarantees of Theorem 1, the inferences for mixing proportions so-obtained can be data-inefficient.

## 4  Experimental Assessment

To complement the theoretical assessment we now provide a detailed empirical assessment, focusing on three important instances of discrete intractable likelihood. First, in Section 4.1 we consider
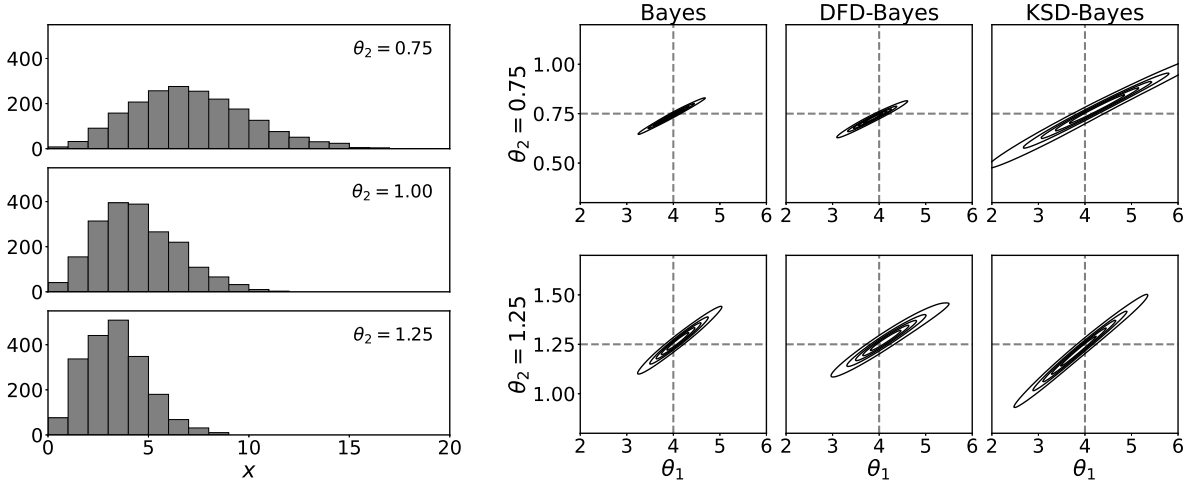
Figure 1: Comparison of standard Bayesian inference with the generalised posteriors from DFD-Bayes and KSD-Bayes on the Conway–Maxwell–Poisson model in the over-dispersed case $\theta_2 = 0.75$ and the under-dispersed case $\theta_2 = 1.25$ for $n = 2,000$.

a relatively simple model for over- and under-dispersed count data, called the Conway–Maxwell–Poisson model. Section 4.2 concerns an application to Ising-type models for discrete spatial data. Finally, we apply DFD-Bayes to perform inference for the parameters of flexible multivariate models for count data in Section 4.3. Source code to reproduce these experiments can be downloaded from `https://github.com/takuomatsubara/Discrete-Fisher-Bayes`.

## 4.1 Conway–Maxwell–Poisson Model

The first model we consider is a generalisation of the Poisson model for over- and under-dispersed count data, due to Conway and Maxwell [1962]. This model is on $\mathcal{X} = \mathbb{N} \cup \{0\}$ (hence $d = 1$ and $\mathrm{card}(\mathcal{X}) = \infty$) and generalises the Poisson distribution through the inclusion of an additional parameter controlling how the data are dispersed. Since the work of Shmueli et al. [2005], this model has been used in a wide range of fields including transport, finance and retail. The model has two parameters $\theta \in \Theta = (0, \infty)^2 \cup ([0, 1] \times \{0\})$ (and hence $p = 2$) and its probability mass function is given by $p_\theta(x) = \tilde{p}_\theta(x) Z_\theta^{-1}$ where $\tilde{p}_\theta(x) = (\theta_1)^x (x!)^{-\theta_2}$. The normalising constant is given by $Z_\theta = \sum_{y=0}^{\infty} \tilde{p}_\theta(y)$, which has no analytical form except for certain special cases of $\theta \in \Theta$, including the case $\theta_2 = 1$ for which the standard Poisson model is recovered.

This model is an ideal test-bed for DFD-Bayes: although the likelihood is formally intractable, it is relatively straightforward to directly approximate the normalising constant[4]. This enables a direct comparison with standard Bayesian inference in the case where the model is well-specified. To this end, we simulated two datasets from the model: (i) an under-dispersed case where $\theta^* = (4, 1.25)$, and (ii) an over-dispersed case where $\theta^* = (4, 0.75)$, shown in Figure 1 (left). Three inference methods were compared: standard Bayesian inference, the KSD-Bayes method of Matsubara et al. [2022], and the DFD-Bayes method we have proposed. The settings of KSD-Bayes are described in Appendix F.1.1. In each case, the prior $\pi$ was taken to be the chi-squared distribution with 3

---

[4]The standard Bayesian inferences reported in this section used the approximation $Z_\theta \approx \sum_{y=0}^{99} \tilde{p}_\theta(y)$ and the associated approximate likelihood. Alternative estimators are available; see Benson and Friel [2021].
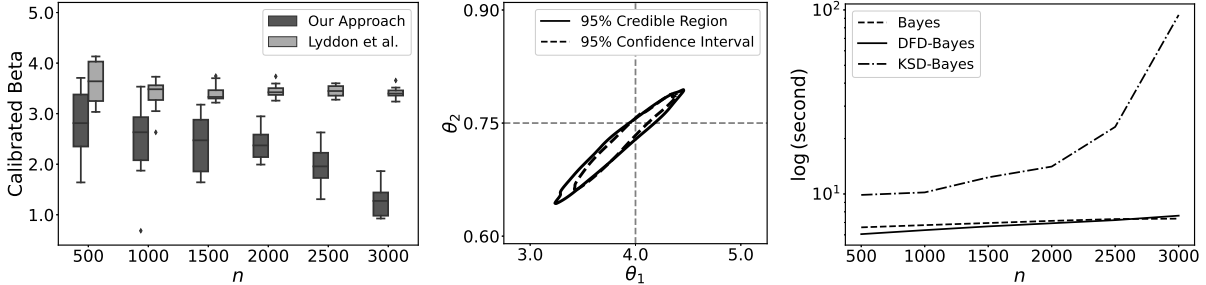
11

Figure 2: Distribution of $\beta_*$ across different realisations of the dataset at each data number $n$ for $\theta_2 = 0.75$ (left), comparison of a 95% credible region of the DFD-Bayes posterior and a 95% confidence interval of the frequentist counterpart for $n = 2000$ (centre), and comparison of computational times of each Metropolis–Hastings algorithm (right). The confidence interval was estimated by a 95% highest probability density region of a kernel density estimator applied to the 100 bootstrap minimisers used in calibration of $\beta$.

degrees of freedom for each of $\theta_1$ and $\theta_2$ independently. A Metropolis–Hastings algorithm was used to sample from all the posteriors; and details can be found in Appendix F.1.2. The weight $\beta$ in DFD-Bayes and KSD-Bayes was calibrated by our approach described in Section 3.3.

Figure 1 (right) illustrates the posteriors, based on typical datasets of size $n = 2,000$. The estimated value of $\beta_*$ was 1.91 for DFD-Bayes and 5.04 for KSD-Bayes in the over-dispersed case $\theta_2 = 0.75$, and 0.46 for DFD-Bayes and 2.51 for KSD-Bayes in the under-dispersed case $\theta_2 = 1.25$. The left panel of Figure 2 displays the distribution of calibrated weight $\beta_*$ as in Section 3.3 over multiple instances of the dataset, along with the values advocated in Lyddon et al. [2019]. For both methods, the calibrated weight is stably estimated.

The inferences obtained using DFD-Bayes resembled those obtained using standard Bayesian inference, irrespective of whether the data were over- or under-dispersed. Those obtained using KSD-Bayes were more conservative than standard Bayes and DFD-Bayes, although the maximum a posteriori estimator approximated the true parameter well. Note that the credible regions of the generalised posteriors can substantially differ from those of standard Bayesian inference; in our approach a credible region of a generalised posterior is calibrated with reference to the distribution of a corresponding frequentist estimator estimated by bootstrapping, leading to approximately correct frequentist coverage as shown in Figure 2 (middle). Calibration led to improved inference outcomes for both DFD-Bayes and KSD-Bayes. In the KSD-Bayes case for example, the value of $\beta_* \geq 1$ intensified the concentration around the true parameter by placing more importance on the loss than the prior. In addition, our approach to calibration is relatively more conservative than Lyddon et al. [2019] because the prior is taken into account.

There is a stark difference in computational cost between DFD-Bayes and KSD-Bayes[5], as demonstrated in the right panel of Figure 2. Indeed, the computational cost of DFD-Bayes is seen to increase linearly with $n$, while the cost of KSD-Bayes increases quadratically.

Finally, to assess performance in a real-world data setting, we apply DFD-Bayes to infer the parameters of a Conway–Maxwell–Poisson model using the sales dataset of Shmueli et al. [2005]. All

---

[5]The cost of standard Bayesian inference in this experiment is entirely determined by the accuracy with which the normalisation constant is approximated; since direct approximation of the normalisation constant is infeasible in general, we do not report this cost.
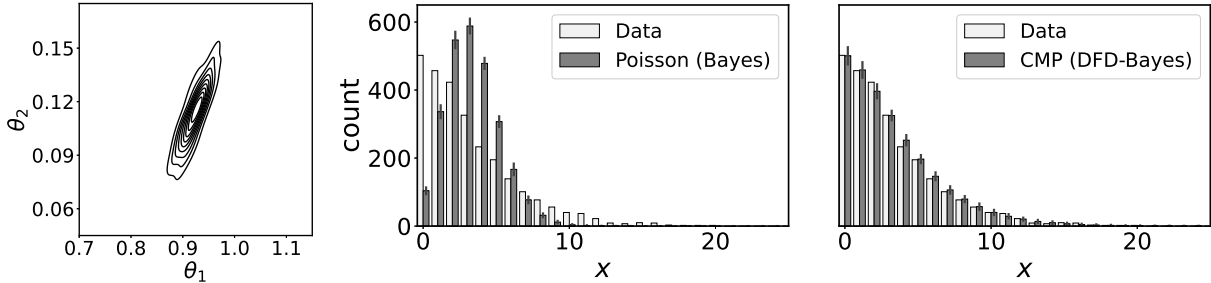
Figure 3: Comparison of DFD-Bayes for the Conway–Maxwell–Poisson model and standard Bayes for the Poisson distribution on the sales data of Shmueli et al. [2005]. Left: The generalised posterior distribution produced using DFD-Bayes. Centre: Posterior predictive distribution, at the level of the data, for a Poisson model with standard Bayesian inference performed. Right: Posterior predictive distribution, at the level of the data, for a Conway–Maxwell–Poisson model with DFD-Bayes inference performed. In both cases, error bars indicate one standard deviation of the posterior predictive distribution.

relevant details are contained in Appendix F.1.3. Figure 3 compares our fitted model to a standard Bayesian analysis using the Poisson distribution, which is the closest analysis one can perform without confronting an intractable likelihood. As observed in the central panel of Figure 3, the Poisson model is not able to capture over-dispersion of the data, whereas the Conway–Maxwell–Poisson model fitted using DFD-Bayes, shown in the right panel, provides a reasonable fit. The DFD-Bayes posterior (left) appears approximately normal, in line with Theorem 1.

## 4.2   Ising Model

The aim of this section is to consider a more challenging instance of discrete intractable likelihood, where the data are high-dimensional (i.e. $d$ is large) and the cardinality of each coordinate domain $S_i$ is small. A small cardinality of $S_i$ is particularly interesting, because the intuition that our difference operators arise from discretisation of continuous differential operators fails to hold. This setting is typified by the Ising model (which has $S_i = \{0, 1\}$), variants of which are used to model diverse phenomena, e.g., the network structure of the amino-acid sequences [Xue et al., 2012]. The computational challenge of performing Bayesian inference for Ising-type models has, to-date, principally been addressed using techniques such as pseudo-likelihood [see the recent survey in Bhattacharyya and Atchade, 2019]. As with the case of generalised Bayesian inference, these do not necessarily lead to the same asymptotic distribution as standard Bayesian inference since the original likelihood is replaced by an approximation [Gong and Samaniego, 1981].

Let $G$ be an undirected graph on a $d$-dimensional vertex set and let $\mathcal{N}_i$ denote the neighbours of node $i$, with self-edges excluded. An Ising model describes a discrete process that assigns each vertex of $G$ either the value 0 or 1, and thus the data domain is $\mathcal{X} = \{0, 1\}^d$. The probability mass function has the exponential family form

$$p_\theta(\boldsymbol{x}) \propto \exp\left(\frac{1}{\theta} \sum_{i=1}^{d} \sum_{j \in \mathcal{N}_i} x_i x_j\right) \tag{11}$$
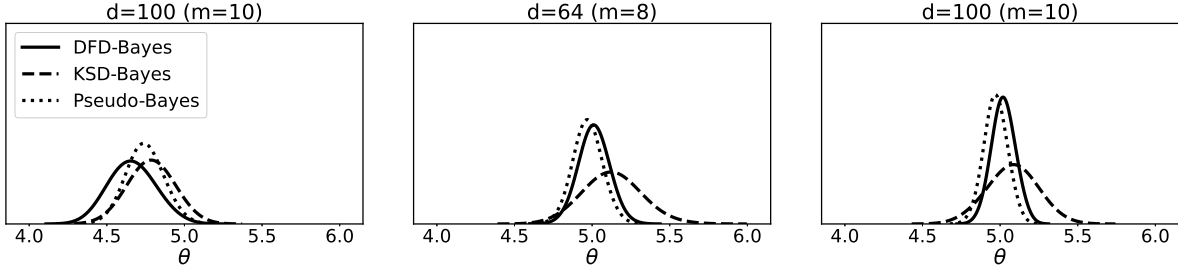
13

Figure 4: Comparison of approximate Bayesian inference based on pseudo-likelihood, DFD-Bayes and KSD-Bayes, applied to the Ising model with $\theta = 5$ for $n = 1,000$ and $d = 10 \times 10$. For all methods, the value $\beta_*$ from Section 3.3 was used.

where $\theta$ is a temperature parameter, controlling the propensity for neighbouring vertices to share a common value. Here we consider the ferromagnetic Ising model, which has $\theta \in (0, \infty)$. To conduct a simulation study we consider the case where $G$ is a $m \times m$ grid. Simulating from Ising models is challenging due to the high-dimensional discrete domain, so here we restrict attention to $m \in \{5, \ldots, 10\}$ to ensure that data were accurately simulated. A total of $n = 1,000$ data points were generated from an Ising model with $\theta = 5$, using an extended run of a Metropolis—Hastings algorithm, the details of which are contained in Appendix F.2.1. A chi-squared prior with degree of freedom 3 was used. Three inference methods were compared: the KSD-Bayes method of Matsubara et al. [2022], the proposed DFD-Bayes method, and standard Bayesian inference based on a the pseudo-likelihood

$$\tilde{p}_\theta(\boldsymbol{x}) = \prod_{i=1}^{d} p_\theta(x_i | \{x_j : j \in \mathcal{N}_i\}),$$

where $p_\theta(x_i | \{x_j : j \in \mathcal{N}_i\})$ is a restriction of the original model (11) to the $i$-th coordinate $x_i$ under the condition $\{x_j : j \in \mathcal{N}_i\}$ that results in a Bernoulli distribution of $x_i$ for each $i = 1, \cdots, d$ [Besag, 1974]. The latter Pseudo-Bayes approach can be viewed as a special case of generalised Bayes inference, since it replaces the original likelihood loss of the model (11) with the pseudo-likelihood loss, and therefore we also applied the proposed calibration procedure to this method. The settings of KSD-Bayes are described in Appendix F.2.2. A Metropolis–Hastings algorithm was also used to sample from all generalised posteriors, the details for which are contained in Appendix F.2.3.

Results are presented for three different datasets of size $n = 1,000$ and dimension $d = 36$ ($m = 6$), $d = 64$ ($m = 8$), and $d = 100$ ($m = 10$) in Figure 4. For the lowest dimension $d = 36$, all the approaches produced similar posteriors. For the higher-dimensional cases, it can be seen that the DFD-Bayes and Pseudo-Bayes posteriors concentrate around the true parameter $\theta = 5$. The KSD-Bayes posterior is more conservative, whilst DFD-Bayes gives a comparable result to Pseudo-Bayes. For $d = 100$, the total computational time required to perform this analysis (including calibration) was 540 seconds for DFD-Bayes, $2,353$ seconds for KSD-Bayes, and $1,053$ seconds for Pseudo-Bayes each in average over 10 independent experiments, confirming that DFD-Bayes incurs a significantly lower computational cost than both alternatives. The value of the weight obtained through our calibration method for $d = 100$ in Figure 4 was 0.013 for DFD-Bayes, 0.157 for KSD-Bayes, and 0.579 for Pseudo-Bayes. These small values of weight indicated that the calibration

14

worked effectively, preventing the over-concentration of each posterior.

## 4.3   Multivariate Count Data

Finally we consider a problem involving multivariate count data. Count data occur in diverse application areas, and variables in such data are rarely independent, yet the literature on statistical modelling of such data is limited. Poisson graphical models and their extensions have emerged as a powerful tool for modelling such data; see the recent review of Inouye et al. [2017]. To the best of our knowledge a complete Bayesian treatment of Poisson graphical models has yet to be attempted[6], and we speculate that this is due to the computational challenges of the associated intractable likelihood. Our aim here is to assess the suitability of DFD-Bayes for learning the parameters of a Poisson graphical model.

Let $G$ be an undirected graph on a vertex set $\{1, \ldots, d\}$ and let $\mathcal{M}_i$ denote the neighbours of node $i$ that are contained in the set $\{i+1, \ldots, d\}$. A Poisson graphical model has probability mass function

$$p_\theta(\boldsymbol{x}) \propto \exp\left( \sum_{i=1}^d \theta_i x_i - \sum_{i=1}^d \sum_{j \in \mathcal{M}_i} \theta_{i,j} x_i x_j - \sum_{i=1}^d \log(x_i!) \right)$$

where the parameters $\theta$ consist of both the linear coefficients $\theta_i \in (-\infty, \infty)$ and the interaction coefficients $\theta_{i,j} \in [0, \infty)$. Our aim is to reproduce an analysis of a breast cancer gene expression dataset described in Inouye et al. [2017], but in a generalised Bayesian framework. For this problem, $n = 878$, $d = 10$, and $p = 64$ which renders the computational cost of $O(n^2 d)$ at every MCMC step and of $O(p^2 n^2 d)$ at calibration associated with KSD-Bayes inefficient. Full details of the dataset are contained in Appendix F.3.1. Independent standard normal priors were employed for each $\theta_i$, and half-normal distributions with scale $(d(d-1)/2)^{-1}$ were employed for each $\theta_{i,j}$. A No-U-Turn Sampler was used to sample from the DFD-Bayes posterior, as described in Appendix F.3.2. The gradient of the discrete Fisher divergence is available whenever $p_\theta(\boldsymbol{x}) = q_\theta(\boldsymbol{x})/C(\theta)$ with $q_\theta(\boldsymbol{x})$ differentiable with respect to $\theta$ at any $\boldsymbol{x} \in \mathcal{X}$; see Appendix F.3.3. The total computational time required for this analysis, including calibration, was $1,896$ seconds. Results, in Figure 5, demonstrate that the Poisson graphical model is in fact a poor fit for these data, which exhibit under-dispersion relative to the standard Poisson model. However, in terms of identifying the best parameter values for this model, DFD-Bayes appears to have performed well.

As a possible improvement, and to further stress-test the DFD-Bayes method, we considered a generalisation of the Poisson graphical model that allows for over- and under-dispersion, analogous to Conway and Maxwell [1962]. This model takes the form

$$p_\theta(\boldsymbol{x}) \propto \exp\left( \sum_{i=1}^d \theta_i x_i - \sum_{i=1}^d \sum_{j \in \mathcal{M}_i} \theta_{i,j} x_i x_j - \sum_{i=1}^d \theta_{0,i} \log(x_i!) \right)$$

where the additional parameters $\theta_{0,i} \in [0, \infty)$ control the dispersion, with $\theta_{0,i} = 1$ recovering the standard Poisson marginal. This time, $p = 74$ as opposed to $p = 64$ for the Poisson-based model. For this Conway–Maxwell–Poisson graphical model, the same priors as the Poisson graphical model

---

[6]A pairwise Markov random field whose marginals are close to being Poisson was used in Roy and Dunson [2020], and a specific generalisation of the Conway-Maxwell-Poisson was used in Piancastelli et al. [2021].
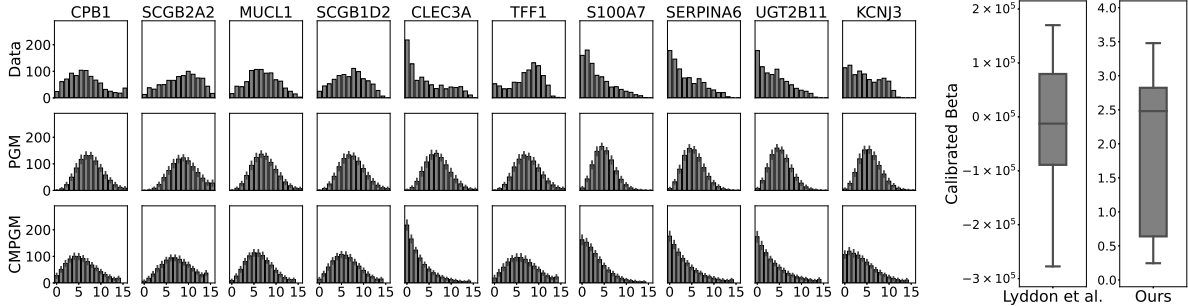
Figure 5: Left: Posterior predictive distributions from the Poisson graphical model and the Conway–Maxwell–Poisson graphical model. Right: Sampling distributions of $\beta_*$ for the Conway–Maxwell–Poisson graphical model by Lyddon et al. [2019] and by the proposed approach, computed using 10 independent realisations of the dataset.

were used for $\theta_i$ and $\theta_{i,j}$, and half-normal priors with scale $1/\sqrt{2}$ were used for each $\theta_{0,i}$. Results in Figure 5 demonstrate an improved fit to the dataset. Indeed, the optimal $\beta$ for the Poisson graphical model was $\beta_* = 0.2150$, which is smaller than the corresponding value $\beta_* = 0.9971$ for the Conway–Maxwell–Poisson graphical model, resulting in a conservative inference outcome when the statistical model is most misspecified and supporting the effectiveness of the proposed approach to calibration.

The right panel of Figure 5 shows the sampling distributions of estimators for the weight $\beta$ in the context of the Conway–Maxwell–Poisson graphical model, computed using bootstrap resampling of the gene expression dataset. It can be seen that the asymptotic approach proposed in Lyddon et al. [2019] is severely numerically unstable and can even lead to a negative weight, while the approach proposed in Section 3.3 remains stable within a reasonable range between 0 and 3.5. The lack of stability of the approach by Lyddon et al. [2019] arises from the need to invert a covariance matrix of derivatives of the loss, which can become numerically singular if the parameter dimension is high. In contrast, our approach involves no matrix inversion. This real-data analysis using flexible parametric models highlights the value in being able to perform rapid and automatic (i.e. free from user-specified degrees of freedom) generalised Bayesian inference for discrete intractable likelihood.

## 5   Conclusion

This paper proposed a novel generalised Bayesian inference procedure for discrete intractable likelihood. The approach, called DFD-Bayes, is distinguished by its lack of user-specified hyperparameters, its suitability for standard Markov chain Monte Carlo algorithms, and its linear (in $n$, the size of the dataset) computational cost per-iteration of the Markov chain. Furthermore, the generalised posterior is consistent and asymptotically normal. This paper also established a novel approach to calibration of generalised Bayesian posteriors which is computationally efficient (through embarrassing parallelism) and numerically stable, even when the parameter of the statistical model is high-dimensional.

This work focused on independent and identically distributed data, meaning that (for example) regression models were not considered. Relaxing the independence and identical distribution assumptions represents a natural direction for future work, and a road map is provided by recent

research in the score-matching literature [Xu et al., 2022].

One of our technical contributions is to present a discrete Fisher divergence applicable to distributions defined on multi-dimensional and countably infinite sets. This divergence can be regarded as a proper local scoring rule, which complements existing methodology developed in the finite domain context in Dawid et al. [2012]. The use of scoring rules as loss functions within a generalised Bayesian framework for continuous data was considered in Giummolè et al. [2019], Pacchiardi and Dutta [2021], and our work can be seen as an analogous approach for discrete data, with particular focus on intractable likelihood.

DFD-Bayes was demonstrated to outperform the comparative approach, KSD-Bayes, in our experiments both in terms of inferential performance and computational cost. However, one of the significant advantages of KSD-Bayes is robustness in the presence of outliers contained in dataset [Matsubara et al., 2022]. This is confirmed through an additional experiment on the Ising model in Appendix D.4 Thus, in settings where robust inference is required, the KSD-Bayes approach may be preferred. Future work could however focus on generalising our construction of the discrete Fisher divergence to allow for further robustness as per the diffusion score-matching framework of Barp et al. [2019].

# References

S. Amari. *Information Geometry and Its Applications*. Springer, 2016.

A. Anastasiou, A. Barp, F.-X. Briol, B. Ebner, R. Gaunt, F. Ghaderinezhad, J. Gorham, A. Gretton, C. Ley, Q. Liu, L. Mackey, C. J. Oates, G. Reinert, and Y. Swan. Stein's method meets computational statistics: A review of some recent developments. *Statistical Science*, 38(1):120–139, 2023.

C. Andrieu and G. Roberts. The pseudo-marginal approach for efficient Monte Carlo computations. *The Annals of Statistics*, 37(2):697–725, 2009.

A. Barp, F.-X. Briol, A. Duncan, M. Girolami, and L. Mackey. Minimum Stein discrepancy estimators. *Advances in Neural Information Processing Systems*, 32, 2019.

A. Benson and N. Friel. Bayesian inference, model selection and likelihood estimation using fast rejection sampling: the Conway-Maxwell-Poisson distribution. *Bayesian Analysis*, 16(3):905–931, 2021.

A. Berlinet and C. Thomas-Agnan. *Reproducing kernel Hilbert spaces in probability and statistics*. Springer, 2011.

J. Besag. Spatial interaction and the statistical analysis of lattice systems. *Journal of the Royal Statistical Society. Series B (Methodological)*, 36(2):192–236, 1974.

A. Bhattacharyya and Y. Atchade. A Bayesian analysis of large discrete graphical models. *arXiv:1907.01170*, 2019.

P. Bissiri, C. Holmes, and S. Walker. A general framework for updating belief distributions. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 78(5):1103, 2016.

B.-E. Cherief-Abdellatif and P. Alquier. MMD-Bayes: Robust Bayesian estimation via maximum mean discrepancy. *Proceedings of the 2nd Symposium on Advances in Approximate Bayesian Inference*, pages 1–21, 2020.

F. Chung and F. Graham. *Spectral graph theory*. American Mathematical Soc., 1997.

R. Conway and W. Maxwell. A queuing model with state dependent service rates. *Journal of Industrial Engineering*, 12:132–136, 1962.

J. Davidson. *Stochastic Limit Theory: An Introduction for Econometricians*. Oxford University Press, 1994.

A. P. Dawid and M. Musio. Bayesian model selection based on proper scoring rules. *Bayesian Analysis*, 10(2):479–499, 2015.

A. P. Dawid, S. Lauritzen, and M. Parry. Proper local scoring rules on discrete sample spaces. *The Annals of Statistics*, 40(1):593–608, 2012.

C. Dellaporta, J. Knoblauch, T. Damoulas, and F.-X. Briol. Robust Bayesian inference for simulator-based models via the MMD posterior bootstrap. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, pages 943–970, 2022.

R. Durrett. *Probability: Theory and Examples (4th Edition)*. Cambridge University Press, 2010.

E. M. Elçi, J. Grimm, L. Ding, A. Nasrawi, T. M. Garoni, and Y. Deng. Lifted worm algorithm for the Ising model. *Physical Review E*, 97(4):042126, 2018.

A. Ghosh and A. Basu. Robust Bayes estimation using the density power divergence. *Annals of the Institute of Statistical Mathematics*, 68:413–437, 2016.

F. Giummolè, V. Mameli, E. Ruli, and L. Ventura. Objective Bayesian inference with proper scoring rules. *Test*, 28(3):728–755, 2019.

G. Gong and F. J. Samaniego. Pseudo Maximum Likelihood Estimation: Theory and Applications. *The Annals of Statistics*, 9(4):861 – 869, 1981.

G. Hooker and A. Vidyashankar. Bayesian model robustness via disparities. *Test*, 23(3):556–584, 2014.

A. Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(24):695–709, 2005.

D. I. Inouye, E. Yang, G. I. Allen, and P. Ravikumar. A review of multivariate distributions for count data derived from the Poisson distribution. *Wiley Interdisciplinary Reviews: Computational Statistics*, 9(3):e1398, 2017.

J. Knoblauch, J. Jewson, and T. Damoulas. An optimization-centric view on bayes' rule: Reviewing and generalizing variational inference. *Journal of Machine Learning Research*, 23(132):1–109, 2022.

B. G. Lindsay, G. Y. Yi, and J. Sun. Issues and strategies in the selection of composite likelihoods. *Statistica Sinica*, pages 71–105, 2011.

D. Lusher, J. Koskinen, and G. Robins. *Exponential random graph models for social networks: Theory, methods, and applications.* Cambridge University Press, 2013.

S. P. Lyddon, C. C. Holmes, and S. G. Walker. General Bayesian updating and the loss-likelihood bootstrap. *Biometrika*, 106(2):465–478, 2019.

A.-M. Lyne, M. Girolami, Y. Atchadé, H. Strathmann, and D. Simpson. On Russian roulette estimates for Bayesian inference with doubly-intractable likelihoods. *Statistical Science*, 30(4): 443–467, 2015.

S. Lyu. Interpretation and generalization of score matching. *Proceedings of the 25th Conference on Uncertainty in Artificial Intelligence*, page 359–366, 2009.

J.-M. Marin, P. Pudlo, C. P. Robert, and R. Ryder. Approximate Bayesian computational methods. *Statistics and Computing*, 22:1167–1180, 2012.

T. Matsubara, J. Knoblauch, F.-X. Briol, and C. J. Oates. Robust generalised Bayesian inference for intractable likelihoods. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 84(3):997–1022, 2022.

J. W. Miller. Asymptotic normality, concentration, and coverage of generalized posteriors. *Journal of Machine Learning Research*, 22(168):1–53, 2021.

M. Moores, G. Nicholls, A. Pettitt, and K. Mengersen. Scalable Bayesian inference for the inverse temperature of a hidden Potts model. *Bayesian Analysis*, 15(1):1–27, 2020.

I. Murray, Z. Ghahramani, and D. MacKay. MCMC for doubly-intractable distributions. *Proceedings of the 22nd Conference on Uncertainty in Artificial Intelligence*, pages 359–366, 2006.

J. Møller, A. Pettitt, R. Reeves, and K. Berthelsen. An efficient Markov chain Monte Carlo method for distributions with intractable normalising constants. *Biometrika*, 93(2):451–458, 2006.

L. Pacchiardi and R. Dutta. Generalized Bayesian likelihood-free inference using scoring rules estimators. *arXiv:2104.03889*, 2021.

J. Park and M. Haran. Bayesian inference in the presence of intractable normalizing functions. *Journal of the American Statistical Association*, 113(523):1372–1390, 2018.

L. S. C. Piancastelli, N. Friel, W. Barreto-Souza, and H. Ombao. Multivariate Conway-Maxwell-Poisson distribution: Sarmanov method and doubly-intractable Bayesian inference. *arXiv:2107.07561*, 2021.

L. Price, C. Drovandi, A. Lee, and D. Nott. Bayesian synthetic likelihood. *Journal of Computational and Graphical Statistics*, 27(1):1–11, 2018.

J. Propp and D. Wilson. Coupling from the past: a user's guide. *Microsurveys in Discrete Probability*, 41:181–192, 1998.

A. Roy and D. B. Dunson. Nonparametric graphical model for counts. *Journal of Machine Learning Research*, 21(229):1–21, 2020.

W. Rudin. *Real and Complex Analysis, 3rd Ed.* McGraw-Hill, Inc., 1987.

J. Shi, Y. Zhou, J. Hwang, M. K. Titsias, and L. Mackey. Gradient estimation with discrete Stein operators. *arXiv: 2202.09497*, 2022.

G. Shmueli, T. P. Minka, J. B. Kadane, S. Borle, and P. Boatwright. A useful distribution for fitting discrete data: revival of the Conway–Maxwell–Poisson distribution. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 54(1):127–142, 2005.

N. Syring and R. Martin. Calibrating general posterior credible regions. *Biometrika*, 106(2):479–486, 2019.

A. van der Vaart. *Asymptotic Statistics*. Cambridge University Press, 1998.

C. Varin, N. Reid, and D. Firth. An overview of composite likelihood methods. *Statistica Sinica*, pages 5–42, 2011.

Y.-W. Wan, G. I. Allen, and Z. Liu. TCGA2STAT: simple TCGA data access for integrated statistical analysis in R. *Bioinformatics*, 32(6):952–954, 11 2015.

L. K. Wenliang and H. Kanagawa. Blindness of score-based methods to isolated components and mixing proportions. In *NeurIPS Workshop "Your Model is Wrong: Robustness and misspecification in probabilistic modeling"*, 2021.

J. Xu, J. Scealy, A. Wood, and T. Zou. Generalized score matching for regression. *arXiv:2203.09864*, 2022.

L. Xue, H. Zou, and T. Cai. Nonconcave penalized composite conditional likelihood estimation of sparse ising models. *The Annals of Statistics*, 40(3):1403–1429, 2012.

J. Yang, Q. Liu, V. Rao, and J. Neville. Goodness-of-fit testing for discrete distributions via Stein discrepancy. *Proceedings of the 35th International Conference on Machine Learning*, pages 5561–5570, 2018.

M. Zhang, O. Key, P. Hayes, D. Barber, B. Paige, and F.-X. Briol. Towards healing the blindness of score matching. In *NeurIPS 2022 Workshop on Score-Based Methods*, 2022.

# SUPPLEMENTARY MATERIAL

This supplementary material is structured as follows: Illustrative analysis of the discrete Fisher divergence and the DFD-Bayes using simple tractable models is presented in Appendix A. The proofs for all theoretical results are contained in Appendix B, with the proof of an auxiliary result reserved for Appendix C. The relationship between discrete Fisher divergence and Stein discrepancies is explored in Appendix D. Detailed calculations for worked examples are provided in Appendix E. Full details on our numerical experiments are provided in Appendix F

## A  Illustrative Analysis with Tractable Models

This section provides illustrative analysis of DFD-Bayes, including comparison with standard Bayesian inference and KSD-Bayes, using simple tractable models. We first demonstrate the calculation of the discrete Fisher divergence using the Bernoulli model. We then compare the properties of DFD-Bayes with standard Bayesian inference and KSD-Bayes, using the same Bernoulli model. We next discuss the influence of model misspecification on each posterior using the Poisson model. Finally, we provide an empirical illustration of the limitations of the discrete Fisher divergence discussed in Section 3.4. The Bernoulli and Poisson models are used for illustration and comparison in this section, since they are tractable and enable standard Bayesian inference to be performed.

### A.1  The Discrete Fisher Divergence for the Bernoulli Model

For $x \in \{0,1\}$, the Bernoulli model can be expressed by

$$p_\theta(x) = \theta^x (1-\theta)^{1-x} \tag{12}$$

where $\theta$ is the probability of $x = 1$. Recall that $p_\theta(1^+) = p_\theta(0)$ and $p_\theta(0^-) = p_\theta(1)$ under our increment/decrement rule. Both the increment and decrement of $p_\theta(1)$ are simply equal to $p_\theta(0)$, and likewise both the increment and decrement of $p_\theta(0)$ are equal to $p_\theta(1)$. Hence, they can be expressed by

$$p_\theta(x^+) = p_\theta(x^-) = \theta^{1-x}(1-\theta)^x, \tag{13}$$

that is $p_\theta(x^+) = \theta$ if $x = 0$ and $p_\theta(x^+) = 1 - \theta$ if $x = 1$. Plugging these into equation (5) in the manuscript with $d = 1$ gives an explicit form of the discrete Fisher divergence:

$$\mathrm{DFD}(p\|q) \stackrel{\theta}{=} \frac{1}{n} \sum_{i=1}^{n} \left( \frac{\theta^{1-x_i}(1-\theta)^{x_i}}{\theta^{x_i}(1-\theta)^{1-x_i}} \right)^2 - 2 \left( \frac{\theta^{x_i}(1-\theta)^{1-x_i}}{\theta^{1-x_i}(1-\theta)^{x_i}} \right) \tag{14}$$

Figure 6 shows the discrete Fisher divergence in (14) computed in three cases where 500 random samples are generated from the Bernoulli model with $\theta = 0.1$, $\theta = 0.5$ and $\theta = 0.9$, comparing the loss surface geometry with that of the negative log-likelihood. Both of the losses identify the parameter correctly in each case.

Although the geometrical shape of (14) is different from the negative log-likelihood, we can observe in Figure 6 that the discrete Fisher divergence is symmetric under the relabelling $y_i = 1 - x_i$
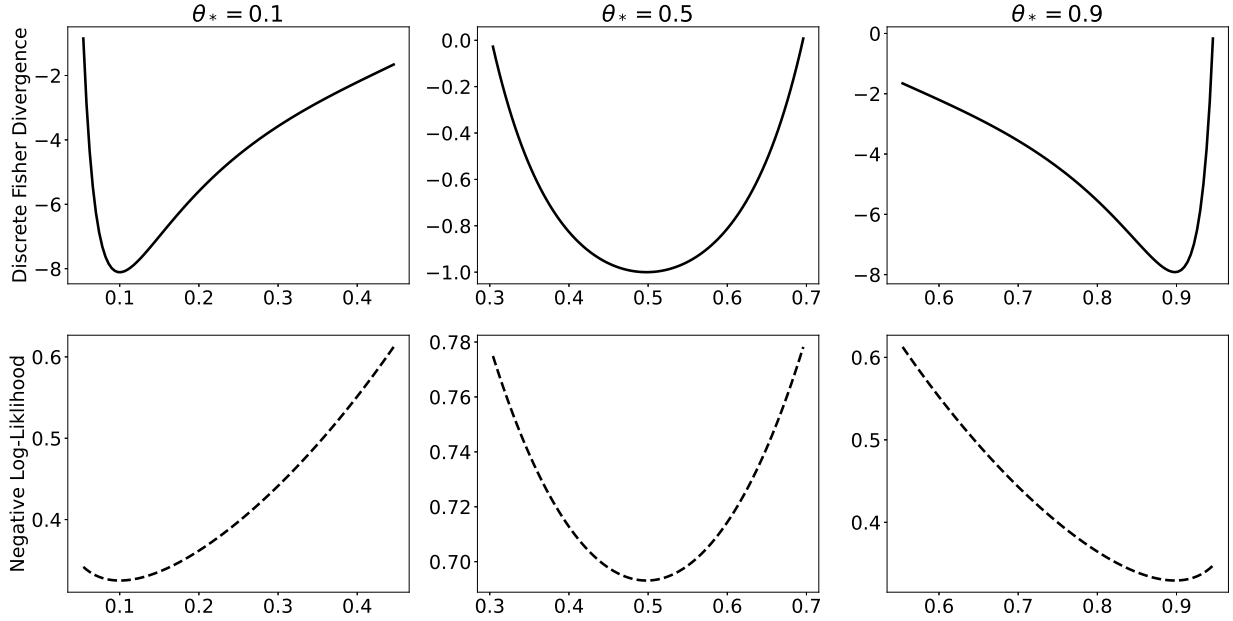
Figure 6: The discrete Fisher divergence (top, solid) and the negative log-likelihood (bottom, dash) between the Bernoulli model and data generated from the Bernoulli model of three different parameters $\theta_* = 0.1$ (left), $\theta_* = 0.5$ (centre), and $\theta_* = 0.9$ (right). They both identify the correct parameter $\theta_*$ in each case albeit the different loss surface geometries.

similarly to the negative log-likelihood in this example. This can indeed be verified as follows. If all data are relabelled, the above formula corresponds to

$$\text{DFD}(p\|q) \overset{\theta}{=} \frac{1}{n} \sum_{i=1}^n \left( \frac{\theta^{y_i}(1-\theta)^{1-y_i}}{\theta^{1-y_i}(1-\theta)^{y_i}} \right)^2 - 2 \left( \frac{\theta^{1-y_i}(1-\theta)^{y_i}}{\theta^{y_i}(1-\theta)^{1-y_i}} \right). \tag{15}$$

With a transform of the parameter $\rho = 1 - \theta$ applied, it further corresponds to

$$\text{DFD}(p\|q) \overset{\theta}{=} \frac{1}{n} \sum_{i=1}^n \left( \frac{\rho^{1-y_i}(1-\rho)^{y_i}}{\rho^{y_i}(1-\rho)^{1-y_i}} \right)^2 - 2 \left( \frac{\rho^{y_i}(1-\rho)^{1-y_i}}{\rho^{1-y_i}(1-\rho)^{y_i}} \right). \tag{16}$$

It is clear from comparison of (14) and (16) here that the discrete Fisher divergence of $\theta$ based on the original data $x_i$ is equivalent to that of $\rho = 1 - \theta$ based on the relabelled data $y_i = 1 - x_i$.

## A.2 Illustrative Comparison of DFD-Bayes with standard Bayes and KSD-Bayes

First, we derive the negative log-likelihood and the kernel Stein discrepancy for the Bernoulli model. The negative log-likelihood is

$$\text{NLL}(p_\theta \| p_n) = -\frac{1}{n} \sum_{i=1}^{n} x_i \log(\theta) + (1 - x_i) \log(1 - \theta). \tag{17}$$

The kernel Stein discrepancy in the discrete context was considered in Yang et al. [2018]. Letting $\rho_-(\theta, x) := p_\theta(x^-)/p_\theta(x) = \theta^{1-2x}(1-\theta)^{-1+2x}$, the kernel Stein discrepancy given a kernel function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is derived as

$$\text{KSD}(p_\theta \| p) \overset{\theta}{=} \frac{1}{n^2} \sum_{i=1}^{n} \sum_{j=1}^{n} (1 - \rho_-(\theta, x_i)) \, k(x_i, x_j) \, (1 - \rho_-(\theta, x_j)) +$$

$$(1 - \rho_-(\theta, x_i)) \left( k(x_i, x_j) - k(x_i, x_j^-) \right) + \left( k(x_i, x_j) - k(x_i^-, x_j) \right) (1 - \rho_-(\theta, x_j)) \Big]. \tag{18}$$

The DFD-Bayes posterior, the standard posterior, and the KSD-Bayes posterior are recovered from generalised posterior (2) built upon losses (14), (17), and (18), where $\beta$ is set to 1 for the standard posterior.

Next, we provide an analytical comparison of the credible regions of each posterior. As discussed in Section 3, a generalised posterior produces a credible region that differs from that of a standard posterior even in the asymptotic regime. For illustration, we derive the asymptotic variance of each posterior for the Bernoulli model. The asymptotic distribution of each posterior (appropriately centred) follows a Gaussian distribution $\mathcal{N}(0, \sigma^2)$ whose variance $\sigma^2$ is the inverse loss-Hessian at the minimiser $\theta_*$. To simplify the derivation, we use the Hamming distance kernel $k(x, x') = \mathbb{1}_{x=x'}$, that is 1 when $x = x'$ and otherwise 0, for the kernel Stein discrepancy. Let $\rho_+(\theta, x) := p_\theta(x)/p_\theta(x^+) = \theta^{-1+2x}(1-\theta)^{1-2x}$. By routine calculation, the second derivatives of each loss in the limit $n \to \infty$ are

$$\frac{\partial^2}{\partial^2 \theta} \text{NLL}(p_\theta \| p) = \mathbb{E}_{X \sim p} \left[ \frac{X}{\theta^2} + \frac{1 - X}{(1 - \theta)^2} \right],$$

$$\frac{\partial^2}{\partial^2 \theta} \text{DFD}(p_\theta \| p) = \mathbb{E}_{X \sim p} \left[ 2\rho_-(\theta, X) \frac{\partial^2}{\partial^2 \theta} \rho_-(\theta, X) + 2 \left( \frac{\partial}{\partial \theta} \rho_-(\theta, X) \right)^2 - 2 \frac{\partial^2}{\partial^2 \theta} \rho_+(\theta, X) \right],$$

$$\frac{\partial^2}{\partial^2 \theta} \text{KSD}(p_\theta \| p) = \mathbb{E}_{X \sim p} \left[ 2\rho_-(\theta, X) \frac{\partial^2}{\partial^2 \theta} \rho_-(\theta, X) + 2 \left( \frac{\partial}{\partial \theta} \rho_-(\theta, X) \right)^2 - 2 \frac{\partial^2}{\partial^2 \theta} \rho_-(\theta, X) \right].$$

For the kernel Stein discrepancy, given that $k(x_1, x_2) - k(x_1, x_2^-)$ and $k(x_1, x_2) - k(x_1^-, x_2)$ are 1 when $x = x'$ and otherwise $-1$, we simplify the expression as

$$\text{KSD}(p_\theta \| p) \overset{\theta}{=} \mathbb{E}_{X_1, X_2 \sim p} \left[ (1 - \rho_-(\theta, X_1)) \, k(X_1, X_2) \, (1 - \rho_-(\theta, X_2)) \right]$$

$$\overset{\theta}{=} \mathbb{E}_{X \sim p} \left[ (1 - \rho_-(\theta, X))^2 \right] \overset{\theta}{=} \mathbb{E}_{X \sim p} \left[ (\rho_-(\theta, X))^2 - 2\rho_-(\theta, X) \right].$$

Suppose that the population loss minimiser is $\theta_* = 0.5$, meaning that the data-generating distribution $p$ is the Bernoulli model with $\theta_* = 0.5$. We then have $\rho_-(\theta_*, x) = 1$, $\frac{\partial}{\partial \theta} \rho_-(\theta_*, x) = 2^2(1 - 2x)$,

$\frac{\partial^2}{\partial^2\theta}\rho_-(\theta_*, x) = 2^4(1 - 2x)^2$, and $\frac{\partial^2}{\partial^2\theta}\rho_+(\theta_*, x) = -2^4(1 - 2x)^2$. These gives us that

$$(\partial^2/\partial^2\theta)\text{NLL}(p_\theta\|p)|_{\theta=\theta_*} = \mathbb{E}_{X\sim p}\left[2^2 \times (X + 1 - X)\right] = 4,$$
$$(\partial^2/\partial^2\theta)\text{DFD}(p_\theta\|p)|_{\theta=\theta_*} = \mathbb{E}_{X\sim p}\left[3 \times 2^5 \times (1 - 2X)^4\right] = 96,$$
$$(\partial^2/\partial^2\theta)\text{KSD}(p_\theta\|p)|_{\theta=\theta_*} = \mathbb{E}_{X\sim p}\left[2 \times 2^4 \times (1 - 2X)^2\right] = 32.$$

By taking the inverse, the asymptotic variance $\sigma^2$ for the standard Bayes, the DFD-Bayes, and the KSD-Bayes is each given by 1/4, 1/96, and 1/32. In this example, the above calculation suggests that the DFD-Bayes has the narrowest credible region. The difference in these values emphasises the importance of calibrating $\beta$, which we do for all of our experiments in the manuscript.

Finally, we empirically demonstrate the difference between the posteriors and the influence of $\beta$. We computed each posterior in cases where (i) $\beta$ is *not* calibrated i.e. $\beta = 1$ and (ii) $\beta$ *is* calibrated (except for the standard posterior, which has $\beta = 1$). A Metropolis–Hastings algorithm was adopted to sample from all the posteriors. A Gaussian random walk proposal with covariance $\sigma^2 = 0.01$ was used. In total, 100 samples were obtained from each posterior by thinning 2,000 samples, after an initial burn-in of length 2,000. Figure 7 shows each posterior computed without $\beta$ calibrated. It confirms that, without calibration of $\beta$, the DFD-Bayes posterior has the narrowest credible region, which agrees with the analytical illustration provided above. Figure 8 shows each posterior computed with $\beta$ calibrated, where the result for the standard posterior is identical to Figure 7 as $\beta = 1$. For the DFD-Bayes and the KSD-Bayes, calibration of $\beta$ was performed by our proposal in Section 3.3, where we used 100 bootstrap minimisers to compute the analytical solution of $\beta_*$ in (10). It demonstrates that calibration of $\beta$ prevents over-concentration of the DFD-Bayes and the KSD-Bayes.
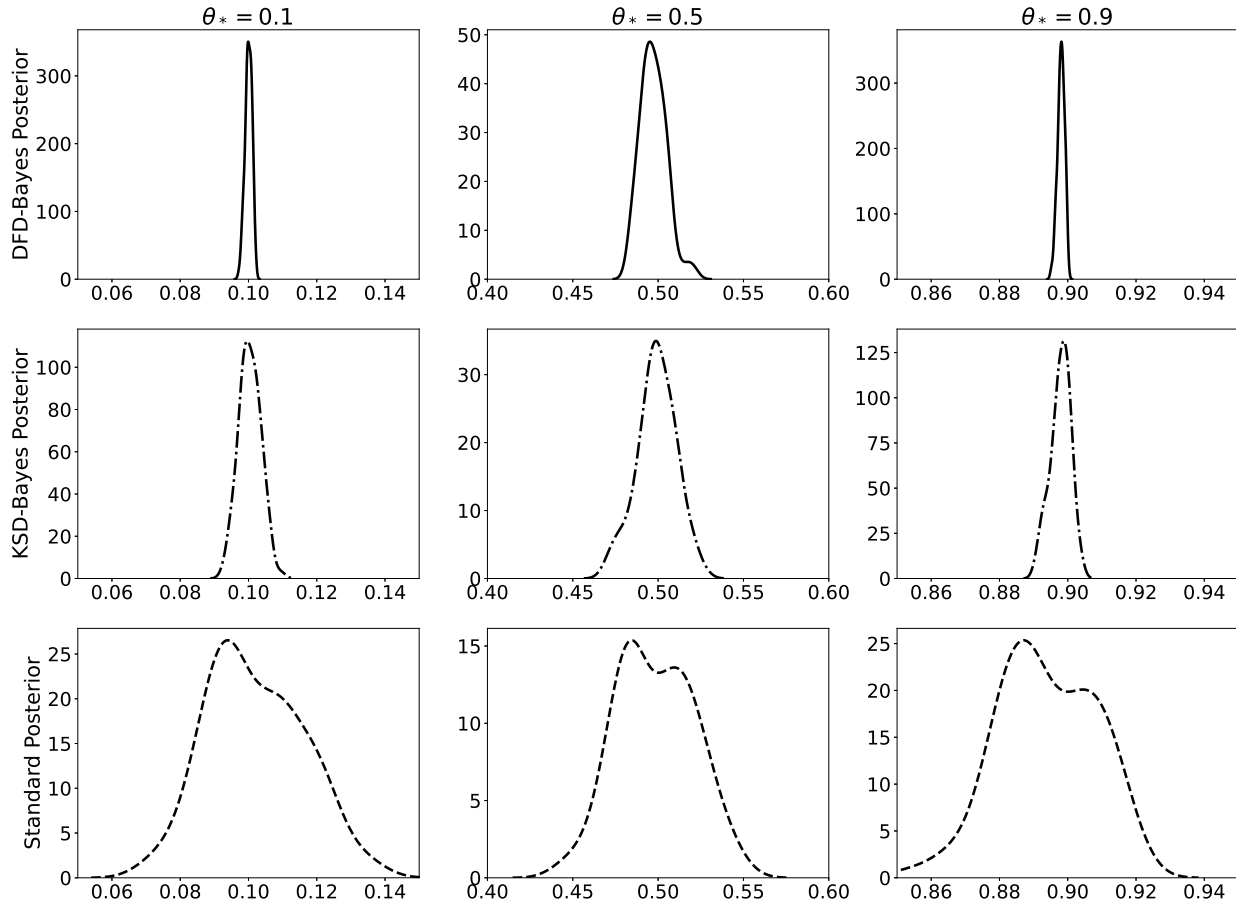
Figure 7: The DFD-Bayes posterior (top, solid), the KSD-Bayes posterior (middle, dash-dot), and the standard posterior (bottom, dash) computed without $\beta$ calibrated, for data generated from the Bernoulli model with three different parameters $\theta_* = 0.1$ (left), $\theta_* = 0.5$ (centre), and $\theta_* = 0.9$ (right). While their scales and geometries are different, all methods identify the correct parameter $\theta_*$.

Figure 8: The DFD-Bayes posterior (top, solid), the KSD-Bayes posterior (middle, dash-dot), and the standard posterior (bottom, dash) computed with $\beta$ calibrated, for data generated from the Bernoulli model with three different parameters $\theta_* = 0.1$ (left), $\theta_* = 0.5$ (centre), and $\theta_* = 0.9$ (right). While their scales and geometries are different, all methods identify the correct parameter $\theta_*$.
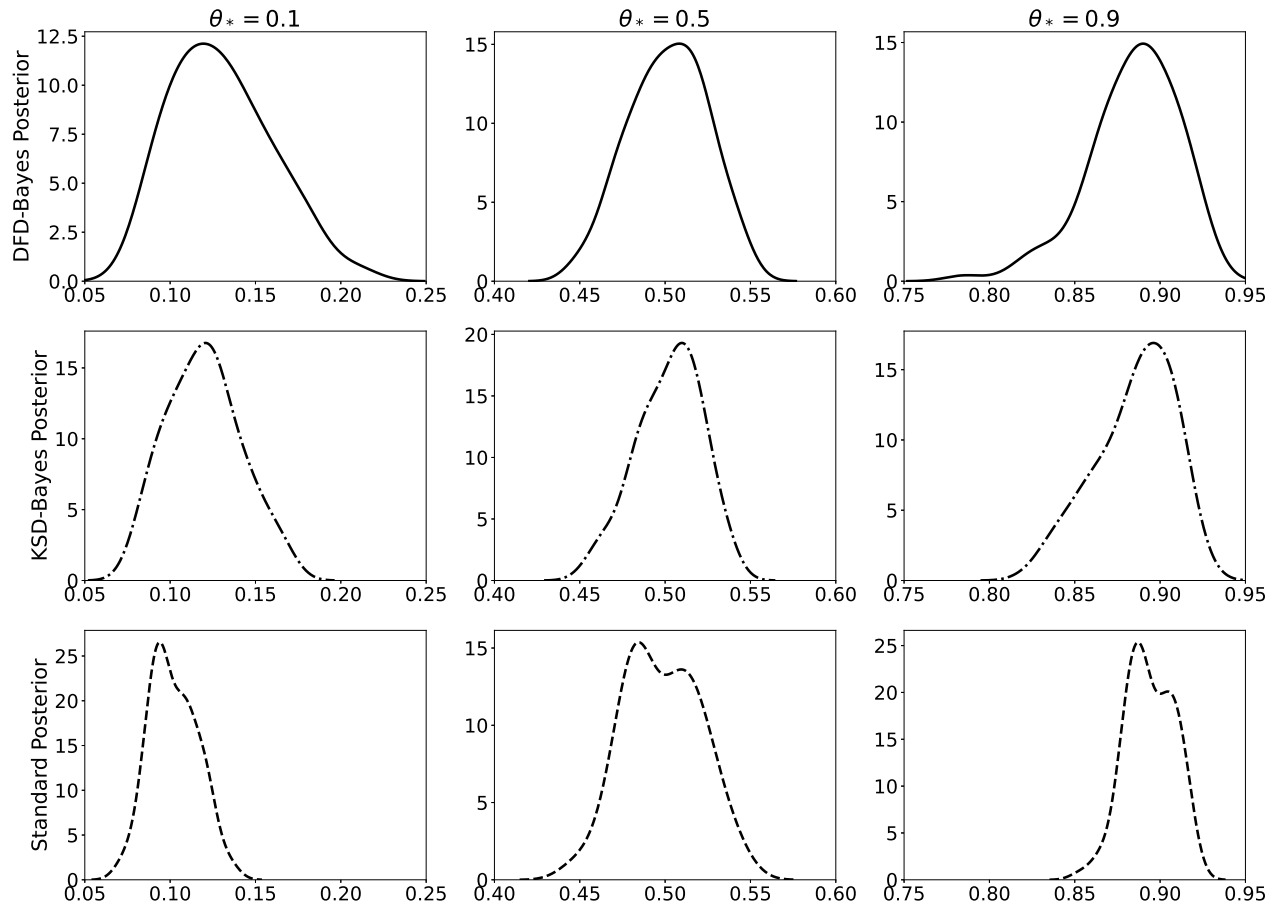
## A.3 Influence of Model Misspecification

Next we turn our attention to the influence of model misspecification on each method. It is convenient to consider the Poisson model to introduce a synthetic model misspecification. For $x \in \mathbb{N}_0$, the Poisson model is

$$p_\theta(x) = \frac{\theta^x \exp(\theta)}{x!}. \tag{19}$$

Then, the negative log-likelihood and the discrete Fisher divergence are

$$\mathrm{NLL}(p_\theta \| p_n) \stackrel{\theta}{=} \frac{1}{n} \sum_{i=1}^n -x_i \log(\theta) + \theta, \tag{20}$$

$$\mathrm{DFD}(p_\theta \| p_n) \stackrel{\theta}{=} \frac{1}{n} \sum_{i=1}^n \left( \frac{x_i}{\theta} \right)^2 - 2\frac{x_i + 1}{\theta}. \tag{21}$$

Letting $\rho_-(\theta, x) := p_\theta(x^-)/p_\theta(x) = x_i/\theta$, the kernel Stein discrepancy is

$$\mathrm{KSD}(p_\theta \| p) \stackrel{\theta}{=} \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \left(1 - \rho_-(\theta, x_i)\right) k(x_i, x_j) \left(1 - \rho_-(\theta, x_j)\right) +$$

$$\left(1 - \rho_-(\theta, x_i)\right) \left( k(x_i, x_j) - k(x_i, x_j^-) \right) + \left( k(x_i, x_j) - k(x_i^-, x_j) \right) \left(1 - \rho_-(\theta, x_j)\right). \tag{22}$$

For the kernel Stein discrepancy, we use a similar choice of kernel to Matsubara et al. [2022], that induces a robustness suitable for this example: $k(x, x') = m(x) \exp(-\mathbb{1}_{x=x'}) m(x')$ where $m(x) = \sigma(15 - x)$ based on a sigmoid function $\sigma(t) = (1 + \exp(-t))^{-1}$.

For illustration, we synthetically introduce model misspecification by mixing outliers into the data. We sampled 500 data points $\{x_i\}_{i=1}^n$ from the Poisson model with the parameter $\theta_* = 5$, and replaced the $100 \times \epsilon$ percent of data with an outlier $y = 20$ that is larger than the 99.9% percentile of the Poisson distribution of $\theta_* = 5$. This causes a synthetic model misspecification because the dataset is generated from a mixture of the Possion model and the Dirac distribution at $y = 20$, which cannot be adequately explained by only the Poisson model. The sensitivity of each posterior to the outlier can be analytically investigated. The standard Bayesian posterior is modestly impacted by the outlier $y$, given that the negative log-likelihood (20) is a linear function of each datum $x_i$. On the other hand, in this example, DFD-Bayes may be more severely impacted, given the discrete Fisher divergence (21) is a quadratic function of each datum $x_i$. The growth rate of the kernel Stein discrepancy with respect to each datum $x_i$ is determined by the choice of kernel $k$. We compute each posterior for two cases when $\epsilon = 0.0$ (no outlier contained) and $\epsilon = 0.1$ (10% outliers contained), to empirically demonstrate the impact of the model misspecification. The Metropolis–Hastings algorithm with the Gaussian random walk proposal of $\sigma^2 = 0.1$ is used to sample from each posterior with calibration applied. In total, 100 samples were obtained from each posterior by thinning 2,000 samples, after an initial burn-in of length 2,000.

Figure 9 demonstrates the sensitivity of the standard Bayesian posterior and DFD-Bayes to the outliers, whlie KSD-Bayes shows insensitivity due to the careful choice of kernel. See also Appendix D.4 for more discussion on robustness of KSD-Bayes. In this example, the sensitivity of the DFD-Bayes to the outlier was higher than the standard Bayesian posterior, as anticipated. Barp et al. [2019] proposed a robust analogue of the Fisher divergence in the continuous case.
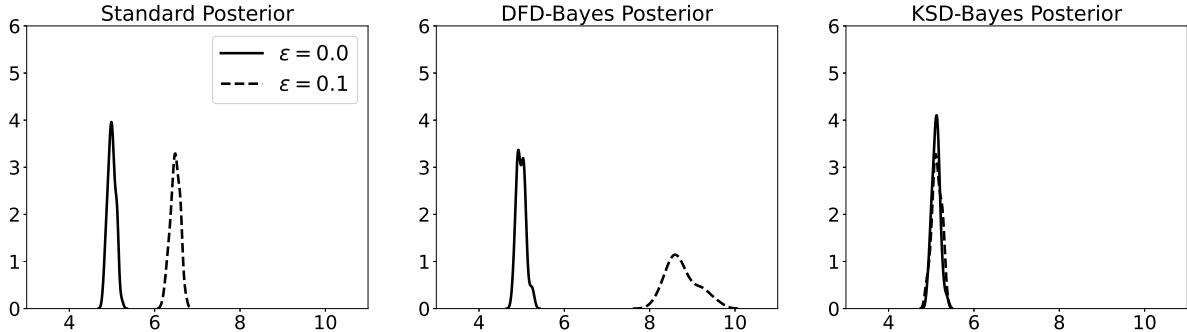
Figure 9: The standard posterior (left), The DFD-Bayes posterior (centre), and the KSD-Bayes posterior (right) computed with $\beta$ calibrated for data when $\epsilon = 0.0$ (solid line) and $\epsilon = 0.1$ (dash line), that is, the 10% of data is replaced with outlier $y$.

Although this is not a focus of this work, a similar approach may be applied to the discrete case when severe model misspecification is anticipated. This would be an interesting avenue for further work, but our present interest is in computation for discrete intractable likelihood.

## A.4 Limitation of DFD-Bayes for Inference of Mixture Parameters

Finally, we provide an empirical illustration of the limitation of score-based methods in Section 3.4. It has been pointed out that score-based methods generally exhibit insensitivity to mixing proportions when mixture components have isolated high-probability regions [Wenliang and Kanagawa, 2021, Zhang et al., 2022]. In the continuous case, this can be observed using a mixture model of two Gaussian distributions $\mathbb{P}_\theta(x) = (1 - \theta) \times \mathcal{N}(-\mu, 1) + \theta \times \mathcal{N}(\mu, 1)$ whose parameter is the mixture ratio. Zhang et al. [2022] illustrated how the Fisher divergence is approximately constant over $\Theta$ if $\mu$ is large enough to isolate the components $\mathcal{N}(-\mu, 1)$ and $\mathcal{N}(\mu, 1)$. We illustrate the same limitation for the discrete Fisher divergence using a mixture model of two Poisson distributions $p_\theta(x) = (1 - \theta) \times q_{\lambda_1}(x) + \theta \times q_{\lambda_2}(x)$, where $q_{\lambda_1}$ and $q_{\lambda_2}$ are the Poisson distributions with rate parameters $\lambda_1 > 0$ and $\lambda_2 > 0$. Figure 10 shows the geometry of the discrete Fisher divergence between the mixture model $p_\theta$ and data generated from the mixture model $p_{\theta_*}$ with the true mixture proportion $\theta_*$, for two cases when the supports of the two Poisson distributions are highly isolated and when they are not isolated. The correct mixture proportion $\theta_*$ was identified only in the latter case, while in the former case the discrete Fisher divergence was approximately constant. See Zhang et al. [2022] for a potential approach to remedy this general limitation of score-based methods.

## B Proofs of Theoretical Results

This section contains the proof of all theoretical results in the paper, including Proposition 1, Theorem 1 and Theorem 2.
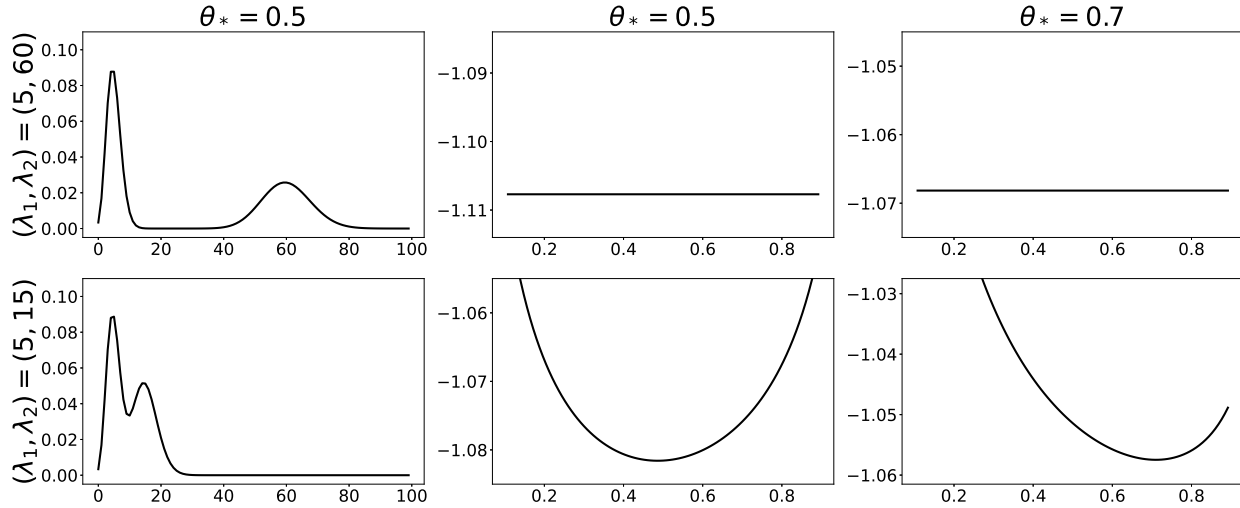
Figure 10: The form of the Poisson mixture model $p_{\theta_*}$ when $\theta_* = 0.5$ (left), the discrete Fisher divergence computed for data generated from the model $p_{\theta_*}$ with $\theta_* = 0.5$ (middle), and the discrete Fisher divergence computed for data generated from the model $p_{\theta_*}$ with $\theta_* = 0.7$ (right), for two cases where $\lambda_1 = 5, \lambda_2 = 60$ (top) and $\lambda_1 = 5, \lambda_2 = 15$ (bottom).

## B.1 Proof of Proposition 1

First we introduce three technical lemmas that will be useful:

**Lemma 1.** *For any $\boldsymbol{x} \in \mathcal{X}$ and $i = 1, \ldots, d$, it holds that $(\boldsymbol{x}^{i-})^{i+} = \boldsymbol{x}$ and $(\boldsymbol{x}^{i+})^{i-} = \boldsymbol{x}$.*

*Proof.* Since $\mathcal{X} = S_1 \times \cdots \times S_d$ from the Standing Assumption,

$$\boldsymbol{x}^{i-} = (x_1, \cdots, x_i^-, \cdots, x_d), \qquad \boldsymbol{x}^{i+} = (x_1, \cdots, x_i^+, \cdots, x_d). \tag{23}$$

It is thus sufficient to show that $(x_i^-)^+ = x_i$ and $(x_i^+)^- = x_i$ for any $i = 1, \ldots, d$. Consider, therefore, a set $S \cong I \subseteq \mathbb{Z}$ with more than one element. Our aim is to establish the identity $(s^-)^+ = s$ and $(s^+)^- = s$ for all $s \in S$. Existence of the least and greatest element $s_{\min}$ and $s_{\max}$ of $S$ determines four qualitatively distinct cases to be checked: (i) neither of them exist; (ii) both of them exists; (iii) only $s_{\min}$ exists; (iv) only $s_{\max}$ exists. Recall that we identify the case (iv) with (iii) without loss of generality by reversing the ordering of $S$. The identity for (i) & (ii) is trivial since the maps $s \mapsto s^-$ is bijective from $S$ to itself with inverse $s \mapsto s^+$. For case (iii), we have $(s^-)^+ = s$ for $s \neq s_{\min}$ and $(s^+)^- = s$ for all $s \in S$, Recalling the definition $s_{\min}^- = \star$ and $\star^+ = s_{\min}$ completes the argument. $\square$

**Lemma 2.** *For any $f, g : \mathcal{X} \to \mathbb{R}$ and any $i = 1, \ldots, d$, suppose $\sum_{\boldsymbol{x} \in \mathcal{X}} |f(\boldsymbol{x}) g(\boldsymbol{x}^{i-})| < \infty$, that is, the series is absolutely convergent. Then we have*

$$\sum_{\boldsymbol{x} \in \mathcal{X}} f(\boldsymbol{x}) g(\boldsymbol{x}^{i-}) = \sum_{\boldsymbol{x} \in \mathcal{X}} f(\boldsymbol{x}^{i+}) g(\boldsymbol{x}). \tag{24}$$

29

*Proof.* Since $\mathcal{X} = S_1 \times \cdots \times S_d$ from the Standing Assumption, the series can be expressed as

$$\sum_{\boldsymbol{x} \in \mathcal{X}} f(\boldsymbol{x})g(\boldsymbol{x}^{i-}) = \sum_{x_1 \in S_1} \cdots \sum_{x_i \in S_i} \cdots \sum_{x_d \in S_d} f(x_1, \cdots, x_i, \cdots, x_d)g(x_1, \cdots, x_i^-, \cdots, x_d),$$

$$\sum_{\boldsymbol{x} \in \mathcal{X}} f(\boldsymbol{x}^{i+})g(\boldsymbol{x}) = \sum_{x_1 \in S_1} \cdots \sum_{x_i \in S_i} \cdots \sum_{x_d \in S_d} f(x_1, \cdots, x_i^+, \cdots, x_d)g(x_1, \cdots, x_i, \cdots, x_d).$$

Holding the coordinates $x_1, \ldots, x_{i-1}, x_{i+1}, \ldots, x_d$ fixed, and exploiting absolute convergence to justify the interchange of summations, the claimed result follows if

$$\sum_{x_i \in S_i} \tilde{f}(x_i)\tilde{g}(x_i^-) = \sum_{x_i \in S_i} \tilde{f}(x_i^+)\tilde{g}(x_i) \tag{25}$$

where $\tilde{f}(x_i) := f(x_1, \cdots, x_i, \cdots, x_d)$ and $\tilde{g}(x_i) := g(x_1, \cdots, x_i, \cdots, x_d)$ are viewed as functions on $S_i$.

Consider, therefore, an arbitrary set $S \cong I \subseteq \mathbb{Z}$, for which we aim to establish the identity $\sum_{s \in S} h(s)k(s^-) = \sum_{s \in S} h(s^+)k(s)$ for any functions $h, k : S \to \mathbb{R}$ s.t. $\sum_{s \in S} |h(s)k(s^-)| < \infty$. From the definition of an order isomorphism, the elements of $S$ can be indexed as $S = \{s_i : i \in I\}$, where $s_i < s_j$ if and only if $i < j$. The identity therefore can be written as $\sum_{i \in I} h(s_i)k(s_i^-) = \sum_{i \in I} h(s_i^+)k(s_i)$, and will be verified for the three qualitatively distinct cases of index set $I$ described in the proof of Lemma 1:

(i) $I = \mathbb{Z}$. The result is immediate, since $(s_i, s_i^-) = (s_i, s_{i-1})$ and $(s_i^+, s_i) = (s_{i+1}, s_i)$ range over the same set for $i \in I$. The series $\sum_{i \in I} h(s_i^+)k(s_i)$ is absolutely convergent since the sets $\{h(s_i)k(s_i^-)\}_{i \in I}$ and $\{h(s_i^+)k(s_i)\}_{i \in I}$ in the two series are equal.

(ii) $I = \{1, \ldots, n\}$ for some $n \in \mathbb{N}$. In this case $s_{\min} = s_1$ and $s_{\max} = s_n$, and it follows from the definition of decrements and increments that

$$\sum_{i \in I} h(s_i)k(s_i^-) = h(s_1)k(s_1^-) + h(s_2)k(s_1) + \cdots + h(s_n)k(s_{n-1})$$

$$= h(s_n^+)k(s_n) + h(s_2)k(s_1) + \cdots + h(s_n)k(s_{n-1}) = \sum_{i \in I} h(s_i^+)k(s_i),$$

where the sets $\{h(s_i)k(s_i^-)\}_{i \in I}$ and $\{h(s_i^+)k(s_i)\}_{i \in I}$ are again equal.

(iii) $I = \{1, 2, \ldots\}$. In this case $s_{\min} = s_1$, and it follows from the definition $s_1^- = \star$ and $k(\star) = 0$ that

$$\sum_{i \in I} h(s_i)k(s_i^-) = \underbrace{h(s_1)k(\star)}_{=0} + h(s_2)k(s_1) + h(s_3)k(s_2) + \cdots$$

$$= h(s_2)k(s_1) + h(s_3)k(s_2) + \cdots = \sum_{i \in I} h(s_i^+)k(s_i).$$

The series $\sum_{i \in I} h(s_i^+)k(s_i)$ is absolutely convergent since the set $\{h(s_i^+)k(s_i)\}_{i \in I}$ is a subset of the absolutely summable set $\{h(s_i)k(s_i^-)\}_{i \in I}$.

This completes the proof. $\qquad\square$

Let $F(\mathcal{X}, S)$ denote the set of all functions $f$ of the form $f : \mathcal{X} \to S$.

**Lemma 3.** *For probability mass function* $p : \mathcal{X} \to (0, \infty)$, *the map* $\mu_p := (\nabla^- p)/p$ *is an injection* $\mu : F(\mathcal{X}, (0, \infty)) \to F(\mathcal{X}, \mathbb{R}^d)$.

*Proof.* It suffices to show that each value $p(\boldsymbol{x})$, for $\boldsymbol{x} \in \mathcal{X}$, can be explicitly recovered from $\mu_p$. Note that, since $p$ takes values in $(0, \infty)$, the embedding $\mu_p$ is well-defined. From the Standing Assumption, we have that $\mathcal{X} = S_1 \times \cdots \times S_d$, where each $S_i \cong I_i \subseteq \mathbb{Z}$ is a set with more than one element. Since the $S_i$ serve only as index sets, we can without loss of generality assume that $S_i$ is a consecutive subset of $\mathbb{Z}$ and that $0 \in S_i$, for each $i = 1, \ldots, d$. The idea of the proof is to demonstrate that each of the quantities $p(\boldsymbol{x})$ can be explicitly expressed in terms of $\mu_p$, $p(\mathbf{0})$ and $\{p(\boldsymbol{y}) : \|\boldsymbol{y}\|_1 < \|\boldsymbol{x}\|_1\}$, where $\|\boldsymbol{x}\|_1 := |x_1| + \cdots + |x_d|$. It would then follow from a simple inductive argument that $p(\boldsymbol{x})$ can be expressed in terms of $\mu_p$ and $p(\mathbf{0})$. Finally, the constraint that $\sum_{\boldsymbol{x} \in \mathcal{X}} p(\boldsymbol{x}) = 1$ uniquely determines $p(\mathbf{0})$, demonstrating that $p(\boldsymbol{x})$ can be explicitly recovered.

Given $\boldsymbol{x} \in \mathcal{X}$, assume $\boldsymbol{x} \neq \mathbf{0}$, for otherwise the claim will trivially hold. Then let $i \in \{1, \ldots, d\}$ be such that $x_i \neq 0$. If $x_i > 0$, then from the definition of $\mu_p(\boldsymbol{x})_i = 1 - p(\boldsymbol{x}^{i-})/p(\boldsymbol{x})$ we have the relation

$$p(\boldsymbol{x}) = \frac{p(\boldsymbol{x}^{i-})}{1 - \mu_p(\boldsymbol{x})_i}$$

where $\|\boldsymbol{x}^{i-}\|_1 = \|\boldsymbol{x}\|_1 - 1$. Conversely, if $x_i < 0$, then using Lemma 1 we have $\mu_p(\boldsymbol{x}^{i+})_i = 1 - p(\boldsymbol{x})/p(\boldsymbol{x}^{i+})$ and we have the relation

$$p(\boldsymbol{x}) = [1 - \mu_p(\boldsymbol{x}^{i+})_i]p(\boldsymbol{x}^{i+})$$

where $\|\boldsymbol{x}^{i+}\|_1 = \|\boldsymbol{x}\|_1 - 1$. The previously described inductive argument completes the proof. $\qquad\square$

Now we prove the main result:

*Proof of Proposition 1.* Expanding the square gives that

$$\mathrm{DFD}(p\|q) = \mathbb{E}_{X \sim q}\left[\sum_{i=1}^{d}\left(\frac{p(X) - p(X^{j-})}{p(X)}\right)^2\right.$$

$$\left. - 2\underbrace{\frac{p(X) - p(X^{j-})}{p(X)}\frac{q(X) - q(X^{j-})}{q(X)}}_{=:(*)} + \left(\frac{q(X) - q(X^{j-})}{q(X)}\right)^2\right].$$

Denote by $\mathrm{supp}(q)$ the support of $q$ i.e. $\{\boldsymbol{x} \in \mathcal{X} \mid q(\boldsymbol{x}) \neq 0\}$. For the term $\mathbb{E}_{X \sim q}[(*)]$, it follows from the definition $\mathbb{E}_{X \sim q}[f(X)] = \sum_{\boldsymbol{x} \in \mathrm{supp}(q)} f(\boldsymbol{x})q(\boldsymbol{x})$ that

$$\mathbb{E}_{X \sim q}[(*)] = \sum_{j=1}^{d} \mathbb{E}_{X \sim q}\left[\frac{p(X) - p(X^{j-})}{p(X)} - \frac{p(X) - p(X^{j-})}{p(X)}\frac{q(X^{j-})}{q(X)}\right]$$

$$= \sum_{j=1}^{d}\left\{\sum_{\boldsymbol{x} \in \mathrm{supp}(q)} \frac{p(\boldsymbol{x}) - p(\boldsymbol{x}^{j-})}{p(\boldsymbol{x})}q(\boldsymbol{x}) - \underbrace{\sum_{\boldsymbol{x} \in \mathrm{supp}(q)} \frac{p(\boldsymbol{x}) - p(\boldsymbol{x}^{j-})}{p(\boldsymbol{x})}q(\boldsymbol{x}^{j-})}_{(**)}\right\},$$

We apply Lemma 2 to the term $(**)$ with $f(\boldsymbol{x}) = (p(\boldsymbol{x}) - p(\boldsymbol{x}^{j-}))/p(\boldsymbol{x})$ and $g(\boldsymbol{x}) = q(\boldsymbol{x})$, where $f(\boldsymbol{x}^{j+})$ is well-defined for all $\boldsymbol{x} \in \mathrm{supp}(q)$ due to the assumption that $p(\boldsymbol{x}^{j+}) > 0$ for $\boldsymbol{x} \in \mathrm{supp}(q)$ and Lemma 2 is thus applicable. This reveals that

$$(**) = \sum_{\boldsymbol{x} \in \mathrm{supp}(q)} \frac{p(\boldsymbol{x}^{j+}) - p(\boldsymbol{x})}{p(\boldsymbol{x}^{j+})} q(\boldsymbol{x})$$

for each $j = 1, \ldots, d$, where Lemma 1 is used to deduce that $(\boldsymbol{x}^{j-})^{j+} = \boldsymbol{x}$. Hence, we have

$$\mathbb{E}_{X \sim q}[(*)] = \mathbb{E}_{X \sim q}\left[ \sum_{i=1}^{d} \frac{p(X) - p(X^{j-})}{p(X)} - \frac{p(X^{j+}) - p(X)}{p(X^{j+})} \right] = \mathbb{E}_{X \sim q}\left[ \sum_{i=1}^{d} -\frac{p(X^{j-})}{p(X)} + \frac{p(X)}{p(X^{j+})} \right].$$

Plugging this equality in the discrete Fisher divergence at the top and completing the expansion establish that

$$\mathrm{DFD}(p\|q) = \mathbb{E}_{X \sim q}\left[ \sum_{i=1}^{d} \left(1 - \frac{p(X^{j-})}{p(X)}\right)^2 + 2\frac{p(X^{j-})}{p(X)} - 2\frac{p(X)}{p(X^{j+})} + \left(1\frac{q(X^{j-})}{q(X)}\right)^2 \right]$$

$$= \mathbb{E}_{X \sim q}\left[ \sum_{i=1}^{d} \left(\frac{p(X^{j-})}{p(X)}\right)^2 - 2\frac{p(X)}{p(X^{j+})} \right] + \underbrace{\mathbb{E}_{X \sim q}\left[ \sum_{i=1}^{d} 1 + \left(1 - \frac{q(X^{j-})}{q(X)}\right)^2 \right]}_{=:C(q)}.$$

Finally we verify that $\mathrm{DFD}(p\|q) = 0$ if and only if $p = q$. From Lemma 3 we have the injective embedding $p \mapsto \mu_p := (\nabla^- p)/p$ of a positive density $p : \mathcal{X} \to (0, \infty)$ into $F(\mathcal{X}, \mathbb{R}^d)$. Since $q > 0$, the map $p \mapsto \mu_p$ is also an injection into $L^2(q, \mathbb{R}^d)$, equipped with the canonical norm $\|\nu\|_{L^2(q, \mathbb{R}^d)} := \mathbb{E}_{X \sim q}[\|\nu(X)\|^2]$, $\forall \nu \in L^2(q, \mathbb{R}^d)$. From (3) we recognise that $\mathrm{DFD}(p\|q) = \|\mu_p - \mu_q\|^2_{L^2(q, \mathbb{R}^d)}$ is the squared distance between $\mu_p$ and $\mu_q$ according to the canonical norm of $L^2(q, \mathbb{R}^d)$. Since $\|\mu_p - \mu_q\|_{L^2(q, \mathbb{R}^d)} = 0$ if and only if $\mu_p = \mu_q$ in $L^2(q, \mathbb{R}^d)$, it follows from injectivity of $p \mapsto \mu_p$ that $\mathrm{DFD}(p\|q) = 0$ if and only if $p = q$, as required. $\square$

## B.2 Proof of Theorem 1

This appendix contains the proof of Theorem 1. Miller [2021] provided sufficient conditions for consistency and asymptotic normality of generalised Bayesian posteriors of the form $\pi_n^D(\mathrm{d}\theta) \propto \exp(-n D_n(\theta)) \pi(\mathrm{d}\theta)$, where $D_n : \Theta \to \mathbb{R}$ is a loss function that may depend on the data $\{\boldsymbol{x}_i\}_{i=1}^n$. These results can be leveraged to analyse DFD-Bayes, by setting

$$D_n(\theta) \overset{\theta}{=} \frac{\beta}{n} \sum_{i=1}^{n} \sum_{i=1}^{d} \left(\frac{p_\theta(\boldsymbol{x}_i^{j-})}{p_\theta(\boldsymbol{x}_i)}\right)^2 - 2\left(\frac{p_\theta(\boldsymbol{x}_i)}{p_\theta(\boldsymbol{x}_i^{j+})}\right). \tag{26}$$

These conditions were refined into more applicable forms in Matsubara et al. [2022]. While Matsubara et al. [2022] focused on their particular case of losses based on kernelised Stein discrepancies, their argument can be directly applied for essentially any arbitrary loss $D_n$. We repeat this argument by modifying it so that it can be applied for any loss $D_n$. Let $B_\epsilon(\theta_*) := \{\theta \in \Theta : \|\theta - \theta_*\| < \epsilon\}$.

**Theorem 3.** *Let $\Theta \subseteq \mathbb{R}^v$ be Borel. Let $D : \Theta \to \mathbb{R}$ be a fixed measurable function and $\{D_n\}_{n=1}^\infty$ be a sequence s.t. $D_n : \Theta \to \mathbb{R}$ is a measurable function dependent on random data $\{\boldsymbol{x}_i\}_{i=1}^n \subset \mathcal{X}$. Let $H_n(\theta) := \nabla_\theta^2 D_n(\theta)$. Suppose that, for some bounded convex open set $U \subseteq \Theta$, the following hold:*

*C1* $D_n$ *a.s. converges pointwise to* $D$;

*C2* $D_n$ *is* $r$ *times continuously differentiable in* $U$ *and* $\limsup_{n\to\infty} \sup_{\theta\in U} \|\nabla_\theta^r D_n(\theta)\| < \infty$ *a.s. for* $r = 1, 2, 3$;

*C3 for all* $n$ *sufficiently large,* $\theta_n \in U$ *for any* $\theta_n \in \arg\min D_n$ *a.s., and a point* $\theta_* \in U$ *uniquely attains* $D(\theta_*) = \inf_{\theta\in\Theta} D(\theta)$.

*C4* $H_n(\theta_*) \overset{a.s.}{\to} H_*$ *for some nonsingular* $H_*$;

*C5* $\pi$ *is continuous and positive at* $\theta_*$.

*Then, for any* $\epsilon > 0$, *the generalised posterior* $\pi_n^D(d\theta) \propto \exp(-nD_n(\theta))\pi(d\theta)$ *satisfies*

$$\int_{B_\epsilon(\theta_*)} \pi_n^D(\theta)\, d\theta \xrightarrow{\text{a.s.}} 1.$$

*Let* $(\theta_n)_{n=1}^\infty \subset \Theta$ *be a sequence s.t.* $\theta_n$ *minimises* $D_n$ *for all* $n$ *sufficiently large. Denote by* $\widetilde{\pi}_n^D$ *a density on* $\mathbb{R}^v$ *of the random variable* $\sqrt{n}(\theta - \theta_n)$, *where* $\theta \sim \pi_n^D$. *Then*

$$\int_{\mathbb{R}^d} \left| \widetilde{\pi}_n^D(\theta) - \frac{1}{Z_*} \exp\left(-\frac{1}{2}\theta \cdot H_*\theta\right) \right| d\theta \xrightarrow{\text{a.s.}} 0$$

*where* $Z_*$ *is the normalising constant of* $\exp(-\frac{1}{2}\theta \cdot H_*\theta)$.

The proof of Theorem 3 is deferred to Appendix C. The main proof of Theorem 1 aims to show that the preconditions C1-C5 of Theorem 3 are satisfied for the particular function $D_n$ in (26), defining the DFD-Bayes generalised posterior.

*Proof of Theorem 1.* Without loss of generality, we will give the proof for $\beta = 1$ for notational convenience.[7] Let $r_{j-}(\boldsymbol{x}, \theta) := p_\theta(\boldsymbol{x}^{j-})/p_\theta(\boldsymbol{x})$ and $r_{j+}(\boldsymbol{x}, \theta) := p_\theta(\boldsymbol{x})/p_\theta(\boldsymbol{x}^{j+})$ for each $j = 1, \ldots, d$. We can write $D_n$ as

$$D_n(\theta) \overset{\theta}{=} \frac{1}{n} \sum_{i=1}^n \underbrace{\sum_{j=1}^d \left(r_{j-}(\boldsymbol{x}_i, \theta)\right)^2 - 2r_{j+}(\boldsymbol{x}_i, \theta)}_{=:R(\boldsymbol{x}_i, \theta)}.$$

In what follows we set $D(\theta) := \mathbb{E}_{X\sim p}[R(X, \theta)]$ and verify that preconditions C1-C5 of Theorem 1 are satisfied. Note that C3 holds directly by Assumption 1 and C5 is also assumed directly in Theorem 1.

**C1:** By the strong law of large numbers [Durrett, 2010, Theorem 2.5.10],

$$D_n(\theta) = \frac{1}{n} \sum_{i=1}^n R(\boldsymbol{x}_i, \theta) \quad \xrightarrow{\text{a.s.}} \quad \mathbb{E}_{X\sim p}[R(X, \theta)] = D(\theta), \tag{27}$$

---

[7]To extend the proof to arbitrary $\beta > 0$, simply replace $D_n(\theta) = \mathrm{DFD}(p_\theta \| p_n)$ in all arguments by $D_n(\theta) = \beta\,\mathrm{DFD}(p_\theta \| p_n)$. All the arguments hold immediately since $\beta$ is a constant.

provided that $\mathbb{E}_{X \sim p}[|R(X, \theta)|] < \infty$ for each $\theta \in \Theta$. Thus we must check that $\mathbb{E}_{X \sim p}[|R(X, \theta)|] < \infty$. By the triangle inequality,

$$
\begin{aligned}
\mathbb{E}_{X \sim p}[|R(X, \theta)|] &= \mathbb{E}_{X \sim p}\left[|R(X, \theta)|\right] + C(p) - C(p) \\
&= \mathbb{E}_{X \sim p}\left[\left|R(X, \theta) + 1 + \left\|\frac{\nabla^- p(X)}{p(X)}\right\|^2\right|\right] + \mathbb{E}_{X \sim p}\left[1 + \left\|\frac{\nabla^- p(X)}{p(X)}\right\|^2\right] \\
&= \mathbb{E}_{X \sim p}\left[\left\|\frac{\nabla^- p_\theta(X)}{p_\theta(X)} - \frac{\nabla^- p(X)}{p(X)}\right\|^2\right] + 1 + \mathbb{E}_{X \sim p}\left[\left\|\frac{\nabla^- p(X)}{p(X)}\right\|^2\right]
\end{aligned}
$$

where the last equality holds from Proposition 1 and both the quantities are finite by Standing Assumption 1. Hence (27) holds for every $\theta \in \Theta$.

**C2:** From Assumption 2, we have that $r_{j+}(\boldsymbol{x}, \theta)$ and $r_{j-}(\boldsymbol{x}, \theta)$ are three times continuously differentiable with respect to $\theta \in U$ for all $\boldsymbol{x} \in \mathcal{X}$, and thus $D_n(\theta)$ is three times continuously differentiable with respect to $\theta \in U$. For any $s \in \{1, 2, 3\}$,

$$
\nabla_\theta^s D_n(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla_\theta^s R(\boldsymbol{x}_i, \theta) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d \nabla_\theta^s(r_{j-}(\boldsymbol{x}_i, \theta)^2) - 2\nabla_\theta^s r_{j+}(\boldsymbol{x}_i, \theta). \tag{28}
$$

By the triangle inequality, we have an upper bound

$$
\sup_{\theta \in U} \|\nabla_\theta^s D_n(\theta)\| \le \frac{1}{n} \sum_{i=1}^n \underbrace{\sum_{j=1}^d \sup_{\theta \in U} \|\nabla_\theta^s(r_{j-}(\boldsymbol{x}_i, \theta)^2)\| + 2\sup_{\theta \in U} \|\nabla_\theta^s r_{j+}(\boldsymbol{x}_i, \theta)\|}_{=:G(\boldsymbol{x}_i)}.
$$

The quantity $\frac{1}{n} \sum_{i=1}^n G(\boldsymbol{x}_i)$ is a random variable dependent on $\{\boldsymbol{x}_i\}_{i=1}^n$. By the strong law of large numbers [Durrett, 2010, Theorem 2.5.10],

$$
\frac{1}{n} \sum_{i=1}^n G(\boldsymbol{x}_i) \xrightarrow{\text{a.s.}} \mathbb{E}_{X \sim p}[G(X)] < \infty
$$

provided that $\mathbb{E}_{X \sim p}[|G(X)|] < \infty$. Indeed, this condition holds since from positivity of $G$

$$
\mathbb{E}_{X \sim p}[|G(X)|] = \sum_{j=1}^d \mathbb{E}_{X \sim p}\left[\sup_{\theta \in U} \|\nabla_\theta^s(r_{j-}(X, \theta)^2)\|\right] + 2\mathbb{E}_{X \sim p}\left[\sup_{\theta \in U} \|\nabla_\theta^s r_{j+}(X, \theta)\|\right],
$$

where the right hand side is finite by Assumption 2. Then

$$
\limsup_{n \to \infty} \sup_{\theta \in U} \|\nabla_\theta^s D_n(\theta)\| \le \limsup_{n \to \infty} \frac{1}{n} \sum_{i=1}^n G(\boldsymbol{x}_i) = \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^n G(\boldsymbol{x}_i) \overset{\text{a.s.}}{=} \mathbb{E}_{X \sim p}[G(X)] < \infty
$$

for any $s \in \{1, 2, 3\}$, which establishes C2.

**C4:** Let $h(\boldsymbol{x}, \theta) := \nabla_\theta^2 R(\boldsymbol{x}, \theta)$. From (28), $H_n(\theta) = \frac{1}{n} \sum_{i=1}^n h(\boldsymbol{x}_i, \theta)$. By the strong law of large numbers [Durrett, 2010, Theorem 2.5.10], we have $H_n(\theta) \xrightarrow{\text{a.s.}} \mathbb{E}_{X \sim p}[h(X, \theta)]$ provided that $\mathbb{E}_{X \sim p}[\|h(X, \theta)\|] < \infty$. Indeed, this condition holds for all $\theta \in U$, since we have the upper bound

$$
\mathbb{E}_{X \sim p}[\|h(X, \theta_*)\|] \le \mathbb{E}_{X \sim p}\left[\sup_{\theta \in U} \|h(X, \theta)\|\right] \le \mathbb{E}_{X \sim p}[|G(X)|] < \infty
$$

34

where the right hand side is bounded by the preceding argument. It remains to verify that $H_* := \lim_{n\to\infty} H_n(\theta_*)$ is equal to $\nabla_\theta^2 \mathrm{DFD}(p_\theta\|p)|_{\theta=\theta_*}$, from which C4 follows since $H_*$ was assumed to be nonsingular in the statement of Theorem 1. By the Lebesgue's dominated convergence theorem, for each $\theta \in U$,

$$\lim_{n\to\infty} H_n(\theta) = \mathbb{E}_{X\sim p}[\nabla_\theta^2 R(\boldsymbol{x},\theta)] = \nabla_\theta^2 \mathbb{E}_{X\sim p}[R(\boldsymbol{x},\theta)] = \nabla_\theta^2 D(\theta).$$

provided that $\mathbb{E}_{X\sim p}[\sup_{\theta\in U}\|\nabla_\theta^2 R(\boldsymbol{x},\theta)\|] < \infty$. This condition holds for all $\theta \in U$ since $\mathbb{E}_{X\sim p}[\sup_{\theta\in U}\|\nabla_\theta^2 R(\boldsymbol{x},\theta)\|]$ $\mathbb{E}_{X\sim p}[|G(X)|] < \infty$. Since $\theta_* \in U$ in particular, $H_* = \nabla_\theta^2 D(\theta)\ |_{\theta=\theta_*} = \nabla_\theta^2 \mathrm{DFD}(p_\theta\|p)|_{\theta=\theta_*}$, as claimed.

Thus preconditions C1-C5 are satisfied and the result follows from Theorem 3. □

## B.3 Proof of Theorem 2

*Proof.* We first calculate the Fisher divergence between the generalised posterior $\pi_n^D$ and an empirical distribution $\delta_\theta^B$ of the bootstrap minimisers $\{\theta_n^{(b)}\}_{b=1}^B$, and then minimise it as a function of the weighting constant $\beta$. Recall that the score-matching divergence [Hyvärinen, 2005] is given by

$$\mathrm{D}(\pi_n^D\|\delta_\theta^B) = \frac{1}{B}\sum_{b=1}^B \underbrace{\left\|\nabla_\theta \log \pi_n^D(\theta_n^{(b)})\right\|^2}_{=(*_1)} + \underbrace{2\,\mathrm{Tr}\left(\nabla_\theta^2 \log \pi_n^D(\theta_n^{(b)})\right)}_{=(*_2)}.$$

The score function of $\pi_n^D$ is given by

$$\nabla_\theta \log \pi_n^D(\theta) = -\beta\nabla_\theta D_n(\theta) + \nabla_\theta \log \pi(\theta),$$

which is independent of the normalising constant of $\pi_n^D$. Similarly, the second derivative is $\nabla_\theta^2 \log \pi_n^D(\theta) = -\beta\nabla_\theta^2 D_n(\theta) + \nabla_\theta^2 \log \pi(\theta)$. Therefore the terms $(*_1)$ and $(*_2)$ in the Fisher divergence can be written as

$$(*_1) = \beta^2\|\nabla_\theta D_n(\theta_n^{(b)})\|^2 - 2\beta\nabla_\theta D_n(\theta_n^{(b)})\cdot\nabla_\theta \log \pi(\theta_n^{(b)}) + \|\nabla_\theta \log \pi(\theta_n^{(b)})\|^2$$
$$(*_2) = -\beta\,\mathrm{Tr}\left(\nabla_\theta^2 D_n(\theta_n^{(b)})\right) + \mathrm{Tr}\left(\nabla_\theta^2 \log \pi(\theta_n^{(b)})\right)$$

Now consider minimising the Fisher divergence $\mathrm{D}(\pi_n^D\|\delta_\theta^B)$ with respect to the weighting constant $\beta$. Plugging the terms $(*_1)$ and $(*_2)$ in the Fisher divergence, we have

$$\mathrm{D}(\pi_n^D\|\delta_\theta^B) = \frac{1}{B}\sum_{b=1}^B \beta^2\|\nabla_\theta D_n(\theta_n^{(b)})\|^2 - 2\beta\nabla_\theta D_n(\theta_n^{(b)})\cdot\nabla_\theta \log \pi(\theta_n^{(b)}) - 2\beta\,\mathrm{Tr}\left(\nabla_\theta^2 D_n(\theta_n^{(b)})\right) + C$$

where we denote any term independent of $\beta$ by $C$ in this proof. Exchanging the order of the summation and the constant $\beta$, the Fisher divergence turns out to be a quadratic function of $\beta$ as follows:

$$\mathrm{D}(\pi_n^D\|\delta_\theta^B) = \beta^2\underbrace{\frac{1}{B}\sum_{b=1}^B \|\nabla_\theta D_n(\theta_n^{(b)})\|^2}_{=(a)} - 2\beta\underbrace{\frac{1}{B}\sum_{b=1}^B \nabla_\theta D_n(\theta_n^{(b)})\cdot\nabla_\theta \log \pi(\theta_n^{(b)}) + \mathrm{Tr}\left(\nabla_\theta^2 D_n(\theta_n^{(b)})\right)}_{=(b)} + C$$

$$= a\beta^2 - 2b\beta + C = a\left(\beta - \frac{b}{a}\right)^2 - \frac{b^2}{4a^2} + C$$

35

where the last equality follows from completing the square. Therefore the Fisher divergence $D(\pi_n^D \| \delta_\theta^B)$ is minimised at $\beta_* = b/a$, that is,

$$\beta_* = \frac{\sum_{b=1}^{B} \nabla_\theta D_n(\theta_n^{(b)}) \cdot \nabla_\theta \log \pi(\theta_n^{(b)}) + \text{Tr}\left(\nabla_\theta^2 D_n(\theta_n^{(b)})\right)}{\sum_{b=1}^{B} \|\nabla_\theta D_n(\theta_n^{(b)})\|^2},$$

as claimed, where the denominator and numerator are positive immediately from the first and second assumption respectively, which assures that $\beta_* > 0$. $\qquad\square$

## C  Proof of Theorem 3: Simplified Conditions for Miller [2021]

Before showing that the preconditions C1-C5 of Theorem 3 are sufficient for [Miller, 2021, Theorem 4], we introduce the following lemma on a.s. uniform convergence used in the proof.

**Lemma 4.** *(a.s. uniform convergence) Suppose that the preconditions C1 and C2 in Theorem 3 holds for $r = 1$. Then $D_n$ a.s. converges uniformly to $D$ on the bounded convex open set $U$ in Theorem 3.*

*Proof.* Davidson [1994, Theorem 21.8] showed that $D_n \xrightarrow{a.s.} D$ uniformly on $U$ if and only if (a) $D_n \xrightarrow{a.s.} D$ pointwise on $U$ and (b) $\{D_n\}_{n=1}^\infty$ is strongly stochastically equicontinuous on $U$. The condition (a) is immediately implied by the precondition C1 of Theorem 3 and hence the condition (b) is shown in the remainder. By Davidson [1994, Theorem 21.10], $\{D_n\}_{n=1}^\infty$ is strongly stochastically equicontinuous on $U$ if there exists a stochastic sequence $\{\mathcal{L}_n\}_{n=1}^\infty$ independent of $\theta$ s.t.

$$|D_n(\theta) - D_n(\theta')| \le \mathcal{L}_n \|\theta - \theta'\|_2, \quad \forall \theta, \theta' \in U \qquad \text{and} \qquad \limsup_{n\to\infty} \mathcal{L}_n < \infty \text{ a.s.}$$

Since $D_n$ is continuously differentiable on the set $U$ by the precondition C2 of Theorem 3 with $r = 1$, the mean value theorem yields that

$$|D_n(\theta) - D_n(\theta')| \le \sup_{\theta \in U} \|\nabla_\theta D_n(\theta)\|_2 \|\theta - \theta'\|_2, \quad \forall \theta, \theta' \in U.$$

Again by the precondition C2 of Theorem 3 with $r = 1$, we have $\limsup_{n\to\infty} \sup_{\theta \in U} \|\nabla_\theta D_n(\theta)\|_2 < \infty$ a.s. Therefore, setting $\mathcal{L}_n = \sup_{\theta \in U} \|\nabla_\theta D_n(\theta)\|_2$ concludes the proof. $\qquad\square$

We now show that [Miller, 2021, Theorem 4] holds a.s. under the preconditions C1-C5 of Theorem 3, which in turn implies Theorem 3 directly. A main argument in the proof is essentially same as that of Matsubara et al. [2022] but that is modified here to allow for an arbitrary loss $D_n$.

*Proof.* In order to apply [Miller, 2021, Theorem 4], we first extend $\pi$ and $D_n$ from $\Theta$ to $\mathbb{R}^v$ by setting $\pi(\theta) = 0$ and $D_n(\theta) = \sup_{\theta \in \Theta} |D_n(\theta)| + 1$ for all $\theta \in \mathbb{R}^v \setminus \Theta$, so that we have $\pi : \mathbb{R}^v \to \mathbb{R}$, $D_n : \mathbb{R}^v \to \mathbb{R}$ and $\pi_n^D : \mathbb{R}^v \to \mathbb{R}$. Note that in Miller [2021, Theorem 4], $\{D_n\}_{n=1}^\infty$ is regarded as a sequence of deterministic functions, while here $\{D_n\}_{n=1}^\infty$ is a sequence of stochastic functions dependent of random data $\{X_i\}_{i=1}^n$. It will be shown that Miller [2021, Theorem 4] holds a.s. for the stochastic sequence $\{D_n\}_{n=1}^\infty$. We hence verify the following prerequisites (1)–(6) of [Miller, 2021, Theorem 4] a.s. hold. Recall that $H_n(\theta) = \nabla_\theta^2 D_n(\theta)$ and $H_* = \lim_{n\to\infty} H_n(\theta_*)$ from Theorem 3:

1. the prior density $\pi$ is continuous at $\theta_*$ and $\pi(\theta_*) > 0$.

2. $\theta_n \overset{a.s.}{\to} \theta_*$.

3. the Taylor expansion $D_n(\theta) = D_n(\theta_n) + (1/2)(\theta - \theta_n) \cdot H_n(\theta_n)(\theta - \theta_n) + r_n(\theta - \theta_n)$ holds on $U$ a.s. where $r_n$ is the reminder term.

4. the remainder $r_n$ of the Taylor expansion satisfies that $|r_n(\theta)| \le C\|\theta\|_2^3$, $\forall \theta \in B_\epsilon(0)$ a.s. for all $n$ sufficiently large and some $\epsilon > 0$.

5. $H_n(\theta_n) \overset{a.s.}{\to} H_*$, $H_n(\theta_n)$ is symmetric for all $n$ sufficiently large and $H_*$ is positive definite.

6. $\liminf_{n \to \infty} \left( \inf_{\theta \in \mathbb{R}^v \backslash \mathcal{B}_\epsilon(\theta_n)} D_n(\theta) - D_n(\theta_n) \right) > 0$ a.s. for any $\epsilon > 0$.

**Part (1):** The precondition C5 of Theorem 3.

**Part (2):** The strong consistency $\theta_n \overset{a.s.}{\to} \theta_*$ is shown by an argument similar to van der Vaart [1998, Theorem 5.7] or essentially same as Matsubara et al. [2022, Lemma 3]. First, it follows from Lemma 4 that $D_n \overset{a.s.}{\to} D$ uniformly on $U$ under the conditions of Theorem 3. Thus, for all $n$ sufficiently large, we can take $\delta > 0$ s.t. $|D_n(\theta) - D(\theta)| < \delta/2$ a.s. over $\theta \in U$, which in turn leads to (a) $D(\theta) < D_n(\theta) + \delta/2$ and (b) $D_n(\theta) < D(\theta) + \delta/2$ a.s. over $\theta \in U$. Then applying both (a) and (b), the following bound on $D(\theta_n)$ holds for all $n$ sufficiently large:

$$D(\theta_n) \overset{(a)}{<} D_n(\theta_n) + \delta/2 \overset{(*)}{\le} D_n(\theta_*) + \delta/2 \overset{(b)}{<} D(\theta_*) + \delta \quad \text{a.s.} \tag{29}$$

where the second inequality $(*)$ follows from the fact that $\theta_n$ is the minimiser of $D_n$. Since $\inf_{\theta \in \mathbb{R}^v} D(\theta) = \inf_{\theta \in \Theta} D(\theta)$ is uniquely attained at $\theta_* \in U$ by Theorem 3 (3), for any $\epsilon > 0$ we have $D(\theta) - D(\theta_*) > 0$ for all $\theta \in \mathbb{R}^v \backslash B_\epsilon(\theta_*)$. Given an arbitrary $\epsilon > 0$, let $\delta = \inf_{\theta \in \Theta \backslash B_\epsilon(\theta_*)} D(\theta) - D(\theta_*) > 0$. It then follows from (29) that, for all $n$ sufficiently large,

$$D(\theta_n) < \inf_{\theta \in \mathbb{R}^v \backslash B_\epsilon(\theta_*)} D(\theta) \quad \text{a.s.}$$

This implies that $\theta_n \in B_\epsilon(\theta_*)$ a.s. for any $\epsilon > 0$ arbitrary small for all $n$ sufficiently large. Therefore $\theta_n \overset{a.s.}{\to} \theta_*$ by definition of convergence.

**Part (3):** From the precondition C2 of Theorem 3, $D_n$ is 3 times continuously differentiable over $U$. Noting that $\nabla_\theta D_n(\theta) = 0$ at a minimiser $\theta_n$ of $D_n$, the Taylor expansion of $D_n$ around the minimiser $\theta_n$ gives that

$$D_n(\theta) = D_n(\theta_n) + \frac{1}{2}(\theta - \theta_n) \cdot H_n(\theta_n)(\theta - \theta_n) + r_n(\theta - \theta_n)$$

where $r_n$ is the remainder of the Taylor expansion.

**Part (4):** Since $r_n$ is the remainder of the Taylor expansion, we have an upper bound

$$|r_n(\theta - \theta_n)| \le \frac{1}{6} \sup_{\theta \in U} \|\nabla_\theta^3 D_n(\theta)\|_2 \|\theta - \theta_n\|_2^3, \quad \forall \theta \in U.$$

The precondition C2 of Theorem 3 guarantees that $\limsup_{n \to \infty} \sup_{\theta \in U} \|\nabla_\theta^3 D_n(\theta)\|_2 < \infty$ a.s. It is thus possible to take some positive constant $C$ s.t. $(1/6)\sup_{\theta \in U} \|\nabla_\theta^3 D_n(\theta)\|_2 \le C$ a.s. for all $n$

sufficiently large. For all $n$ sufficiently large, there exists some open $\epsilon$-neighbour $B_\epsilon(\theta_n)$ contained in the open set $U$ since $\theta_n \in U$. Combining these two facts concludes that

$$|r_n(\theta - \theta_n)| \le C\|\theta - \theta_n\|_2^3, \quad \forall \theta \in B_\epsilon(\theta_n) \quad \implies \quad |r_n(\theta)| \le C\|\theta\|_2^3, \quad \forall \theta \in B_\epsilon(0)$$

holds for some $\epsilon > 0$.

**Part (5):** We first show that $\|H_n(\theta_n) - H_*\|_2 \overset{a.s.}{\to} 0$. By the triangle inequality,

$$\|H_n(\theta_n) - H_*\|_2 \le \|H_n(\theta_n) - H_n(\theta_*)\|_2 + \|H_n(\theta_*) - H_*\|_2.$$

For the first term, it follows from the mean value theorem that

$$\|H_n(\theta_n) - H_n(\theta_*)\|_2 \le \sup_{\theta \in U} \|\nabla_\theta H_n(\theta)\|_2 \|\theta_n - \theta_*\|_2 = \sup_{\theta \in U} \|\nabla_\theta^3 D_n(\theta)\|_2 \|\theta_n - \theta_*\|_2.$$

The precondition C2 of Theorem 3 guarantees that $\limsup_{n\to\infty} \sup_{\theta \in U} \|\nabla_\theta^3 D_n(\theta)\|_2 < \infty$ a.s. It is thus possible to take some positive constant $C'$ s.t. $\|H_n(\theta_n) - H_n(\theta_*)\|_2 \le C'\|\theta_n - \theta_*\|_2$ for all $n$ sufficiently large. Then we have $\|H_n(\theta_n) - H_n(\theta_*)\|_2 \overset{a.s.}{\to} 0$ by the preceding part (2) $\theta_n \overset{a.s.}{\to} \theta_*$. For the second term, it is directly implied by the precondition C4 of Theorem 3 that $\|H_n(\theta_*) - H_*\|_2 \overset{a.s.}{\to} 0$. Combining these two facts concludes that $\|H_n(\theta_n) - H_*\|_2 \overset{a.s.}{\to} 0$. We next show that $H_n(\theta_n)$ is symmetric. The $(i,j)$ entry of $H_n(\theta) = \nabla_\theta^2 D_n(\theta)$ is given by the partial derivative $(\partial^2/\partial\theta_i\partial\theta_j)D_n(\theta)$ with respect to $i$-th and $j$-th entry of $\theta$. Since $D_n$ is twice continuously differentiable by the precondition C2 of Theorem 3, the Schwartz's theorem implies that the commutation $(\partial^2/\partial\theta_i\partial\theta_j)D_n(\theta) = (\partial^2/\partial\theta_j\partial\theta_i)D_n(\theta)$ holds and therefore $H_n(\theta)$ is symmetric for any $\theta \in \Theta$. Finally we show positive definiteness of $H_*$. For all $n$ sufficiently large, $H_n(\theta_n)$ is positive semi-definite by the fact that $\theta_n$ is the minimiser of $D_n$ and accordingly the limit $H_*$ is positive semi-definite. Then $H_*$ is positive definite since $H_*$ is nonsingular by the precondition C4 of Theorem 3.

**Part (6):** It holds for any sequence $a_n, b_n \in \mathbb{R}$ that $\liminf_{n\to\infty}(a_n - b_n) \ge \liminf_{n\to\infty} a_n + \liminf_{n\to\infty}(-b_n)$. Furthermore from the property that $\liminf_{n\to\infty}(-b_n) = -\limsup_{n\to\infty} b_n$, we have $\liminf_{n\to\infty}(a_n - b_n) \ge \liminf_{n\to\infty} a_n - \limsup_{n\to\infty} b_n$. Applying this, we have

$$\liminf_{n\to\infty}\left(\inf_{\theta\in\mathbb{R}^v\backslash B_\epsilon(\theta_n)} D_n(\theta) - D_n(\theta_n)\right) = \underbrace{\liminf_{n\to\infty} \inf_{\theta\in\mathbb{R}^v\backslash B_\epsilon(\theta_n)} D_n(\theta)}_{=:(*_1)} - \underbrace{\limsup_{n\to\infty} D_n(\theta_n)}_{=:(*_2)}.$$

For the first term $(*_1)$, it is obvious from the way of extending $D_n$ from $\Theta$ to $\mathbb{R}^v$ that

$$(*_1) = \liminf_{n\to\infty} \inf_{\theta\in\mathbb{R}^v\backslash B_\epsilon(\theta_n)} D_n(\theta) \ge \liminf_{n\to\infty} \inf_{\theta\in\Theta\backslash B_\epsilon(\theta_n)} D_n(\theta) \quad \text{a.s.}$$

For any set $A \subset \mathbb{R}^v$ and function $g : \mathbb{R}^v \to \mathbb{R}$, define $\inf_{\theta\in A\backslash B_\epsilon(\theta_n)} g(\theta) := \sup_{\theta\in A} g(\theta)$ if $A \backslash B_\epsilon(\theta_n)$ is empty. Decomposing $\Theta$ into two sets $U$ and $\Theta \backslash U$ leads to

$$(*_1) \ge \liminf_{n\to\infty} \inf_{\theta\in\Theta\backslash B_\epsilon(\theta_n)} D_n(\theta) \ge \min\left(\underbrace{\liminf_{n\to\infty} \inf_{\theta\in U\backslash B_\epsilon(\theta_n)} D_n(\theta)}_{=:(*_{11})}, \underbrace{\liminf_{n\to\infty} \inf_{\theta\in\Theta\backslash(U\cup B_\epsilon(\theta_n))} D_n(\theta)}_{=:(*_{12})}\right) \quad \text{a.s.}$$

For the term $(*_{11})$, since $D_n \overset{a.s.}{\to} D$ uniformly on $U$ by Lemma 4 and $\theta_n \overset{a.s.}{\to} \theta_*$ by the preceding part (2),

$$(*_{11}) = \liminf_{n \to \infty} \inf_{\theta \in U \backslash B_\epsilon(\theta_n)} D_n(\theta) = \lim_{n \to \infty} \inf_{\theta \in U \backslash B_\epsilon(\theta_n)} D_n(\theta) = \inf_{\theta \in U \backslash B_\epsilon(\theta_*)} D(\theta) \quad \text{a.s.}$$

For the term $(*_{12})$, since the global minimiser $\theta_n$ of $D_n$ is contained in $U$ a.s. for all $n$ sufficiently large by the precondition C3 of Theorem 3,

$$(*_{12}) = \liminf_{n \to \infty} \inf_{\theta \in \Theta \backslash (U \cup B_\epsilon(\theta_n))} D_n(\theta) > \liminf_{n \to \infty} \inf_{\theta \in U} D_n(\theta) = \inf_{\theta \in U} D(\theta) = D(\theta_*) \quad \text{a.s.}$$

where the second equality follows from the a.s. uniform convergence of $D_n$ on $U$ by Lemma 4. For the second term $(*_2)$, again since $D_n \overset{a.s.}{\to} D$ uniformly on $U$ and $\theta_n \overset{a.s.}{\to} \theta_*$, we have

$$(*_2) = \limsup_{n \to \infty} D_n(\theta_n) = \lim_{n \to \infty} D_n(\theta_n) = D(\theta_*) \quad \text{a.s.}$$

The original term $(*_1) - (*_2)$ is lower bounded by $(*_1) - (*_2) \geq \min((*_{11}) - (*_2), (*_{12}) - (*_2))$ a.s., and both the term $(*_{11}) - (*_2)$ and $(*_{12}) - (*_2)$ are then further lower bounded by

$$(*_{11}) - (*_2) = \inf_{\theta \in U \backslash B_\epsilon(\theta_*)} D(\theta) - D(\theta_*) > 0 \quad \text{and} \quad (*_{12}) - (*_2) > D(\theta_*) - D(\theta_*) = 0 \quad \text{a.s.}$$

where the first inequality follows from the precondition C3 of Theorem 3 indicating that $\inf_{\theta \in \Theta} D(\theta)$ is uniquely attained at $\theta_* \in U$. Therefore we have $(*_1) - (*_2) \geq \min((*_{11}) - (*_2), (*_{12}) - (*_2)) > 0$ a.s., which concludes the proof. $\square$

# D  Relation to Stein Discrepancies

Fisher divergences can be related to a more general class of divergences called Stein discrepancies. Since their introduction, Stein discrepancies have demonstrated utility over a range of statistical applications, including hypothesis testing, parameter estimation, variational inference, and post-processing of Markov chain Monte Carlo; see Anastasiou et al. [2023] for a review.

This appendix clarifies the sense in which discrete Fisher divergence can be seen as a special case of a discrete Stein discrepancy with an $L^2$-based Stein set. The continuous case was previously covered by Theorem 2 in Barp et al. [2019]. As a consequence, we deduce that the discrete Fisher divergence is stronger than the popular class of Stein discrepancies based on reproducing kernels.

## D.1  Background on Stein Discrepancies

Let $\mathcal{X}_*$ be a locally compact Hausdorff space. For a set $\mathcal{H}$ of functions $f : \mathcal{X}_* \to \mathbb{R}^d$, an operator $S_p : \mathcal{H} \to L^1(p, \mathbb{R}^m)$ depending on a probability distribution $p$ on $\mathcal{X}_*$ is called a *Stein operator* if $\mathbb{E}_{X \sim p}[S_p[h](X)] = 0$ for all $h \in \mathcal{H}$. In these circumstances, we refer to $\mathcal{H}$ as a *Stein set*. The next proposition defines a particular Stein operator that arises naturally when considering discrete domains $\mathcal{X}_* = \mathcal{X}$, where we recall that $\mathcal{X}$ is a countable space in Standing Assumption 1. The reader is referred to Shi et al. [2022] for discussion of alternative Stein operators in the discrete context. Define the forward divergence operator $\nabla^+ \cdot$ for a $\mathbb{R}^d$-valued function $h : \mathcal{X}_* \to \mathbb{R}^d$ by $\nabla^+ \cdot h(\boldsymbol{x}) = \sum_{j=1}^d h(\boldsymbol{x}^{j+}) - h(\boldsymbol{x})$.

**Proposition 2.** *Let $p$ be a positive probability distribution on $\mathcal{X}$, such that $(\nabla^- p)/p \in L^2(p, \mathbb{R}^d)$. Define an operator $S_p$, acting on functions $h \in L^2(p, \mathbb{R}^d)$, by*

$$S_p[h](\boldsymbol{x}) := \frac{\nabla^- p(\boldsymbol{x})}{p(\boldsymbol{x})} \cdot h(\boldsymbol{x}) + \nabla^+ \cdot h(\boldsymbol{x}). \tag{30}$$

*Then it holds that $\mathbb{E}_{X \sim p}[S_p[h](X)] = 0$.*

*Proof.* By positivity of $p$ and Cauchy–Schwarz, observe that

$$\sum_{\boldsymbol{x} \in \mathcal{X}} |p(\boldsymbol{x}^{i-})h_i(\boldsymbol{x})| = \sum_{\boldsymbol{x} \in \mathcal{X}} p(\boldsymbol{x})\frac{p(\boldsymbol{x}^{i-})}{p(\boldsymbol{x})}|h_i(\boldsymbol{x})| = \mathbb{E}_{X \sim p}\left[\frac{p(X^{i-})}{p(X)}|h_i(X)|\right] \tag{31}$$

$$\leq \mathbb{E}_{X \sim p}\left[\frac{p(X^{i-})^2}{p(X)^2}\right] \mathbb{E}_{X \sim p}\left[h_i(X)^2\right] < \infty$$

where the first and second term are implied to be finite for any $i = 1, \ldots, d$ since $h \in L^2(p, \mathbb{R}^d)$ and $(\nabla^- p)/p \in L^2(p, \mathbb{R}^d)$ which implies $[\nabla^- p(\boldsymbol{x})/p(\boldsymbol{x})]_i = 1 - p(\boldsymbol{x}^{i-})/p(\boldsymbol{x})$ is square integrable with respect to $p$.

Now, using the definition of $\nabla^-$ and $\nabla^+ \cdot$, the Stein operator $S_p$ can be simplified as

$$S_p[h](\boldsymbol{x}) = \sum_{i=1}^{d} h_i(\boldsymbol{x}^{i+}) - \frac{p(\boldsymbol{x}^{i-})}{p(\boldsymbol{x})}h_i(\boldsymbol{x}). \tag{32}$$

The expectation of interest can then be expressed as

$$\mathbb{E}_{X \sim p}[S_p[h](X)] = \sum_{\boldsymbol{x} \in \mathcal{X}} p(\boldsymbol{x})S_p[h](\boldsymbol{x}) = \sum_{i=1}^{d} \sum_{\boldsymbol{x} \in \mathcal{X}} p(\boldsymbol{x})h_i(\boldsymbol{x}^{i+}) - \sum_{\boldsymbol{x} \in \mathcal{X}} p(\boldsymbol{x}^{i-})h_i(\boldsymbol{x}),$$

where we have used the absolute convergence of the series, established in (31), to justify the re-ordering of terms. The result is then immediate from Lemma 2. □

The Stein operator (30) can be considered a discrete analogue of the Langevin Stein operator for continuous domains; see Yang et al. [2018].

Given a Stein operator $S_p$ and Stein set $\mathcal{H}$, the *Stein discrepancy* between probability distributions $p$ and $q$ on $\mathcal{X}$ is defined as the maximum deviation between expectations of the test functions $S_p[h]$ for $h \in \mathcal{H}$:

$$\text{SD}(p\|q) := \sup_{h \in \mathcal{H}} |\mathbb{E}_{X \sim q}[S_p[h](X)] - \mathbb{E}_{X \sim p}[S_p[h](X)]| = \sup_{h \in \mathcal{H}} |\mathbb{E}_{X \sim q}[S_p[h](X)]| \tag{33}$$

The final equality follows from Proposition 2, and our discussion in this appendix implicitly assumes all relevant quantities are well-defined. The Stein discrepancy is computable[8] without knowing the normalising constant of $p$ since it depends on $p$ only through the ratio $(\nabla^- p)/p$, in a similar manner to discrete Fisher divergence in the main text.

---

[8]That is, the expectations do not involve the normalising constant; whether the supremum over the Stein set is computable depends on how the Stein set is selected.

## D.2 Discrete Fisher Divergence as a Stein Discrepancy

We now establish that the discrete Fisher divergence, introduced in the main text, is in fact a Stein discrepancy, corresponding to the Stein operator in Proposition 2 and a Stein set equal to the unit ball of $L^2(q, \mathbb{R}^d)$. This observation will allow us to conclude, in Appendix D.3, that discrete Fisher divergence is stronger than popular kernel Stein discrepancies.

**Proposition 3.** *Let $p$ and $q$ be positive distributions on $\mathcal{X}$, such that $(\nabla^- p)/p, (\nabla^- q)/q \in L^2(q, \mathbb{R}^d)$. Consider a Stein discrepancy whose Stein operator is* (30) *and whose Stein set is $\mathcal{H} = \{h : \mathcal{X} \to \mathbb{R}^d \mid \sum_{i=1}^d \mathbb{E}_{X \sim q}[h_i(X)^2] \leq 1\}$. Then*

$$\mathrm{SD}(p\|q) = \sqrt{\mathrm{DFD}(p\|q)}. \tag{34}$$

*Proof.* From (30) and (33), we have that

$$\mathrm{SD}(p\|q) = \sup_{h \in \mathcal{H}} \left| \mathbb{E}_{X \sim q} \left[ \frac{\nabla^- p(X)}{p(X)} \cdot h(X) - \frac{\nabla^- q(X)}{q(X)} \cdot h(X) \right] \right|.$$

Note that $L^2(q, \mathbb{R}^d)$ is a Hilbert space when equipped with the inner product $\langle f, g \rangle_{L^2(q, \mathbb{R}^d)} := \mathbb{E}_{X \sim q}[f(X) \cdot g(X)]$. Thus we can view $\mathrm{SD}(p\|q)$ as the maximum of the inner product

$$\mathrm{SD}(p\|q) = \sup_{h \in \mathcal{H}} \left| \left\langle \frac{\nabla^- p}{p} - \frac{\nabla^- q}{q}, h \right\rangle_{L^2(q, \mathbb{R}^d)} \right|, \tag{35}$$

which is well-defined since $u := (\nabla^- p)/p - (\nabla^- q)/q \in L^2(q, \mathbb{R}^d)$. Let $\|\cdot\|_{L^2(q, \mathbb{R}^d)}$ denote the norm of $L^2(q, \mathbb{R}^d)$, so that $\mathcal{H}$ is the set of $f \in L^2(q, \mathbb{R}^d)$ for which $\|f\|_{L^2(q, \mathbb{R}^d)} \leq 1$. By the Cauchy–Schwarz inequality, the inner product in (35) attains its supremum at $h = u/\|u\|_{L^2(q, \mathbb{R}^d)} \in \mathcal{H}$. Therefore

$$\mathrm{SD}(p\|q) = \sup_{h \in \mathcal{H}} |\langle u, h \rangle_{L^2(q, \mathbb{R}^d)}| = \|u\|_{L^2(q, \mathbb{R}^d)} = \sqrt{\mathbb{E}_{X \sim q} \left[ \left\| \frac{\nabla^- p(X)}{p(X)} - \frac{\nabla^- q(X)}{q(X)} \right\|^2 \right]},$$

which concludes the proof. $\qquad\square$

## D.3 The Fisher Divergence Dominates the Kernel Stein Discrepancy

A popular choice of Stein set $\mathcal{H}$, that can lead to a closed form Stein discrepancy, is the unit ball of a reproducing kernel Hilbert space. The resulting *kernel Stein discrepancy* was recently considered in the discrete context in Yang et al. [2018]. In this appendix we establish that our discrete Fisher divergence, introduced in the main text, is a stronger notion of divergence than kernel Stein discrepancy. This may render the discrete Fisher divergence more statistically efficient in applications where a statistical model is well-specified, in addition to the computational advantage (Remark 1) and the non-reliance on a user-specified kernel discussed in the main text.

A symmetric, positive definite function $k : \mathcal{X}_* \times \mathcal{X}_* \to \mathbb{R}$ is called a kernel. For every kernel $k$, there exists a unique associated Hilbert space of real-valued functions on $\mathcal{X}_*$, called a reproducing kernel Hilbert space, denoted $\mathcal{H}_k$; see e.g. Berlinet and Thomas-Agnan [2011] for background. Let $\mathcal{H}_k^d := \mathcal{H}_k \times \cdots \times \mathcal{H}_k$, that is, a space of functions $h : \mathcal{X}_* \to \mathbb{R}^d$ whose each $i$-th output-coordinate $h_i : \mathcal{X}_* \to \mathbb{R}$ belongs to $\mathcal{H}_k$. Yang et al. [2018] studied the Stein discrepancy for a discrete space

$\mathcal{X}_* = \mathcal{X}$, of finite cardinality only, using the Stein operator (30) and a Stein set $\mathcal{H} = \{h \in \mathcal{H}_k^d : \sum_{i=1}^d \|h_i\|_{\mathcal{H}_k}^2 \leq 1\}$. Here we first establish that the Stein set $\{h \in \mathcal{H}_k^d : \sum_{i=1}^d \|h_i\|_{\mathcal{H}_k}^2 \leq 1\}$ constructed from $\mathcal{H}_k^d$ is contained in another Stein set $\{h \in L^2(q, \mathbb{R}^d) : \|h\|_{L^2(q,\mathbb{R}^d)}^2 \leq 1\}$ constructed from $L^2(q, \mathbb{R}^d)$ for any general domain $\mathcal{X}_*$, under a standard condition on the reproducing kernel. This in turn shows that the discrete Fisher divergence dominates the kernel Stein discrepancy.

**Proposition 4.** *Let $q$ be a probability distribution on $\mathcal{X}_*$. Let $k : \mathcal{X}_* \times \mathcal{X}_* \to \mathbb{R}$ be a kernel such that $k(\boldsymbol{x}, \boldsymbol{x}) \leq 1$ for all $\boldsymbol{x} \in \mathcal{X}_*$. Then the unit ball of $\mathcal{H}_k^d$ is contained in the unit ball of $L^2(q, \mathbb{R}^d)$.*

*Proof.* First let $f : \mathcal{X}_* \to \mathbb{R}^d$ be any element of $\mathcal{H}_k^d$, where its $i$-th output-coordinate $f_i : \mathcal{X}_* \to \mathbb{R}$ belongs to $\mathcal{H}_k$ each. From the reproducing property of $\mathcal{H}_k$, followed by the Cauchy–Schwartz inequality, the norm of $f$ in $L^2(q, \mathbb{R}^d)$ is upper bounded as follows:

$$
\begin{aligned}
\|f\|_{L^2(q,\mathbb{R}^d)}^2 = \sum_{i=1}^d \mathbb{E}_{X \sim q}[f_i(X)^2] &= \sum_{i=1}^d \mathbb{E}_{X \sim q}[\langle f_i(\cdot), k(X, \cdot)\rangle_{\mathcal{H}_k}^2] \\
&\leq \sum_{i=1}^d \mathbb{E}_{X \sim q}\left[\|f_i\|_{\mathcal{H}_k}^2 \|k(X, \cdot)\|_{\mathcal{H}_k}^2\right] = \sum_{i=1}^d \mathbb{E}_{X \sim q}\left[\|f_i\|_{\mathcal{H}_k}^2 k(X, X)\right] \\
&= \left(\sum_{i=1}^d \|f_i\|_{\mathcal{H}_k}^2\right) \mathbb{E}_{X \sim q}\left[k(X, X)\right] = \|f\|_{\mathcal{H}_k^d}^2 \, \mathbb{E}_{X \sim q}\left[k(X, X)\right].
\end{aligned}
$$

The continuous embedding of $\mathcal{H}_k^d$ in $L^2(q, \mathbb{R}^d)$ therefore holds, and moreover the embedding constant is at most one, since $\mathbb{E}_{X \sim q}[k(X, X)] \leq 1$ due to the assumption that $k(\boldsymbol{x}, \boldsymbol{x}) \leq 1$ for all $\boldsymbol{x} \in \mathcal{X}_*$. In particular, it follows that the unit ball of $\mathcal{H}_k^d$ is contained in the unit ball of $L^2(q, \mathbb{R}^d)$. $\qquad\square$

Built upon Proposition 4, we can immediately show the discrete Fisher divergence dominates the kernel Stein discrepancy for the case where $\mathcal{X}_* = \mathcal{X}$.

**Proposition 5.** *Let $p$ and $q$ be positive distributions on $\mathcal{X}$, such that $(\nabla^- p)/p, (\nabla^- q)/q \in L^2(q, \mathbb{R}^d)$. Let $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ be a kernel such that $k(\boldsymbol{x}, \boldsymbol{x}) \leq 1$ for all $\boldsymbol{x} \in \mathcal{X}$. Let $\mathcal{S}_p$ be a Stein operator in (30). Then the kernel Stein discrepancy, denoted $\mathrm{SD}_k$, satisfies $\mathrm{SD}_k(p\|q) \leq \sqrt{\mathrm{DFD}(p\|q)}$.*

*Proof.* From (33) (which in turn relies on Proposition 2), it is straightforward to see that

$$
\mathrm{SD}_k(p\|q) = \sup_{\|h\|_{\mathcal{H}_k^d} \leq 1} |\mathbb{E}_{X \sim q}\left[S_p[h](X)\right]| \leq \sup_{\|h\|_{L^2(q,\mathbb{R}^d)} \leq 1} |\mathbb{E}_{X \sim q}\left[S_p[h](X)\right]| = \sqrt{\mathrm{DFD}(p\|q)},
$$

where the inequality follows from Proposition 4 immediately and the final equality is Proposition 3. $\qquad\square$

This argument is not restricted to the discrete case but is immediately applicable for the continuous case. One of the most common Stein operator for a continuous domain $\mathcal{X}_* = \mathbb{R}^d$ is

$$
\mathcal{S}_p[h](\boldsymbol{x}) = \nabla \log p(\boldsymbol{x}) \cdot h(\boldsymbol{x}) + \nabla \cdot h(\boldsymbol{x}) \tag{36}
$$

The Fisher divergence $\mathrm{FD}(p\|q) = \mathbb{E}_{X \sim q}[\|\nabla \log p(X) - \nabla \log q(X)\|^2]$ for densities $p, q$ on $\mathcal{X}_*$ dominates the kernel Stein discrepancy constructed from the above Stein operator and the kernel on $\mathcal{X}_*$.

**Proposition 6.** *Let $\mathcal{X}_* = \mathbb{R}^d$. Let $p$ and $q$ be positive continuously differentiable densities on $\mathcal{X}_*$, such that $\nabla \log p, \nabla \log q \in L^2(q, \mathbb{R}^d)$. Let $k : \mathcal{X}_* \times \mathcal{X}_* \to \mathbb{R}$ be a kernel such that $k(\boldsymbol{x}, \boldsymbol{x}) \leq 1$ for all $\boldsymbol{x} \in \mathcal{X}_*$. Let $\mathcal{S}_p$ be a Stein operator in (36) Then the kernel Stein discrepancy, denoted $\mathrm{SD}_k$, satisfies $\mathrm{SD}_k(p\|q) \leq \sqrt{\mathrm{FD}(p\|q)}$.*

*Proof.* We repeat the same argument as Proposition 5. From [Barp et al., 2019, Theorem 2], $\mathrm{FD}(p\|q)$ can be written as the Stein discrepancy constructed by the Stein set $\{h \in L^2(q, \mathbb{R}^d) : \|h\|^2_{L^2(q, \mathbb{R}^d)} \leq 1\}$. Then from (33) (which in turn relies on Proposition 2), it is straightforward to see that

$$\mathrm{SD}_k(p\|q) = \sup_{\|h\|_{\mathcal{H}_k^d} \leq 1} |\mathbb{E}_{X \sim q}[S_p[h](X)]| \leq \sup_{\|h\|_{L^2(q, \mathbb{R}^d)} \leq 1} |\mathbb{E}_{X \sim q}[S_p[h](X)]| = \sqrt{\mathrm{FD}(p\|q)},$$

where the inequality follows from Proposition 4 immediately and the final equality is Proposition 3. $\square$

An interesting recent observation in Shi et al. [2022] was that alternative Stein operators [such as Gibbs and Barker operators; see Table 1 of Shi et al., 2022] gave rise to kernel Stein discrepancies that performed better in their particular context (low-variance gradient estimation). It would be interesting to explore the analogous alternatives to discrete Fisher divergence that would result from such operators, but this is left to future work.

## D.4   Robustness of the Kernel Stein Discrepancy

Appendix D.3 indicates statistical efficiency of the discrete Fisher divergence over the kernel Stein discrepancy. If one's model is well-specified, minimising the discrete Fisher divergence leads us to a correct model faster than the kernel Stein discrepancy. However this does not mean that the use of the discrete Fisher divergence is always better than the kernel Stein discrepancy. In particular, the kernel Stein discrepancy can be equipped with strong robustness by choosing an appropriate kernel. To demonstrate this, we compare three posteriors of Pseudo-Bayes, DFD-Bayes, and KSD-Bayes for the same Ising model as Section 4.2 with $d = 100$ ($m = 10$) in a setting where dataset contains extreme outliers with a proportion $\epsilon$.

We approximately draw 1000 samples $\{x_i\}_{i=1}^{1000}$ from the Ising model $p_\theta$ with $\theta = 5$ by the same Metropolis–Hastings algorithm as Section 4.2. To study the robustness of the posteriors, we replaced a proportion $\epsilon = 0.1$ of the data with the vector $(1, 1, \cdots, 1)$ corresponding to the extreme value in $\mathcal{X}$ that is rarely drawn from the model. Matsubara et al. [2022] showed that KSD-Bayes can satisfy strong qualitative robustness called "global bias-robustness" by choosing a kernel appropriately. For this example, we use the same choice of kernel as Matsubara et al. [2022] below:

$$k(\boldsymbol{x}, \boldsymbol{x}') = m(\boldsymbol{x}) \exp\left( -\frac{1}{d} \sum_{i=1}^d \mathbb{1}(x_i - x_i') \right) m(\boldsymbol{x})$$

where $m(x) = \sigma(90 - |\sum_i x_i|)$ based on a sigmoid function $\sigma(t) = (1 + \exp(-t))^{-1}$. This is indeed a proper choice of kernel, and the function $m(\boldsymbol{x})$ in the definition of kernel is designed to restrict the influence of extreme data whose norm is closer to or larger than 90.

In Figure 11 demonstrated that KSD-Bayes offered a correct inference outcome even when the dataset contains outliers, being less affected by the outliers. On the other hand, the Pseudo-Bayes
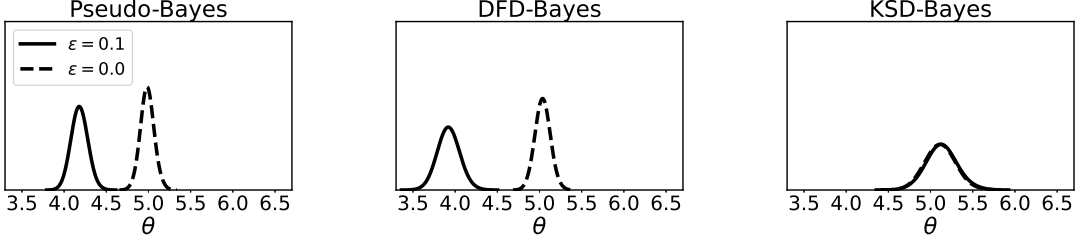
Figure 11: Posteriors of Pseudo-Bayes (left), DFD-Bayes (centre), and KSD-Bayes (right) for the Ising model in the presence of outlier with $\epsilon = 0.1$ and no outlier with $\epsilon = 0.0$.

and DFD-Bayes posteriors placed the majority of the probability mass on smaller $\theta$ than the correct value $\theta = 5$. The extreme value $(1, 1, \cdots, 1)$ of the outliers is more likely to be drawn from the model of $\theta \ll 1$; the posteriors of Pseudo-Bayes and DFD-Bayes were thus pulled in the direction of smaller $\theta$.

# E  Calculations for Worked Examples

## E.1  Assumption 2 for Example 1

The aim of this section is to establish when Assumption 2 is satisfied for the exponential family model in Example 1. For better presentation, let $T_{j-}(\boldsymbol{x}) := T(\boldsymbol{x}^{j-}) - T(\boldsymbol{x})$ and $b_{j-}(\boldsymbol{x}) := b(\boldsymbol{x}^{j-}) - b(\boldsymbol{x})$ to see that $r_{j-}(\boldsymbol{x}, \theta) = \exp(\eta(\theta) \cdot T_{j-}(\boldsymbol{x}) + b_{j-}(\boldsymbol{x}))$. In addition, let $T_{j+}(\boldsymbol{x}) := T(\boldsymbol{x}) - T(\boldsymbol{x}^{j+})$ and $b_{j+}(\boldsymbol{x}) := b(\boldsymbol{x}) - b(\boldsymbol{x}^{j+})$ to see that $r_{j-}(\boldsymbol{x}^{j+}, \theta) = \exp(\eta(\theta) \cdot T_{j+}(\boldsymbol{x}) + b_{j+}(\boldsymbol{x}))$. It is straightforward to see that, for any $\boldsymbol{x} \in \mathcal{X}$,

$$\nabla_\theta r_{j-}(\boldsymbol{x}^{j+}, \theta) = \nabla_\theta \eta(\theta) \cdot T_{j+}(\boldsymbol{x}) \exp(\eta(\theta) \cdot T_{j+}(\boldsymbol{x}) + b_{j+}(\boldsymbol{x}))$$
$$= \nabla_\theta \eta(\theta) \cdot T_{j+}(\boldsymbol{x}) \exp(\eta(\theta) \cdot T_{j+}(\boldsymbol{x})) \exp(b_{j+}(\boldsymbol{x}))$$
$$\nabla_\theta (r_{j-}(\boldsymbol{x}, \theta)^2) = 2 r_{j-}(\boldsymbol{x}, \theta) \nabla_\theta r_{j-}(\boldsymbol{x}, \theta)$$
$$= 2 \nabla_\theta \eta(\theta) \cdot T_{j-}(\boldsymbol{x}) \exp(2\eta(\theta) \cdot T_{j-}(\boldsymbol{x})) \exp(2 b_{j-}(\boldsymbol{x}))$$

By assumption, $T_{j-}(\boldsymbol{x})$ is bounded over all $\boldsymbol{x} \in \mathcal{X}$, which in turn shows that $T_{j+}(\boldsymbol{x}) = T_{j-}(\boldsymbol{x}^{j+})$ is bounded over all $\boldsymbol{x} \in \mathcal{X}$ since $\boldsymbol{x}^{j+} \in \mathcal{X}$. Further, by assumption, $\sup_{\theta \in U} \|\nabla_\theta \eta(\theta)\| < \infty$ and $\sup_{\theta \in U} \|\eta(\theta)\| < \infty$. Let $M$ be a constant that upper bounds all the terms $\sup_{\boldsymbol{x} \in \mathcal{X}} \|T_{j-}(\boldsymbol{x})\|$, $\sup_{\boldsymbol{x} \in \mathcal{X}} \|T_{j+}(\boldsymbol{x})\|$, $\sup_{\theta \in U} \|\nabla_\theta \eta(\theta)\|$ and $\sup_{\theta \in U} \|\eta(\theta)\|$. Then we have

$$\sup_{\theta \in U} \|\nabla_\theta r_{j-}(\boldsymbol{x}^{j+}, \theta)\| \leq M^2 \exp\left(M^2\right) \exp(b_{j+}(\boldsymbol{x})),$$
$$\sup_{\theta \in U} \|\nabla_\theta (r_{j-}(\boldsymbol{x}, \theta)^2)\| \leq 2M^2 \exp\left(2M^2\right) \exp(2 b_{j-}(\boldsymbol{x})).$$

Taking the expectations,

$$\mathbb{E}_{X \sim p}\left[\sup_{\theta \in U} \|\nabla_\theta r_{j-}(X^{j+}, \theta)\|\right] \leq M^2 \exp(M^2) \mathbb{E}_{X \sim p}\left[\exp(b_{j+}(X))\right], \tag{37}$$

$$\mathbb{E}_{X \sim p}\left[\sup_{\theta \in U} \|\nabla_\theta (r_{j-}(X, \theta)^2)\|\right] \leq 2M^2 \exp(2M^2) \mathbb{E}_{X \sim p}\left[\exp(2 b_{j-}(X))\right]. \tag{38}$$

44

By assumption $\mathbb{E}_{X \sim p}\left[\exp(2b_{j-}(X))\right] = \mathbb{E}_{X \sim p}\left[\exp(b_{j-}(X))^2\right] < \infty$ , and we now argue that this also implies $\mathbb{E}_{X \sim p}\left[\exp(b_{j+}(X))\right] < \infty$. Indeed, from Lemma 2,

$$
\begin{aligned}
\mathbb{E}_{X \sim p}\left[\exp(b_{j+}(X))\right] &= \sum_{\boldsymbol{x} \in \mathcal{X}} p(\boldsymbol{x}) \exp(b(\boldsymbol{x}) - b(\boldsymbol{x}^{j+})) = \sum_{\boldsymbol{x} \in \mathcal{X}} p(\boldsymbol{x}^{j-}) \exp(b(\boldsymbol{x}^{j-}) - b(\boldsymbol{x})) \\
&= \sum_{\boldsymbol{x} \in \mathcal{X}} p(\boldsymbol{x}) \frac{p(\boldsymbol{x}^{j-})}{p(\boldsymbol{x})} \exp(b(\boldsymbol{x}^{j-}) - b(\boldsymbol{x})) = \mathbb{E}_{X \sim p}\left[\frac{p(X^{j-})}{p(X)} \exp(b_{j-}(X))\right].
\end{aligned}
$$

Now, using the Cauchy–Schwartz inequality,

$$
\mathbb{E}_{X \sim p}\left[\exp(b_{j+}(X))\right] \leq \mathbb{E}_{X \sim p}\left[\frac{p(X^{j-})^2}{p(X)^2}\right] \mathbb{E}_{X \sim p}\left[\exp(2b_{j-}(X))\right]. \tag{39}
$$

Existence of the first term in (39) is implied by the Standing Assumption $(\nabla^- p)/p \in L^2(p, \mathbb{R}^d)$, while existence of the second term in (39) was assumed. Therefore we have shown that (37) and (38) exist. Repeating an essentially identical argument, it is straightforward to see also that $\mathbb{E}_{X \sim p}\left[\sup_{\theta \in U} \|\nabla_\theta^s r_{j-}(X^{j+}, \theta)\|\right] < \infty$ and $\mathbb{E}_{X \sim p}\left[\sup_{\theta \in U} \|\nabla_\theta^s (r_{j-}(X, \theta)^2)\|\right] < \infty$ for $s = 2, 3$ as claimed.

## E.2    Derivatives of (10) for Example 1

Automatic differentiation is an attractive and promising choice to compute (10) whenever it is available. Nonetheless, it is still straightforward for a majority of parametric models to compute the loss derivatives used in (10). This section aims to demonstrate a form of loss derivatives for a model in Example 1. The optimal $\beta$ of (10) depends on the first and second derivative of a loss $D$ specified by users. Consider the discrete Fisher divergence $D_n$ that this paper established. The discrete Fisher divergence $D_n(\theta) = \mathrm{DFD}(p_\theta \| p_n)$ between a model $p_\theta$ in Example 1 and an empirical distribution $p_n$ of data $\{\boldsymbol{x}_i\}_{i=1}^n$ is given as

$$
D_n(\theta) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d (r_{j-}(\boldsymbol{x}, \theta))^2 - 2r_{j-}(\boldsymbol{x}^{j+}, \theta)
$$

For simplicity, let $\eta(\theta) = \theta$ here. Then $r_{j-}(\boldsymbol{x}, \theta) = \exp(\theta \cdot T_{j-}(\boldsymbol{x}) + b_{j-}(\boldsymbol{x}))$ and $r_{j-}(\boldsymbol{x}^{j+}, \theta) = \exp(\theta \cdot T_{j+}(\boldsymbol{x}) + b_{j+}(\boldsymbol{x}))$ using the notations in Appendix E.1. Therefore the derivatives are

$$
\begin{aligned}
\nabla_\theta r_{j-}(\boldsymbol{x}, \theta) &= T_{j-}(\boldsymbol{x}) \exp(\theta \cdot T_{j-}(\boldsymbol{x}) + b_{j-}(\boldsymbol{x})), \\
\nabla_\theta^2 r_{j-}(\boldsymbol{x}, \theta) &= T_{j-}(\boldsymbol{x}) \otimes T_{j-}(\boldsymbol{x}) \exp(\theta \cdot T_{j-}(\boldsymbol{x}) + b_{j-}(\boldsymbol{x})), \\
\nabla_\theta r_{j-}(\boldsymbol{x}^{j+}, \theta) &= T_{j+}(\boldsymbol{x}) \exp(\theta \cdot T_{j+}(\boldsymbol{x}) + b_{j+}(\boldsymbol{x})), \\
\nabla_\theta^2 r_{j-}(\boldsymbol{x}^{j+}, \theta) &= T_{j+}(\boldsymbol{x}) \otimes T_{j+}(\boldsymbol{x}) \exp(\theta \cdot T_{j+}(\boldsymbol{x}) + b_{j+}(\boldsymbol{x}))
\end{aligned}
$$

where $\otimes$ denotes outer product. Built upon these components, we have the required first derivatives of $D_n(\theta)$

$$\nabla_\theta D_n(\theta) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d \nabla_\theta \left( r_{j-}(\boldsymbol{x}, \theta)^2 \right) - 2\nabla_\theta r_{j-}(\boldsymbol{x}^{j+}, \theta)$$

$$= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d 2r_{j-}(\boldsymbol{x}, \theta)\nabla_\theta r_{j-}(\boldsymbol{x}, \theta) - 2\nabla_\theta r_{j-}(\boldsymbol{x}^{j+}, \theta)$$

$$= \frac{2}{n} \sum_{i=1}^n \sum_{j=1}^d T_{j-}(\boldsymbol{x}) \, \exp(\theta \cdot T_{j-}(\boldsymbol{x}) + b_{j-}(\boldsymbol{x}))^2 - T_{j+}(\boldsymbol{x}) \, \exp(\theta \cdot T_{j+}(\boldsymbol{x}) + b_{j+}(\boldsymbol{x}))$$

as well as the second derivative of $D_n(\theta)$

$$\nabla_\theta^2 D_n(\theta) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d \nabla_\theta \left( 2r_{j-}(\boldsymbol{x}, \theta)\nabla_\theta r_{j-}(\boldsymbol{x}, \theta) \right) - \nabla_\theta \left( 2\nabla_\theta r_{j-}(\boldsymbol{x}^{j+}, \theta) \right)$$

$$= \frac{2}{n} \sum_{i=1}^n \sum_{j=1}^d \nabla_\theta r_{j-}(\boldsymbol{x}, \theta) \otimes \nabla_\theta r_{j-}(\boldsymbol{x}, \theta) + r_{j-}(\boldsymbol{x}, \theta)\nabla_\theta^2 r_{j-}(\boldsymbol{x}, \theta) - \nabla_\theta^2 r_{j-}(\boldsymbol{x}^{j+}, \theta)$$

$$= \frac{2}{n} \sum_{i=1}^n \sum_{j=1}^d 2T_{j-}(\boldsymbol{x}) \otimes T_{j-}(\boldsymbol{x}) \, \exp(\theta \cdot T_{j-}(\boldsymbol{x}) + b_{j-}(\boldsymbol{x}))^2$$

$$- T_{j+}(\boldsymbol{x}) \otimes T_{j+}(\boldsymbol{x}) \, \exp(\theta \cdot T_{j+}(\boldsymbol{x}) + b_{j+}(\boldsymbol{x})).$$

Plugging these derivatives $\nabla_\theta D_n(\theta)$ and $\nabla_\theta^2 D_n(\theta)$ and a given $\nabla_\theta \log \pi(\theta)$ in (10), the optimal $\beta$ is computed.

### E.3   Assumption 2 for Poisson, Ising, and Conway-Maxwell-Poisson Models

Assumption 2 for the Poisson and Ising models used in the experiments can be verified as a special case of Example 1. Any Poisson model can be written in the form

$$p_\theta(x) \propto \exp\left( \log(\theta_1) \, x - \sum_{k=1}^x \log(k) \right).$$

This falls into a class of exponential family in Example 1 by setting $\eta(\theta) = \log(\theta)$, $T(x) = x$, and $b(x) = -\sum_{k=1}^x \log(k)$. This gives that $T(x-1) - T(x) = -1$ and $b(x-1) - b(x) = \log(x)$. The condition in Example 1 is satisfied provided that $\mathbb{E}_{X\sim p}[\exp(\log(X))^2] = \mathbb{E}_{X\sim p}[X^2] < \infty$, i.e. $p$ has a second moment. Similarly, any Ising model can be written in the form

$$p_\theta(x) \propto \exp(\theta \cdot T(\boldsymbol{x}))$$

where $T : \mathcal{X} \to \mathbb{R}^k$ is a vector of summary statistics that define the model. For Ising models, $\mathcal{X}$ is of finite cardinality and $T(\boldsymbol{x})$ is hence bounded for any $\boldsymbol{x} \in \mathcal{X}$. The conditions in Example 1 are then automatically satisfied.

The Conway-Maxwell-Poisson model falls into a class of exponential family but it is beyond the simplified case of Example 1. Nonetheless, Assumption 2 is still verifiable. Recall that the

Conway-Maxwell-Poisson model has the form $p_\theta(x) \propto (\theta_1)^x (x!)^{-\theta_2}$ whose ratio function is given by $r_{j-}(x, \theta) = p_\theta(x-1)/p_\theta(x) = x^{\theta_2}/\theta_1$ where $\theta_1, \theta_2 \in (0, \infty)$. The derivative of the ratio with respect to $\theta = (\theta_1, \theta_2)$ is then given by

$$\nabla_\theta r_{j-}(x+1, \theta) = \left( -\frac{(x+1)^{\theta_2}}{\theta_1^2}, \frac{(x+1)^{\theta_2} \log(x+1)}{\theta_1} \right), \quad \nabla_\theta (r_{j-}(x, \theta))^2 = \left( -\frac{x^{2\theta_2}}{\theta_1^3}, \frac{x^{2\theta_2} \log x}{\theta_1^2} \right).$$

Note that the term $x^{2\theta_2} \log x$ in $\nabla_\theta (r_{j-}(x, \theta))^2$ is well-defined even at $x = 0$ since it converges to 0 as $x \to 0$ if $\theta_2 > 0$ despite the individual term $\log x$ alone is not well-defined for $x = 0$. Let $M_1$ and $M_2$ be the infimum value of $\theta_1$ and the supremum value of $\theta_2$ for $(\theta_1, \theta_2)$ in the bounded set $U$ to see that

$$\sup_{\theta \in U} \|\nabla_\theta r_{j-}(x+1, \theta)\| = \left| \frac{(x+1)^{M_2}}{M_1^2} \right| + \left| \frac{(x+1)^{M_2} \log(x+1)}{M_1} \right|,$$

$$\sup_{\theta \in U} \|\nabla_\theta (r_{j-}(x, \theta))^2\| = \left| \frac{x^{2M_2}}{M_1^3} \right| + \left| \frac{x^{2M_2} \log x}{M_1^2} \right|.$$

We can derive the same quantity up to constants in the power exponent of each term for the second and third derivative. Then Assumption 2 imposes that expectations of these quantities with respect to the data generating distribution $x \sim p$ are finite. For example, the expectations for the first derivatives are

$$\mathbb{E}_{X \sim p} \left[ \sup_{\theta \in U} \|\nabla_\theta r_{j-}(X+1, \theta)\| \right] = \frac{1}{M_1^2} \mathbb{E}_{X \sim p} \left[ |(x+1)^{M_2}| \right] + \frac{1}{M_1} \mathbb{E}_{X \sim p} \left[ |(x+1)^{M_2} \log(x+1)| \right],$$

$$\mathbb{E}_{X \sim p} \left[ \sup_{\theta \in U} \|\nabla_\theta (r_{j-}(x, \theta))^2\| \right] = \frac{1}{M_1^3} \mathbb{E}_{X \sim p} \left[ |x^{2M_2}| \right] + \frac{1}{M_1^2} \mathbb{E}_{X \sim p} \left[ |x^{2M_2} \log x| \right],$$

where the boundedness is translated into the moment condition of $p$ as above.

# F   Details of Experimental Assessment

This appendix contains full details for the experiments that were reported in the main text.

## F.1   Conway–Maxwell–Poisson Model

### F.1.1   Settings for KSD-Bayes

KSD-Bayes is a generalised posterior constructed by taking a kernel Stein discrepancy as a loss function; see [Matsubara et al., 2022]. The approach requires us to specify a kernel function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$, based on which the kernel Stein discrepancy is constructed. In these experiments, we adopted a kernel recommended by Yang et al. [2018] for the kernel Stein discrepancy in discrete domains $\mathcal{X}$ given by

$$k(\boldsymbol{x}, \boldsymbol{x}') = \exp \left( -\frac{1}{d} \sum_{i=1}^{d} \mathbb{1}(x_i = x_i') \right)$$

where $\mathbb{1}$ is an indicator function, taking values in $\{0, 1\}$. The effect of kernel choice is difficult to predict in the discrete context; for example, Yang et al. [2018] found that the closely related kernel

$k(\boldsymbol{x}, \boldsymbol{x}') = \sum_{i=1}^{d} \mathbb{1}(x_i = x_i')$, can perform poorly in moderate-to-high dimensions $d$ when employed in a Stein discrepancy. General principles for kernel choice in the discrete setting have not yet been established. Thus, one of the advantages of DFD-Bayes is absence of any user-specified parameters of the method.

### F.1.2 Markov Chain Monte Carlo

A Metropolis–Hasting algorithm was employed to sample from the standard Bayesian posterior, as well as KSD-Bayes and DFD-Bayes. For computational convenience, the parametrisation $\tilde{\theta}_1 = \log(\theta_1)$ and $\tilde{\theta}_2 = \log(\theta_2)$ was applied so that parameters are defined on an unbounded domain $\tilde{\theta} = (\tilde{\theta}_1, \tilde{\theta}_2) \in \mathbb{R}^2$. An isotropic Gaussian random walk proposal with covariance $\sigma^2 I$ was employed, with $\sigma = 0.1$ used for all experiments. The convergence of the Markov chain was diagnosed using univariate Gelman–Rubin statistics for each $\theta_1$ and $\theta_2$ computed from 10 independent chains. In total, 500 samples were obtained from each chain by thinning 5,000 samples, all after an initial burn-in of length 5,000. In all cases, the univariate Gelman–Rubin statistics were below 1.02, respectively for $\theta_1$ and $\theta_2$.

### F.1.3 Sales Dataset of Shmueli et al. [2005]

This dataset consists of quarterly sales figures for a particular item of clothing, taken across the different stores of a large national retailer. The original dataset is publicly available at `https://www.stat.cmu.edu/COM-Poisson/Sales-data.html`; see Shmueli et al. [2005]. Quarterly sales at each store can be small and result in a large proportion of 0 entries in the dataset, so that the Conway–Maxwell–Poisson model has a clear advantage against the standard Poisson model.

To obtain a maximum *a posteriori* estimate for the parameters of the Conway–Maxwell–Poisson model for this sales dataset, Shmueli et al. [2005] considered a prior $\pi$ defined by

$$\pi(\theta) \propto \theta_1^{a-1} \exp(-b\theta_2) \left( \sum_{j=1}^{\infty} \theta_1^j / (j!)^{\theta_2} \right)^{-c} \kappa(a, b, c) \tag{40}$$

where $(a, b, c)$ is the hyper-parameter and $\kappa(a, b, c)$ is the normalising constant of $\pi$. The motivation to use this prior is conjugacy, since the resulting posterior takes the same form as (40). However, the prior itself contains the intractable terms $(\sum_{j=1}^{\infty} \theta_1^j / (j!)^{\theta_2})^{-c}$ and $\kappa(a, b, c)$. To avoid this additional intractability, which is not a focus of the present work, we considered a simpler chi-squared prior distribution in the main text.

## F.2 Ising Model

### F.2.1 Simulating Data from the Ising Model

Samples from the Ising model were obtained using the same Metropolis–Hasting algorithm used in Yang et al. [2018]. First, all coordinates $x_i$ of $\boldsymbol{x}$ were randomly initialised to either $-1$ or $1$ with equiprobability $1/2$. Then, at each iteration, we randomly select one coordinate $x_i$ of $\boldsymbol{x}$ and flip the value of $x_i$ either from $-1$ to $1$ or from $1$ to $-1$, where the flipped value $\tilde{x}_i$ is accepted with probability $\min(1, \exp(-2\tilde{x}_i \sum_{j \in \mathcal{N}_i} x_j / \theta))$ and otherwise rejected. For the experiments in this paper we ran $n = 1,000$ chains in parallel, in each case taking the final state at iteration $100,000$.

This algorithm was used due to its implementational simplicity, rather than its efficiency, and we note that more sophisticated Markov chain Monte Carlo algorithms are available [e.g. Elçi et al., 2018].

### F.2.2   Settings for KSD-Bayes

The same choice of kernel as Appendix F.1.1 is used.

### F.2.3   Markov Chain Monte Carlo

The same Metropolis–Hasting algorithm as Appendix F.1.2 was used, in this case in dimension $p = 1$ with proposal standard deviation $\sigma = 0.1$. The convergence of the Markov chain was again diagnosed using univariate Gelman–Rubin statistics computed from 10 independent chains. In total, 100 samples were obtained after thinning from 2000 samples, with an initial burn-in of length 2000. In all cases, the univariate Gelman–Rubin statistics were below 1.002.

## F.3   Multivariate Count Data

### F.3.1   Description of the Dataset

The original data were gathered by the Cancer Genome Atlas Program, run by the National Cancer Institute in the United States, who have built large-scale genomic profiles of cancer patients with the aim to discover the genetic substructures of cancer [Wan et al., 2015]. It contains molecular profiles of biological samples of more than 30 cancer types e.g. measured via RNA sequencing technology. The raw data were pre-processed using the TCGA2STAT software developed by Wan et al. [2015]. Inouye et al. [2017] studied a subset of these data relevant to breast cancer, consisting of a total count of each gene profile found in biological samples. They applied a "log-count" transform, a common preprocessing technique for RNA sequencing data, for every datum, that is a floor function of a log transformed value of the datum. Gene profiles were then sorted by variance of the counts in descending order, with the top 10 gene profiles constituting the final dataset. The preprocessed data studied in Inouye et al. [2017] can be found in
`https://github.com/davidinouye/sqr-graphical-models`.

### F.3.2   Markov Chain Monte Carlo

The Metropolis-Hasting Markov Chain Monte Carlo was applied for this experiment. The detail for the Conway–Maxwell–Poisson graphical model is described first as the Poisson graphical model is the special case. For computational convenience, we work with the square of the interaction and dispersion parameters, i.e. $\tilde{\theta}_{i,j} := \theta_{i,j}^2$ and $\tilde{\theta}_{0,i} = \theta_{0,i}$, which modify the model as

$$p_\theta(\boldsymbol{x}) \propto \exp\left(\sum_{i=1}^d \theta_i x_i - \sum_{i=1}^d \sum_{j \in \mathcal{M}_i} \tilde{\theta}_{i,j}^2 x_i x_j - \sum_{i=1}^d \tilde{\theta}_{0,i}^2 \log(x_i!)\right)$$

The domain of each original parameter $\theta_{i,j}$ and $\theta_{0,j}$ is $[0, \infty)$. With this modification, $\tilde{\theta}_{i,j}$ and $\tilde{\theta}_{0,i}$ can be extended to $\mathbb{R}$, making the model $p_\theta(\boldsymbol{x})$ differentiable with respect to $\theta \in \mathbb{R}^v$. The derivatives of the corresponding DFD-Bayes posterior is then available to implement an efficient gradient-based Markov chain Monte Carlo method. We place a standard normal distribution as a

prior on each $\theta_i$, a normal distribution with mean 0 and scale $(d(d-1)/2)^{-1}$ as a prior on each $\tilde{\theta}_{i,j}$, and a standard normal distribution as a prior on each $\tilde{\theta}_{0,i}$, that corresponds to the original priors of each $\theta_i$, $\theta_{i,j}$, and $\theta_{0,j}$. The small scale of the half normal distribution prior on $\tilde{\theta}_{i,j}$ was chosen to suppress rapid increase of the quadratic term $x_i x_j$ as opposed to the linear term $x_i$ in the first summation. After the Markov chain finished, the absolute value was taken for the sampled values of $\tilde{\theta}_{i,j}$ and $\tilde{\theta}_{0,i}$ to convert them as the original parameters $\theta_{i,j}$ and $\theta_{0,j}$. The same setting is applied for the Poisson graphical model by fixing the dispersion parameter $\theta_{0,i} = \theta_{0,i} = 1$.

A No-U-Turn Sampler was used to approximate the DFD-Bayes posterior of both the models. In total, 100 points were obtained thinning from $5,000$ samples, with an initial burn-in of length $5,000$. The posterior predictive of each model $p_\theta(\boldsymbol{x})$ was computed by generating $500,000$ samples from $p_\theta(\boldsymbol{x})$ at every $\theta$ sampled from the DFD-Bayes posterior. Each $500,000$ predictive samples were thinned to 878 points to make it comparable with the original data of $n = 878$. The number of bootstrap minimisers $B$ used to calibrate $\beta$ for this experiment was $B = 100$.

### F.3.3  Gradient of the Discrete Fisher Divergence

For a model $p_\theta(\boldsymbol{x})$, denote the normalisation constant by $C(\theta)$ and the non-normalised part by $q_\theta(\boldsymbol{x})$, so that $p_\theta(\boldsymbol{x}) = q_\theta(\boldsymbol{x})/C(\theta)$. The discrete Fisher divergence is differentiable whenever the non-normalised part $q_\theta(\boldsymbol{x})$ is differentiable with respect to $\theta$ at any $\boldsymbol{x} \in \mathcal{X}$. Indeed, the discrete Fisher divergence between a model $p_\theta$ and data $p_n$ is given by

$$
\mathrm{DFD}(p_\theta \| p_n) \overset{\theta}{=} \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{d} \left( \frac{p_\theta(\boldsymbol{x}_i^{j-})}{p_\theta(\boldsymbol{x}_i)} \right)^2 - 2 \left( \frac{p_\theta(\boldsymbol{x}_i)}{p_\theta(\boldsymbol{x}_i^{j+})} \right)
$$

$$
\overset{\theta}{=} \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{d} \left( \frac{q_\theta(\boldsymbol{x}_i^{j-})}{q_\theta(\boldsymbol{x}_i)} \right)^2 - 2 \left( \frac{q_\theta(\boldsymbol{x}_i)}{q_\theta(\boldsymbol{x}_i^{j+})} \right)
$$

where the $\theta$-independent term is ignored and the equality holds because the normalising constant $C(\theta)$ is cancelled out. By routine calculation, the gradient of the discrete Fisher divergence is further given by

$$
\nabla_\theta \mathrm{DFD}(p_\theta \| p_n) = \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{d} 2 \left( \frac{q_\theta(\boldsymbol{x}_i^{j-})}{q_\theta(\boldsymbol{x}_i)} \right) \left( \frac{\nabla_\theta q_\theta(\boldsymbol{x}_i^{j-}) q_\theta(\boldsymbol{x}_i) - q_\theta(\boldsymbol{x}_i^{j-}) \nabla_\theta q_\theta(\boldsymbol{x}_i)}{q_\theta(\boldsymbol{x}_i)^2} \right)
$$

$$
- 2 \left( \frac{\nabla_\theta q_\theta(\boldsymbol{x}_i) q_\theta(\boldsymbol{x}_i^{j+}) - q_\theta(\boldsymbol{x}_i) \nabla_\theta q_\theta(\boldsymbol{x}_i^{j+})}{q_\theta(\boldsymbol{x}_i^{j+})^2} \right).
$$

Therefore, the gradient of the discrete Fisher divergence is well-defined as long as $q_\theta(\boldsymbol{x}_i) \neq 0$ and $q_\theta(\boldsymbol{x}_i^{j+}) \neq 0$, which in any case are prerequisites for computation of the discrete Fisher divergence.