# Discriminative Calibration: Check Bayesian Computation from Simulations and Flexible Classifier

Yuling Yao Flatiron Institute, New York, NY 10010. yyao@yyao.dev

Justin Domke
University of Massachusetts,
Amherst, MA 01002.
domke@cs.umass.edu

#### **Abstract**

To check the accuracy of Bayesian computations, it is common to use rank-based simulation-based calibration (SBC). However, SBC has drawbacks: The test statistic is somewhat ad-hoc, interactions are difficult to examine, multiple testing is a challenge, and the resulting p-value is not a divergence metric. We propose to replace the marginal rank test with a flexible classification approach that learns test statistics from data. This measure typically has a higher statistical power than the SBC test and returns an interpretable divergence measure of miscalibration, computed from classification accuracy. This approach can be used with different data generating processes to address simulation-based inference or traditional inference methods like Markov chain Monte Carlo or variational inference. We illustrate an automated implementation using neural networks and statistically-inspired features, and validate the method with numerical and real data experiments.

# 1 Introduction

Simulation based calibration (SBC) is a default approach to diagnose Bayesian computation. SBC was originally designed to validate if computer software accurately draws samples from the exact posterior inference, such as Markov chain Monte Carlo (MCMC, [4, 38, 26]) and variational inference [41]. With recent advances in amortized and simulation-based inferences [6] and growing doubt on the sampling quality [23, 18], there has also been an increasing trend to apply SBC to likelihood-free inference such as approximate Bayesian computation [42] and normalizing-flow-based [31, 1] neural posterior sampling [22, 20], with a wide range of applications in science [15, 7, 34].

Bayesian computation tries to sample from the posterior distribution  $p(\theta|y)$  given data Y. We work with the general setting where it may or may not be possible to evaluate the likelihood  $p(y|\theta)$ . Suppose we have an inference algorithm or software  $q(\theta|y)$  that attempts to approximate  $p(\theta|y)$ , and we would like to assess if this q is calibrated, meaning if  $q(\theta|y) = p(\theta|y)$  for all possible  $\theta$ , Y. Simulation based calibration involves three steps: First, we draw a  $\theta$  from the prior distribution  $p(\theta)$ . Second, we simulate a synthetic observation Y from the data model  $p(y|\theta)$ . Third, given Y we draw a size-M posterior sample  $\theta_1, \ldots, \theta_M$  from the inference engine  $q(\theta|y)$  that we need to diagnose. SBC traditionally computes the rank statistic of the prior draw  $\theta$  among the q samples, i.e. r = M = 1 ( $\theta \le \theta_m$ ). If the inference q is calibrated, then given y both  $\theta$  and  $\theta$  are from the same distribution  $p(\theta|y)$ , hence with repeated simulations of  $(\theta, y)$ , we should expect such rank statistics t to appear uniform, which can be checked by a histogram visualization or a formal uniformity test.

Despite its popularity, rank-based SBC has limitations: (i) We only compute the rank of univariate parameters. In practice,  $\theta$  and Y are high dimensional. We typically run SBC on each component of  $\theta$  separately, this creates many marginal histograms and does not diagnose the joint distribution or interactions. We may compare ranks of some one-dimensional test statistics, but there is no method to

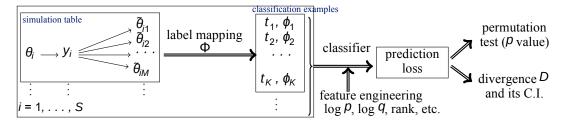


Figure 1: Our discriminate calibration framework has three modules (a) generate simulation table  $(\theta, y, \tilde{\theta})$ , and map it into classification examples with some label (t) and feature (t), (b) train a classifier to predict labels, (c) from the learned classifier, perform hypothesis testing and estimate a divergence.

find the best summary test statistic. (ii) As long as we test multiple components of  $\theta$  or test statistics, directly computing the uniformity  $\rho$ -values is invalid (subject to false discovery) unless we make a *multiple-testing* adjustment, which drops the test power (subject to false negative) in high dimensions. (iii) Often we know inference is approximate, so the final goal of diagnostic is not to reject the null hypothesis of perfect calibration but to measure the *degree of miscalibration*. The  $\rho$ -value is not such a measure: neither can we conclude an inference with  $\rho = .02$  is better than  $\rho = .01$ , nor connect the  $\rho$ -value with the posterior inference error. The evidence lower bound, a quantity common in variational inference, also does not directly measure the divergence due to the unknown entropy.

**Heuristic:** calibration via classification. To address all these drawbacks, while maintaining versatility in various computing tasks, we propose *discriminative calibration*, a pragmatic and unified framework for Bayesian computing calibration and divergence estimation. The intuitive heuristic behind discriminative calibration is to use a classifier to perform similarity tests—when two or several distributions are similar, we cannot distinguish samples from them so the classification error is large. In Bayesian computation, we compare conditional distributions, the true  $p(\theta|y)$  and inferred  $q(\theta|y)$ , to which we have access via simulations, and in some occasions explicit densities. It is natural to try to classify the samples drawn from p v.s. from q, and the ability to distinguish them would suggest a miscalibration between p and q.

To formulate such heuristics into an algorithm, we design a family of "label mapping" that prepares the simulations into classification examples that contain the label and feature. In Sec. 2, we first give four concrete label mapping examples, where the rank-based SBC becomes essentially a special case. Sec. 3 states the general theory: if we train a classifier to predict the labels, then the prediction ability yields a computable divergence from  $p(\theta|y)$  to  $q(\theta|y)$ . In Sec. 4, we illustrate the practical implementation to get the divergence estimate, its confidence interval, and a valid hypothesis testing p-value. We explain why the learned classifier helps statistical power. In Sec. 4.1, we discuss the classifier design and incorporate extra known information such as ranks, likelihood, and log densities, whenever available, as features. We also show our method is applicable to MCMC without waste from thinning. We illustrate numerical and cosmology data examples in Sec. 5. We review other related posterior validation approaches and discuss the limitation and future direction in Sec. 6.

## 2 Generate labels, run a classifier, and obtain a divergence

As with traditional SBC, we generate a simulation table by repeatedly sampling parameters  $\theta$  and synthetic data y from the target model p, i.e. draws  $(\theta, y) \sim p(\theta, y)$ . Then, for each  $(\theta, y)$ , we run the inference routine q to obtain a set of p IID approximate posterior samples p and p approximate posterior samples p and p are p and p and p are to obtain a set of p IID approximate posterior samples p and p are to obtain a set of p and p are the following samples and the classification example-sets can be created in several ways and these produce different divergences. We generalize the four examples and the claims on divergence in Sec. 3.

**Example 1: Binary classification with full features.** An intuitive way to estimate this closeness between  $q(\theta|y)$  and  $p(\theta|y)$  is to train a binary classifier.

Imagine creating a binary classification dataset of  $(t, \phi)$  pairs, where t is a binary label, and  $\phi$  are features. For each  $(\theta, y)$  simulated from P, M+1 pairs of examples are created. In the first, t=0 and  $\phi=(\theta, y)$ . In the others, t=1, and  $\phi=(\theta_m, y)$ ,  $1 \le m \le M$ . Collecting all data across

Now imagine drawing a validation set in the same way, and evaluating the log predictive density of the learned classifier. We will show below (Thm. 1) that the expected log predicted density [10] ELPD = E log  $Pr(t|\phi)$  of the classifier on validation data is a lower bound to a divergence  $D_1$  between  $p(\theta|y)$  and  $q(\theta|y)$  up to the known constant h(w),

$$\mathsf{ELPD}_1 - h(w) \le D_{-1}(p, q) := \mathsf{wKL} \left( p(\theta|y) \mid r(\theta|y) \right) + (1 - \mathsf{w}) \mathsf{KL} \left( q(\theta|y) \mid r(\theta|y) \right), \tag{1}$$

where w = 1/(M+1) and  $r(\theta|y) = wp(\theta|y) + (1-w)q(\theta|y)$  is a mixture of posterior density. If the classifier C is optimal ( $c(t|\phi) = Pr(t|\phi)$  in the distribution) then the bound in the above equation is tight  $\max_{\text{classifiers}} \text{ELPD}_1 - h(w) = D_{-1}(p, q)$ . Here  $\text{KL}(p(\theta|y) | q(\theta|y))$  denotes a standard conditional Kullback-Leibler divergence. By optimizing the classifier,  $\max_{\text{classifier}} \text{ELPD}_1 - h(w)$  becomes a commutable divergence, and its approaching zero is a necessary and sufficient condition for perfect calibration since  $D_1(p, q) = 0$  if and only if  $p(\theta|y) = q(\theta|y)$  almost everywhere.

**Example 2: Binary classification without** Y. Similar to Example 1, from each simulation draw we generate M+1 pairs of  $(t, \phi)$  examples, except that the feature  $\phi$  only contains the parameters  $\theta$ , not Y. A binary classifier is then trained to predict t given  $\phi$ . The ELPD of this classifier on validation data is a lower bound to a generalized divergence  $D_2$  between the prior  $p(\theta)$  and  $q(\theta)$ , up to a known constant h(w)

$$\begin{array}{ccccc} h & t & \varphi \\ e & 0 & \theta \\ t & 1 & \tilde{\theta}_1 \\ ed & 1 & \cdot \\ & 1 & \tilde{\theta}_M \end{array}$$

$$\mathsf{ELPD}_2 - h(w) \le D_{-2}(p, q) := \mathsf{wKL} \left( p(\theta) \mid r(\theta) \right) + (1 - \mathsf{w}) \mathsf{KL} \left( q(\theta) \mid r(\theta) \right), \tag{2}$$

where  $w = (M+1)^{-1}$ ,  $r(\theta) = wp(\theta) + (1-w)q(\theta)$  is the prior mixture, and the bound is tight when the classifier is optimal. A large ELPD reveals the difference between the inference  $q(\theta|y)$  and  $p(\theta|y)$ . But  $D_2$  is only a generalized divergence:  $D_2 = 0$  is necessary not sufficient for  $q(\theta|y) = p(\theta|y)$ .

**Example 3: Binary classification with ranks (where the classical SBC is a special case).** Instead of classification using full  $(\theta, y)$ , we construct a *feature*: the rank statistics. From each simulation draw we generate M+1 pairs of  $(t, \phi)$  examples. The first pair is t=0 and  $\phi = {M \choose m=1} 1(\theta < \tilde{\theta}_m)$ , the rank statistics of the prior draw. In the others, t=1 and  $\phi = {M \choose m'=1} 1(\tilde{\theta}_m < \tilde{\theta}_m') + 1(\tilde{\theta}_m < \theta)$ ,  $1 \le m' \le M$  are the rank statistics of the inferred samples. A binary

Out	<b>put</b> : $M + 1$ examples:
t	Ρ.,
0	$P_{M^{m=1}}^{M} 1(\theta \leq \theta_{m})$ $P_{M'=1}^{M} 1(\tilde{\theta}_{1} \leq \tilde{\theta}_{m'})$
1	$m'_{m'=1}$ $1(\tilde{\theta}_1 \leq \tilde{\theta}_{m'})$
1	Р
_1	$ \mathring{\theta}_{m'} _{=1} 1(\mathring{\theta}_{M} \leq \mathring{\theta}_{m'})$

classifier is then trained to predict t given  $\phi$ . The ELPD of this classifier is a lower bound to a generalized divergence  $D_3$  between  $p(\theta|y)$  and  $q(\theta|y)$  up to a known constant h(w)

ELPD<sub>3</sub> – 
$$h(w) \le D_{-3}(p, q) := D_{-2}(Z(p, q) \parallel \text{Uniform}(0, 1)), \quad w = 1/(M+1),$$
 (3)

and again the bound is tight if the classifier is optimal. Here Z(p, q) is a random variable defined by  $Z = Q(\theta|y)$ ,  $(\theta, y) \sim p(\theta, y)$ , where Q is the cumulative distribution function of  $q(\theta|y)$ .

Training a "generative" classifier on this rank-based label-generating map is similar to testing for uniformity in rank statistics, as done in traditional SBC which estimates the distribution of r|t=0 by histograms (See Appendix A.1 for precise correspondence between SBC and the naive Bayes classifier). The success of SBC suggests the usefulness of ranks, which motivates us to include ranks or more generally feature engineering in the classifier. However,  $D_3$  is only a generalized divergence:  $D_3 = 0$  is necessary but not sufficient for  $p(\theta|y) = q(\theta|y)$ . If inference always returns the prior,  $q(\theta|y) = p(\theta)$ , then  $D_3(p, q) = 0$ , a known counterexample of when rank-based SBC fails [32].

**Example 4: Multi-class classification.** We go beyond binary labeling. Given a simulation run  $(y, \theta, \tilde{\theta}_1, \dots, \tilde{\theta}_M)$ , we now create an (M+1)-class classification dataset with M+1 pairs of  $(t, \phi)$  examples. In each one, the features are  $\phi = (y, \theta^*)$  where  $\theta^*$  is a permutation of  $(\theta, \tilde{\theta}_1, \dots, \tilde{\theta}_M)$  that moves  $\theta$  into a given location and  $t \in [0, \dots, M]$  indicates the location of  $\theta$  in the permutation

Outp	ut: $M+1$	examples:
t	4	)
0	$(\theta, \theta_1, \theta_2, \cdot)$ $(\tilde{\theta}_1, \theta, \tilde{\theta}_2 \cdot)$	· · · $p_M$ , $y$ )
1	$(\tilde{\theta}_1, \theta, \tilde{\theta}_2)$	$\cdot\cdot,\tilde{\theta}_{M}, y)$
•		
М	$(\theta_1, \theta_2 \cdots$	$\tilde{\theta}_{M}$ , $\theta$ , $y$ )

<sup>&</sup>lt;sup>1</sup>Standard notation for conditional divergences [5] is that  $KL(p(\theta|y)||q(\theta|y)) := E_{p(y,\theta)} \log \frac{p(\theta|y)}{q(\theta|y)}$ . Conditional divergence is *not* the divergence of conditional distributions. We interprep (y) = q(y).

(See the table on the right). We train a (M+1)-class classifier to predict t from  $\phi$ . The ELPD on a validation set is a lower bound to the following divergence  $D_4$  between  $p(\theta|y)$  and  $q(\theta|y)$  up to a known constant:

$$\mathsf{ELPD}_{4} + \mathsf{log}(M+1) \leq D_{4}(p, q) := \mathsf{KL}^{(p)} p(\theta_{0}) \prod_{k=1}^{q} q(\theta_{k}), \quad \frac{1}{M+1} \prod_{m=0}^{m} p(\theta_{m}) \prod_{k=m}^{q} q(\theta_{k})^{(p)} \cdot (4)$$

Again the bound is tight if the classifier is optimal. The divergence  $0 \le D_4 \le \log(M+1)$ , and  $D_4 = 0$  if and only if  $p(\theta|y) \stackrel{\text{a.e.}}{=} q(\theta|y)$ , necessary and sufficient for calibration. In Theorem 3 we shows that as  $M \to \infty$ ,  $D_4(p, q)$  converges to  $KL(p(\theta|y), q(\theta|y))$  at an O(1/M) convergence rate.

## **3** Theory on calibration divergence

To generalize the previous examples, we define a "label mapping"  $\Phi: (y, \theta, \theta_1, \dots, \theta_M)$   $7 \rightarrow \{(t_1, \phi_1), \dots, (t_t, \phi_L)\}$ . that maps one simulation run $(y, \theta, \theta_1, \dots, \theta_M)$  into a K-class classification example-set containing L pairs of labels t and features  $\Phi$ . The label  $t_1 \in \{0, 1, \dots, K-1\}$  is deterministic and only depends on t. The features  $\Phi_t$  can depend on t and t are t and t and t and t and t are t and t and t and t are t and t and t are t and t and t are t are t and t are t are t and t are t and t are t are t and t are t are t and t are t and t are t are t and t are t are t and t are t and t are t are t and t are t are t are t and t are t and t are t are t are t and t are t are t are t and t are t and t are t are t and t are t and t are t and t are t are

We train a classifier on the classification examples created by  $\Phi$  collected across repeated sampling. The classifier performance is measured by its expected log predictive density (ELPD), or equivalently the negative cross-entropy loss. Given P, Q, and the mapping  $\Phi$ , let  $c(\phi)$  be any probabilistic classifier that uses  $\phi$  to predict the label t, and  $Pr(k|\phi,c)$  is the predicted k-th class probability. Taking expectations over features and labels with  $(\theta,y) \sim p$  and  $\theta_m \sim q(\theta|y)$  reaches the ELPD of the classifier,

$$\mathsf{ELPD}(\Phi, c) := \mathsf{E}_{t,\phi} \; \mathsf{log} \; \mathsf{Pr}(t = k | \phi, c), \quad (t, \phi) = \Phi(y, \theta, \theta_1, \dots, \theta_M)$$
 (5)

We then define the **prediction ability** $D(p, q, \Phi, c)$ to be the ELPD plus the entropy of a categorical distribution, i.e.

$$D(p, q, \Phi, c) = \text{ELPD}(\Phi, c) - \bigvee_{k=0}^{k \times 1} W_k \log w_k, \text{ where } w_k = \frac{1}{L} \bigvee_{l=1}^{k} 1(t_l = k).$$
 (6)

The optimal classifier is the C that achieves the highest prediction ability in the population:

$$D^{\text{opt}}(p, q, \Phi) := \max_{c \in C} D(p, q, \Phi, c)$$
, where  $C$  is the set of all probabilistic classifiers. (7)

The next theorem is the basic theory of our method: as long as we pass the simulation draws to a label mapping  $\Phi \in F$ , and train a classifier on the classification examples, then  $D^{\text{opt}}(p, q, \Phi)$  is a generalized divergence between P and Q.

**Theorem 1** (Prediction ability yields divergence). Given any P, Q and feature mapping  $\Phi \in F$ , the optimal prediction ability  $D^{\text{opt}}(p, q, \Phi)$  is a generalized divergence from P to Q in the sense that  $D^{\text{opt}}(p, q, \Phi) \ge 0$ , and  $p(\theta|y) = q(\theta|y)$  almost everywhere implies  $D^{\text{opt}}(p, q, \Phi) = 0$ . This generalized divergence is reparametrization invariant and uniformly bounded. For any classifier P,

$$0 \le D(p, q, \Phi, c) \le D^{\text{opt}}(p, q, \Phi) \le - \sum_{k=0}^{K \times 1} W_k \log w_k;$$

$$p(\theta|y) \stackrel{\text{a.e.}}{=} q(\theta|y) \Rightarrow D^{\text{opt}}(p, q, \Phi) = 0.$$
(8)

That is, any label mapping  $\Phi$  produces a generalized divergence  $D^{\text{opt}}(p, q, \Phi)$  the prediction ability of the optimal classifier. The prediction ability  $D(p, q, \Phi, c)$  of any classifier C estimated on validation data is a lower bound to  $D^{\text{opt}}(p, q, \Phi)$  Further, this generalized divergence  $D^{\text{opt}}(p, q, \Phi)$  is always a proper Jensen–Shannon divergence in the projected feature-label space (Theorem 1 in the Appendix).

When  $p(\theta|y)/=q(\theta|y)$ , to increase the statistical power, we wish that the generalized divergence can be as "strong" as possible such that we can detect the miscalibration. In the four examples in Section 2, we have used  $D_1$ ,  $D_2$ ,  $D_3$ ,  $D_4$  to denote the (generalized) divergence they yield. The next theorem shows there is a deterministic domination order among these four metrics. Moreover,  $P_4$  is the largest possible classification divergence from any given simulation table.

**Theorem 2** (Strongest divergence). For any given P, Q and any  $\Phi \in F$ , (1)  $D_4 \ge D_1 \ge D_3 \ge D_2$ . (2)  $D_4$  and  $D_1$  are proper divergences. They attain 0 if and only  $i\hat{p}(\theta|y) = q(\theta|y)$  almost everywhere. They attain the corresponding upper bound in (8) if and only  $p(\theta|y)$  are  $q(\theta|y)$  are disjoint, i.e.,  $\bigcap_{A} p(\theta|y)q(\theta|y)d\theta = 0$  for any measurable set  $\widehat{A}$  and almost surely  $\widehat{Y}$ . (3) For any  $\widehat{P}$ ,  $\widehat{Y}$ ,  $\widehat{Y}$  and  $\widehat{Y}$   $\widehat{Y}$  $(\Phi \text{ can have an arbitrary example size } L \text{ and class size } K), D_4(p, q) \ge D^{\text{ opt}}(p, q, \Phi).$ 

The following result shows that the divergence  $D_4$  in Eq. (4) approaches the "mode-spanning" KL divergence in the limit that the number of posterior draws  $M \to \infty$ . This is appealing because for many inference routes, increasing the number of posterior draws is cheap. Thus,  $D_4$  provides an accessible approximation to the KL divergence that is otherwise difficult to compute from samples.

**Theorem 3** (Big M limit and rate). For any P, Q generate the simulation  $\{(y, \theta, \tilde{\theta}_i, \dots \tilde{\theta}_M)\}$  and train the multiclass-classier, then  $P_4(p, q) - \text{KL}(p(\theta|y)|| q(\theta|y))_h \to 0$ , as  $M \to \infty$ . If further  $p(\theta|y)$  and  $q(\theta|y)$  have the same support, and if  $E_{p(\theta|y)} = \frac{p(\theta|y)}{q(\theta|y)} < \infty$  for a.s. Y, then

$$D_4(p, q) = \mathsf{KL}(p(\theta|y) || q(\theta|y)) - \frac{1}{2M} X^2(q(\theta|y) || p(\theta|y)) + o(M^{-1}).$$

where  $X^2(\cdot||\cdot)$  is the conditional chi-squared divergence.

# **Practical implementation**

```
Algorithm 1: Proposed method: Discriminative calibration
```

**input**: The ability to sample from  $p(\theta, y)$  and  $q(\theta|y)$ , and a label mapping  $\Phi$ . **output**: (i) estimate of a divergence between  $p(\theta|y)$  and  $q(\theta|y)$ ; (ii) p-value for testing  $p(\theta|y) = q(\theta|y)$ . for(i = 1:S)

Sample  $(\theta_i, y_i) \sim p(\theta, y)$ , and sample  $\tilde{\theta}_{i1}$ ,  $\tilde{\theta}_{i2}$ ... $\tilde{\theta}_{iM} \sim q(\theta|y_i)$ ; Simulation table Generate a batch of L examples of  $(t, \phi) = \Phi(y_i, \theta_i, \tilde{\theta}_{i1}, \dots, \tilde{\theta}_{iM})$ ,  $0 \le t \le K - 1$ ; Dabel mapping Randomly split the LS classification examples  $(t, \phi)$  into training and validation sets (all L examples for a given i go to either training or validation);

Train a K-class classifier to predict label t on the training examples, incorporating useful features; Compute the validation log predictive density LPD val in (9), obtain an estimate of the divergence (7) and its bootstrap confidence intervals; divergence

for (b = 1 : B)

Randomly permute the label t in the validation set within each batch;

Compute LPD  $_{\text{val}}^{b}$  on the permutated validation set; Compute the calibration p-value p = 1/B  $_{b=1}^{B}$  1(LPD  $_{b}^{\text{val}} \ge \text{LPD}^{\text{val}}$ ).

Pfrequentist test

**Workflow for divergence estimate.** We repeat the simulation of  $(y, \theta, \tilde{\theta}_1, \dots, \tilde{\theta}_M)$  for S times. Each time we sample  $(y, \theta) \sim p(\theta, y)$  and  $\theta_{1:M} \sim q(\theta|y)$  and generate a *batch* of L examples through a label mapping  $\Phi \in F$ . In total, we obtain SL pairs of  $(t, \phi)$  classification examples. We recommend using the binary and multiclass labeling schemes (Examples 1 and 4, where L = M + 1 and K = 2 or M+1) such that we can obtain a proper divergence estimate. We split the classification example-set  $\{(t, \phi)\}$  into the training and validation set (do not split batches) and train a K -class classifier c on the training set to minimize cross-entropy. Denote  $I_{\text{val}}$  to be the validation index,  $\Pr(t = t_j | \phi_j, c)$  to be the learned class probability for any validation example  $(t_i, \phi_i)$ , we compute the ELPD (5) by the validation set log predictive density:

ELPD(
$$\Phi$$
,  $c$ )  $\approx$  LPD <sup>val</sup> ( $\Phi$ ,  $c$ )  $\coloneqq$   $|I|_{\text{val}} |I|^{-1} \log \Pr(t = t_j | \phi_j, c).$  (9)

For any  $^{C}$ , LPD  $^{\text{val}}(\Phi, c) = \sum_{k=0}^{P} W_k \log W_k$  becomes a lower bound estimate of the (generalized) divergence  $^{D \text{ opt}}(p, q, \Phi)$  in Thm. 1, and an estimate of  $^{D \text{ opt}}(p, q, \Phi)$  itself when the classier  $^{C}$  is good enough. In addition to the point estimate of the divergence, we can compute the confidence interval. It is straightforward to obtain the standard error of the sample mean (9). To take into account the potentially heavy tail of the log predictive densities, we can also use Bayesian bootstrap [37] that reweights the sum in (9) by a uniform Dirichlet weight.

**Hypothesis testing.** Our approach facilitates rigorous frequentist hypothesis testing. The null hypothesis is that the approximate inference matches the exact posterior, i.e.,  $p(\theta|y) = q(\theta|y)$  almost everywhere. We adopt the permutation test: We train the classifier  $^{C}$  once on the training set and keep it fixed, and evaluate the validation set log predictive density LPD  $^{\text{val}}$  ( $\Phi$ , c) in (9). Next, permutate the validation set  $^{B}$  times: at time  $^{b}$ , keep the features unchanged and randomly permutate the validation labels  $^{t}$  within each batch of examples ( $\Phi$  generates a batch of  $^{L}$  examples each time), and reevaluate the validation set log predictive density (9) op permutated labels, call it LPD  $^{\text{val}}_{b}$ . Then we compute the one-sided permutation  $^{b}$ -value as  $p = {}^{b}_{b=1}$  1(LPD  $^{\text{val}}_{b} \ge \text{LPD}$   $^{\text{val}}_{b}$ )/ $^{B}$ . Given a significance level, say 0.05, we will reject the null if  $^{b}$  < 0.05 and conclude a miscalibration.

**Theorem 4** (Finite sample frequentist test). For any finite simulation size S and posterior draw size M, and any classifier C, under the null hypothesis  $p(\theta|y) = q(\theta|y)$  almost everywhere, the permutation test is valid as the P-value computed above is uniformly distributed on [0, 1]

Our test is exact when the simulation size S and M is finite, while the original SBC [4] relied on asymptotic approximation. Further, we *learn* the test statistic via the classifier, and our test is valid regardless of the dimension of  $\theta$  and there is no need to worry about post-learning testing [3], while the usual SBC rank test will suffer from low power due to multiple testing.

Our test is always valid even if the classifier C is not optimal. Why does our *learning* step help? For notional brevity, here we only reason for the binary classification. For any P, Q we apply the binary classification as described in Example 1, t = 0 or 1 and  $\phi = (\theta, y)$ .

**Theorem 5** (Sufficiency). Let  $\hat{c}(\theta, y) = \Pr(t = 1 | \theta, y)$  be the probability of label 1 in the optimal classifier as per (7), and let  $\Pi_c^p$  and  $\Pi_c^q$  be the one-dimensional distributions of this  $\hat{c}(\theta, y)$  when  $(\theta, y)$  is sampled from  $p(\theta, y)$  or from  $p(y)q(\theta|y)$  respectively, then (i) Conditional on the summary statistic  $\mathcal{C}$ , the label  $\mathcal{C}$  is independent of features  $\phi = (\theta, y)$ . (ii) Under regularity conditions, there is no loss of information in divergence as the joint divergence is the same as the projected divergence in the one-dimensional  $\mathcal{C}$ -space  $\mathcal{D}_1(p, q) = \mathcal{D}_1(\Pi_c^p, \Pi_c^q)$ .

That is, the best prediction  $^{\mathbf{c}}$  entails the best one-dimensional summary statistics of the high dimensional  $\theta \times y$  space. The enhancement of the test power from using the sufficient statistics is then assured by the Neyman-Pearson lemma [27].

## 4.1 Feature engineering: use rank, log density and symmetry

Reassuringly, whatever the classifier C, the prediction ability  $D(p, q, \Phi, c)$  is always a lower bound to the corresponding divergence  $D^{\text{opt}}(p, q, \Phi)$  and the permutation test is always exact. But we still wish to train a "good" classifier in terms of its out-of-sample performance, for a tighter bound in divergence, and a higher power in testing. We will typically use a flexible parametric family such as a multilayer perceptron (MLP) network to train the classifier.

The oracle optimal probabilistic classifier is the true label probability, and in the binary and multiclass classification (Example 1 and 4), the oracle has closed-form expressions, although we cannot evaluate:

$$\Pr_{\text{binary}}(t=0|\theta,y) = \frac{p(\theta|y)}{p(\theta|y) + q(\theta|y)}, \quad \Pr_{\text{multi}}(t|\theta_0,\dots,\theta,y) = \Pr_{\substack{M \\ k=0}} \frac{p(\theta_t,y)/q(\theta_t|y)}{p(\theta_k,y)/q(\theta_k|y)}. \quad (10)$$

**Statistically-meaningful feature.** Depending on the inference task, we have more information than just the sample points and should use them in the classifier. In light of the shape and component of the optimal classifiers (10), the following statistically-meaningful features are useful whenever available: (i) The log target density  $\log p(\theta|y)$ . As proxies, the log joint density  $\log p(\theta, y)$  and the log likelihood  $\log p(y|\theta)$  are often known to traditional Bayesian models. (ii) The log approximate density  $\log q(\theta|y)$ , known to volitional and normalizing flows. (iii) When the log approximate density is unknown, a transformation is to integrate the density and obtain the posterior CDF,  $Q(\theta|y) = E_{x-q(x|y)} \quad 1(x > \theta)$ , and this CDF can be approximated by the rank statistics among the approximate draws up to a rescaling  $p(\theta, y) := \frac{P_{MS}}{m=1, i=1} \quad 1(y_i = y) \quad 1(\theta < \tilde{\theta}_{im})$ . See Appendix Table 1 for the usage and extension of these features in binary and multiclass classifiers.

**Linear features.** We call the log likelihood, the log approximate density, or log prior (whenever available) linear features, and denote them to be  $l(\theta, y)$ . For example if both likelihood and q density

is known, then  $l(\theta, y) = (\log p(y|\theta), \log q(\theta|y))$ . Because they appear in the oracle 10, we will keep linear features in the last layer of the network (followed by a softmax). We recommend parameterizing the binary classification in the following form:

$$Pr(t = 1 | (\theta, y)) = inv \quad logit[MLP(\theta, y) + w^{T} l(\theta, y)].$$
 (11)

**Symmetry in multiclass classifier.** The multi-class classifier is harder to train. Luckily, we can use the symmetry of the oracle classifier (10): the probability of class k is proportional to a function of  $(\theta_k, y)$  only, hence we recommend parameterizing the multiclass probability as

$$\Pr(t = k | (\theta_0, \theta_1, \dots, \theta_n, y)) = \Pr \frac{\exp(g(\theta_k, y))}{\frac{M}{k' = 0} \exp(g(\theta_{k'}, y))}, \ g(\theta, y) = \text{MLP}(\theta, y) + w^T I(\theta, y), \ (12)$$

where  $I(\theta, y)$  is available linear features. We only need to learn a reduced function from  $\theta \times y$  to R, instead of from  $\theta^{M+1} \times y$  to R, reducing the complexity while still keeping the oracle (10) attainable.

**Zero waste calibration for MCMC.** If  $\tilde{\theta}_1, \dots, \tilde{\theta}_M \sim q(\theta|y)$  are produced by MCMC sampler, typically  $\tilde{\theta}$  has autocorrelations. Although classical SBC originated from MCMC application, the rank test requires independent samples, so SBC can only handle autocorrelation by thinning: to subsample  $\tilde{\theta}_m$  and hope the thinned  $\tilde{\theta}$  draws are independent. Thinning is a waste of draws, inefficient when the simulations are expensive, and the thinned samples are never truly independent. With MCMC draws  $\tilde{\theta}$ , our method could adopt thinning as our Thm. 1 and 4 are valid even when M = 1.

Yet we can do better by using all draws. The "separable" network architecture (12) is ready to use for MCMC samples. For example, we sample  $(\theta, y) \sim p(\theta, y)$ , and sample  $(\tilde{\theta}_1, \dots, \tilde{\theta}_M)$  from a MCMC sampler whose stationary distribution we believe is  $q(\theta|y)$ , and generate examples from the multiclass permutation (Example 4). Then we run a separable classifier (12) to predict t. Intuitively, the separable network design (12) avoids the interaction between  $\tilde{\theta}_m$  with  $\tilde{\theta}_m$ , and disallows the network to predict t based on the autocorrelation or clustering of  $\tilde{\theta}$ . The next theorem states the validity of our method in MCMC settings without thinning.

**Theorem 6** (MCMC). Suppose we sample  $(\theta, y) \sim p(\theta, y)$ , and sample  $(\theta_1, \dots, \theta_M)$  from a MCMC sampler whose stationary distribution we believe is  $q(\theta|y)$  (i.e., marginally  $\theta_i$  is from  $q(\theta|y)$ ), and generate examples  $((t_1, \theta_1), \dots, (t_{M+1}, \theta_{M+1}))$  from the multiclass permutation, such that  $\phi = (\theta_0, \theta_1, \dots, \theta_M)$ . Then we run an exchangeable classifier (12) in which g is any  $g \times g \to g$  mapping. Denote  $g \to g$  (p, q) to be the predictive ability of the optimal classifier among all separable classifiers (12), then  $g \to g$  (p, q) =  $g \to g$  (p, q).

**Dimension reduction and nuisance parameter.** Sometimes we only care about the sampling quality of one or a few key dimensions of the parameter space, then we only need to restrict the classifier to use these targeted dimensions, as a result of Theorem 1. For example, in binary classification, if we reduce the feature  $\phi = (\theta, y)$  to  $\phi = (h(\theta), y)$  in the classifier, where  $h(\theta)$  can be a subset of  $\theta$  dimensions, then the resulting classification divergence becomes projected divergence between  $h(\theta)|y$ ,  $\theta \sim p(\theta|y)$  and  $h(\theta)|y$ ,  $\theta \sim q(\theta|y)$ , and other nuisance parameters do not impact the diagnostics.

Weighing for imbalanced binary classification. When M is big, the binary classification (Example 1) can suffer from imbalanced labels, and the divergence in (1) degenerates:  $D_1(p, q) \rightarrow 0$  as  $M \rightarrow \infty$ . One solution is to use the multiclass classification which has a meaningful limit (Thm. 3). Another solution is to reweigh the loss function or log predictive density by the label t. If the weights of class 1 and class 0 examples are  $\frac{M+1}{2M}$  and  $\frac{M+1}{2}$ , used in both training and validation log prediction density, then regardless of M, the weighted classification is equivalent to balanced classification, and the resulting divergence is the symmetric Jensen-Shannon (JS) divergence  $\frac{1}{2}$  KL[ $p(\theta|y)$ || $r(\theta|y)$ ] +  $\frac{1}{2}$  KL[ $p(\theta|y)$ || $r(\theta|y)$ || , where  $r(\theta|y) = \frac{1}{2}[p(\theta|y) + q(\theta|y)]$ . See Appendix A for proofs.

## 5 Experiments

**Closed-form example.** Consider a multivariate normal parameter prior  $\theta \in \mathbb{R}^d \sim \mathsf{MVN}(\mathbf{0}, \mathsf{Id} \ d)$  and a normal data model  $y|\theta \sim \mathsf{MVN}(\theta, \Sigma)$ , so the exact posterior  $p(\theta|y)$  is explicit. In the experiments,

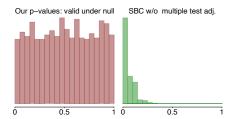


Figure 2: Our test is valid as it yields uniformly distributed P-values under the null  $p(\theta|y) = q(\theta|y)$ . We check all dimensions at once.

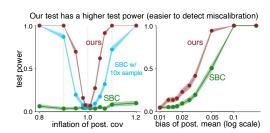


Figure 3: Our test has a uniformly higher power than SBC rank test. We simulate wrong inference 9 by multiplying the posterior covariance or adding biases to the posterior mean.

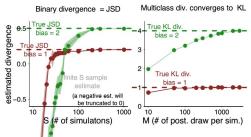


Figure 4: In the Gaussian example with bias = 1 or 2, with a moderately big simulation sample size S, the divergence estimated from a binary (left panel) or multiclass (right panel) classifier matches its theory quantity (dashed): the Jensen-Shannon or Kullback-Leibler divergence in the big M limit.

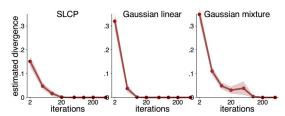


Figure 5: We apply binary classifier calibration to posterior inferences on three models from the SBI benchmark, and check the sampling quality after some iterations. The x-axis is the number of MCMC iterations. The y-axis is the estimated JS divergence (± standard error) between the true target and sampled distribution at a given number of iterations, indicating gradual convergence.

we use neural nets to train binary (11) or symmetric multiclass classifiers (12), which we generally recommend. We assume the inference posterior  $\log q(\theta|y)$  is known and added to the classifier.

Legit macy of testing. First, to validate our hypothesis test under the null, we use true inference  $q(\theta|y) = p(\theta|y)$ . With d = 16, we simulate S = 500 draws of  $(\theta, y) \sim p(\theta, y)$ , then generate true posterior  $\theta_m \sim p(\theta|y)$ , and run our binary-classifier and obtain a permutation p-value. We repeat this testing procedure 1000 times to obtain the distribution of the values under the null, which is uniform as shown in Fig. 2, in agreement with Thm. 4. Because  $\theta$  has d dimensions, a SBC rank test on each margin separately is invalid without adjustment, and would require Bonferroni corrections.

*Power.* We consider two sampling corruptions: (i) bias, where we add a scalar noise (0.01 to 0.2) to each dimension of the true posterior mean, (ii) variance, where we inflate the true posterior covariance matrix by multiplying a scalar factor (0.8 to 1.2). We compare our discriminative tests to a Bonferroni-corrected SBC chi-squared rank-test. In both settings, we fix a 5% type-I error and compute the power from 1000 repeated tests. Fig. 3 shows that our method has uniformly higher power than SBC, sometimes as good as SBC with a 10 times bigger simulation sample size.

Divergence estimate. The left panel of Fig. 4 validates Thm. 1: We run a weighted binary classifier on the Gaussian simulations with a scalar bias 1 or 2 added to the posterior mean in Q. As the number of simulations S grows, the learned classifier quickly approaches the optimal, and the prediction ability matches the theory truth (dashed line): the Jensen-Shannon divergence between  $Q(\theta|y)$  and  $Q(\theta|y)$ . The right panel validates Thm. 3: With a fixed simulation size  $Q(\theta|y)$ , we increase  $Q(\theta|y)$  and run a multiclass classification. When  $Q(\theta|y)$  is big, the estimated divergence converges from below to the dashed line, which is the theory limit  $Q(\theta|y)$ ,  $Q(\theta|y)$  in Thm. 3.

**Benchmark examples.** Next, we apply our calibration to three models from the SBI benchmark [23]: the simple likelihood complex posterior (SLCP), the Gaussian linear, and the Gaussian mixture model. In each dataset, we run adaptive No-U-Turn sampler (NUTS) and check the quality of the sampled distribution after a fixed number of iterations, varying from 2 to 2000 (we use equal number of iterations for warm-up and for sampling, and the warm-up samples were thrown away). At each given MCMC iterations, we run our classifier calibration, and estimate the JS divergence, as reported

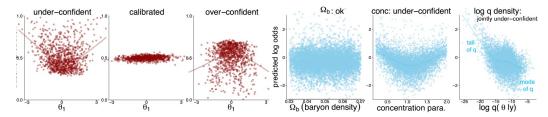


Figure 6: Scatter plots of the classifier prediction v.s. Figure 7: In the cosmology example we visually check a one-dimensional parameter we need to check can the classifier predicted log odds v.s. cosmology parametereal marginal over- or under-confidence in inference eter  $\Omega_b$  and CONC, and  $log q(\theta|y)$ . Underconfidence is  $q(\theta|y)$  from a U or inverted-U shape.

detected in CONC marginally, and more so in the joint.

in Fig. 5. In all three panels, we are able to detect the inference flaws at early iterations and observe a gradual convergence to zero.

**Visual check.** Our diagnostic outputs rigorous numerical estimates, but it also facilities visual checks. We make a scatter plot of the binary classifier prediction  $Pr(t = 1|\phi)$ , a proxy of  $p(\theta|y)/q(\theta|y)$ , against any one-dimensional parameter we need to check: If that parameter is under- or over-confident, this scatter plot will display a U- or inverted-U-shaped trend. Compared with SBC rank histograms, our visualization can further check the magnitude of mismatch (how far awayPr( $t = 1|\phi$ ) is from 0.5), tail behavior (small  $q(\theta|y)$ ), and several dimensions jointly. Fig 6 is a visual check in the Gaussian example when we multiply the posterior covariance in inference q by 0.8, 1, or 1.2.

**Galaxy clustering.** We would like to model galaxy spatial clustering observations in a Lambda

	model 1, $d_y = 38$	model 2, $d_y = 151$	model 3, $d_y = 1031$
1-flow	0.015 (0.004)	0.076 (0.007)	0.068 (0.007)
5-ens.	0.013 (0.004)	0.044 (0.006)	0.035 (0.006)

cold dark matter framework [34], where observations y correspond to statistical descriptors of galaxy clustering, and a 14-dimensional parameter  $\theta$  encodes key cosmological information about the Universe. The forward simulation  $y|\theta$  involves expensive N-body simulations (each single y simulation takes around 5000 cpu hours). We consider three different physical models P: they contain the same cosmology parameter  $\theta$  but three types of observations whose dimension  $d_y$  varies from 38 to 1031, reflecting three power-spectrum designs. We apply simulation based inference to sample from each of the models, where we try various normalizing flow architectures and return either the best architecture or the ensemble of five architectures. We now apply our discriminative calibration to these six approximate inferences using a weighted binary classifier trained from S = 2000, M = 500simulations for each inference (that is 1 million examples). We have added the log densities  $q(\theta|y)$ as extra features in the classifier, as they are known in the normalizing-flows. The table above shows the estimated Jensen-Shannon distances and standard deviation. Compared this table with the estimates from blackbox MLP classifier (Appendix Fig. 8), using statistical features greatly improves the label prediction and tightens the divergence estimate. From the table above, all inferences q have a significantly non-zero but small divergence, among which the 5-ensemble that uses a mixture of 5 flows always has a lower divergence so is more trustworthy. Further visualization of the classification log odds (Fig. 7) reveals that q is under-confident in the parameter that encodes the galaxy concentration rate. To interpret the tail behavior, the predicted log odds are negatively correlated with log joint normalizing-flow density  $q(\theta|y)$ , suggesting  $q(\theta|y) > p(\theta|y)$  in the tail of  $q(\theta|y)$ , another evidence of q underconfidence.

We share Jax implementation of our binary and multiclass classifier calibration in Github.

#### 6 Discussions

This paper develops a classifier-based diagnostic tool for Bayesian computation, applicable to MCMC, variational and simulation based inference, or even their ensembles. We learn test statistics from data using classification. Through a flexible label mapping, our method yields a (family of) computable divergence metric and an always valid testing p-value, and the statistical power is typically higher than the traditional rank based SBC. Various statically-meaningful features are available depending on the task and we include them in the classifier. Visual exploration is also supported.

<sup>&</sup>lt;sup>2</sup>https://github.com/yao-yl/DiscCalibration

**Related diagnostics.** The early idea of simulation based calibration dates back to [13] who compared the prior and the data-averaged posterior, i.e., to check  $p(\theta) = q(\theta|y)p(y)dy$  using simulations. [42] further developed techniques for the first and second moment estimate of the data-averaged posterior using the law of total variance. Our method includes such moment comparison as special cases: using the no-y binary labeling (Example 2), then comparing the empirical moments of the q sample the p sample can be achieved through a naive Bayes classifier or a linear discriminant analysis.

The rank-based SBC [4, 38] can be recovered by ours using binary labels and taking ranks to be features (Example 3). The rank statistic is central to SBC, which follows the classical frequentist approach—using the tail probability under the posited model quantifies the extremeness of the realized value[11] is only a convenient way to locate the observations in the reference distribution, especially in the past when high dimensional data analysis was not developed—it is the use of modern learning tools that sharpens our diagnostic. Recently, SBC has developed various heuristics for designing (one-dimensional) test statistics  $\phi(\theta, y)$ . Such rank tests are recovered by our method by including the rank of  $\phi$  in our features (Sec. 4.1). For example, [26] proposed to test the rank of the likelihood  $\phi = p(\theta|y)$  in MCMC, [20] looked at the rank of proposal density  $q(\theta|y)$ , [22] used the q-probability integral transformation in normalizing flows. In light of our Theorem 5, the optimal test statistic is related to the density ratio and hence problem-specific, which is why our method includes all known useful features in the classifier and learns the test statistic from data.

Classifier two-sample test. Using classification to compare two-sample closeness is not a new idea. The classifier two-sample test (C2ST) has inspired the generative adversarial network (GAN) [6] and conditional GAN [25]. [33] has developed GAN-typed inference tools for SBI. In the same vein, [19] used classifiers to diagnose multiple-run MCMC, and [24] used classifiers to learn the likelihood ratio for frequentist calibration, which is further related to using regression to learn the propensity score [36] in observational studies. The theory correspondence between binary classification loss and distribution divergence has been studied in [21, 28, 35]. This present paper not only applies this classifier-for-two-sample-test idea to amortized Bayesian inference to obtain a rigorous test, but also advances the classifier framework by developing the theory of the "label mapping", while the traditional GAN-type approach falls in the one-class-per-group category and deploy binary classifiers as in Example 1. Our extension is particularly useful when samples are overlapped (multiple  $\theta$  per y), autocorrelated, and imbalanced.

KL divergence estimate from two samples. As a byproduct, our proposed multiclass-classifier provides a consistent KL divergence estimate (Thm. 3) from two samples. Compared with existing two-sample KL estimate tools from the f-divergence representation [29] or the Donsker-Varadhan representation [2], our multiclass-classifier estimate appears versatile for it applies to samples with between-sample dependence or within-sample auto-correlation. It is plausible to apply our multiclass-classifier approach to other two-sample divergence estimation tasks such as the f-GAN [30], which we leave for future investigation.

**Limitations and future directions.** Like traditional SBC, our method assesses the difference between  $p(\theta|y)$  and  $q(\theta|y)$ , averaged over y. This is a "global" measure relevant to developing algorithms and statistical software. But sometimes the concern is how close P and Q are for some particular observation y = y obs. "Local" diagnostics have been developed for MCMC [14, 12, 40, 17], variational inference [41, 8] and simulation-based inference [42, 22] that try to assess  $q(\theta|y^{\text{obs}})$  only. It would be valuable to extend our approach to address these cases. Another future direction would be to extend our approach to posterior predictive checks [11, 9, 39] that diagnose how well the statistical model P fits the observed data.

## Acknowledgement

The authors thank Bruno Regaldo-Saint Blancard for his help on galaxy clustering data and simulations. The authors thank Andreas Buja, Andrew Gelman, and Bertrand Clarke for discussions.

### References

[1] Agrawal, A., Sheldon, D. R., and Domke, J. (2020). Advances in black-box VI: Normalizing flows, importance weighting, and optimization. In *Advances in Neural Information Processing Systems*.

- [2] Belghazi, M. I., Baratin, A., Rajeswar, S., Ozair, S., Bengio, Y., Courville, A., and Hjelm, R. D. (2018). MINE: mutual information neural estimation. In *International Conference on Machine Learning*.
- [3] Berk, R., Brown, L., Buja, A., Zhang, K., and Zhao, L. (2013). Valid post-selection inference. Annals of Statistics, 41(2):802–837.
- [4] Cook, S. R., Gelman, A., and Rubin, D. B. (2006). Validation of software for Bayesian models using posterior quantiles. *Journal of Computational and Graphical Statistics*, 15(3):675–692.
- [5] Cover, T. M. and Thomas, J. A. (1991). Elements of Information Theory. John Wiley & Sons, 2rd edition.
- [6] Cranmer, K., Brehmer, J., and Louppe, G. (2020). The frontier of simulation-based inference. Proceedings of the National Academy of Sciences, 117(48):30055–30062.
- [7] Dax, M., Green, S. R., Gair, J., Macke, J. H., Buonanno, A., and Schölkopf, B. (2021). Real-time gravitational wave science with neural posterior estimation. *Physical Review Letters*, 127(24):241103.
- [8] Domke, J. (2021). An easy to interpret diagnostic for approximate inference: Symmetric divergence over simulations. *arXiv*:2103.01030.
- [9] Gabry, J., Simpson, D., Vehtari, A., Betancourt, M., and Gelman, A. (2019). Visualization in bayesian workflow. *Journal of Royal Statistical Society: Series A*, 182(2):389–402.
- [10] Gelman, A., Hwang, J., and Vehtari, A. (2014). Understanding predictive information criteria for Bayesian models. *Statistics and Computing*, 24(6):997–1016.
- [11] Gelman, A., Meng, X.-L., and Stern, H. (1996). Posterior predictive assessment of model fitness via realized discrepancies. *Statistical Sinica*, 6:733–760.
- [12] Gelman, A. and Rubin, D. B. (1992). Inference from iterative simulation using multiple sequences. *Statistical Science*, 7(4):457–472.
- [13] Geweke, J. (2004). Getting it right: Joint distribution tests of posterior simulators. *Journal of American Statistical Association*, 99(467):799–804.
- [14] Geyer, C. J. (1992). Practical Markov chain Monte Carlo. Statistical Science, 7(4):473–483.
- [15] Gonçalves, P. J., Lueckmann, J.-M., Deistler, M., Nonnenmacher, M., Öcal, K., Bassetto, G., Chintaluri, C., Podlaski, W. F., Haddad, S. A., and Vogels, T. P. (2020)Training deep neural density estimators to identify mechanistic models of neural dynamics. *Elife*, 9:e56261.
- [16] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial networks. In *Advances in Neural Information Processing Systems*.
- [17] Gorham, J. and Mackey, L. (2017). Measuring sample quality with kernels. In *International Conference on Machine Learning*.
- [18] Hermans, J., Delaunoy, A., Rozet, F., Wehenkel, A., Begy, V., and Louppe, G. (2022) A trust crisis in simulation-based inference? Your posterior approximations can be unfaithful. *Transactions on Machine Learning Research*.
- [19] Lambert, B. and Vehtari, A. (2022). R\*: A robust MCMC convergence diagnostic with uncertainty using decision tree classifiers. *Bayesian Analysis*, 17(2):353–379.
- [20] Lemos, P., Coogan, A., Hezaveh, Y., and Perreault-Levasseur, L. (2023)Sampling-based accuracy testing of posterior estimators for general inference. *arXiv:2302.03026*.
- [21] Liese, F. and Vajda, I. (2006). On divergences and informations in statistics and information theory. *IEEE Transactions on Information Theory*, 52(10):4394–4412.
- [22] Linhart, J., Gramfort, A., and Rodrigues, P. L. (2022). Validation diagnostics for SBI algorithms based on normalizing flows. *arXiv*:2211.09602.
- [23] Lueckmann, J.-M., Boelts, J., Greenberg, D., Goncalves, P., and Macke, J. (2021). Benchmarking simulation-based inference. In *International Conference on Artificial Intelligence and Statistics*.
- [24] Masserano, L., Dorigo, T., Izbicki, R., Kuusela, M., and Lee, A. B. (2022). Simulation-based inference with WALDO: Perfectly calibrated confidence regions using any prediction or posterior estimation algorithm.
- [25] Mirza, M. and Osindero, S. (2014). Conditional generative adversarial nets. arXiv:1411.1784.

- [26] Modrák, M., Moon, A. H., Kim, S., Bürkner, P., Huurre, N., Faltejsková, K., Gelman, A., and Vehtari, A. (2022). Simulation-based calibration checking for Bayesian computation: The choice of test quantities shapes sensitivity. arXiv:2211.02383.
- [27] Neyman, J. and Pearson, E. (1936). Sufficient statistics and uniformly most powerful tests of statistical hypotheses. In Statistical Research Memoirs, pages 1–37. Cambridge University Press.
- [28] Nguyen, X., Wainwright, M. J., and Jordan, M. I. (2009). On surrogate loss functions and f -divergences. Annals of Statistics, 37(2):876–904.
- [29] Nguyen, X., Wainwright, M. J., and Jordan, M. I. (2010). Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56(11):5847–5861.
- [30] Nowozin, S., Cseke, B., and Tomioka, R. (2016). f-GAN: Training generative neural samplers using variational divergence minimization. Advances in Neural Information Processing Systems.
- [31] Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S., and Lakshminarayanan, B. (2021). Normalizing flows for probabilistic modeling and inference. *Journal of Machine Learning Research*, 22(1):2617–2680.
- [32] Prangle, D., Blum, M., Popovic, G., and Sisson, S. (2014). Diagnostic tools of approximate Bayesian computation using the coverage property. *Australian & New Zealand Journal of Statistics*, 56(4):309–329.
- [33] Ramesh, P., Lueckmann, J.-M., Boelts, J., Tejero-Cantero, Á., Greenberg, D. S., Gonçalves, P. J., and Macke, J. H. (2022). GATSBI: Generative adversarial training for simulation-based inference*International Conference on Learning Representations*.
- [34] Régaldo-Saint Blancard, B., Hahn, C., Ho, S., Hou, J., Lemos, P., Massara, E., Modi, C., Moradinezhad Dizgah, A., Parker, L., , Yao, Y., and Eickenberg, M. (2023). SimBIG: Galaxy clustering analysis with the wavelet scattering transform. *arXiv:2310.15250*.
- [35] Reid, M. and Williamson, R. (2011). Information, divergence and risk for binary experiments. *Journal of Machine Learning Research*.
- [36] Rosenbaum, P. R. and Rubin, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55.
- [37] Rubin, D. B. (1981). The Bayesian bootstrap. Annals of Statistics, 9(1):130–134.
- [38] Talts, S., Betancourt, M., Simpson, D., Vehtari, A., and Gelman, A. (2018). Validating Bayesian inference algorithms with simulation-based calibration. *arXiv*:1804.06788.
- [39] Vandegar, M., Kagan, M., Wehenkel, A., and Louppe, G. (2021). Neural empirical Bayes: Source distribution estimation and its applications to simulation-based inference. In *International Conference on Artificial Intelligence and Statistics*.
- [40] Vehtari, A., Gelman, A., Simpson, D., Carpenter, B., and Bürkner, P.-C. (2021). Rank-normalization, folding, and localization: An improved for assessing convergence of MCMC. Bayesian Analysis, 16(2):667–718.
- [41] Yao, Y., Vehtari, A., Simpson, D., and Gelman, A. (2018)Yes, but did it work?: Evaluating variational inference. In *International Conference on Machine Learning*.
- [42] Yu, X., Nott, D. J., Tran, M.-N., and Klein, N. (2021). Assessment and adjustment of approximate inference algorithms using the law of total variance. *Journal of Computational and Graphical Statistics*, 30(4):977–990.

# **Appendices to "Discriminative Calibration"**

The Appendices are organized as follows. In Section A, we provide several extra theory results for the main paper. In Section B, we discuss our experiment and implementations. In Section C, we provide proofs for the theory claims.

# A Additional theory results

**Theorem 7** (Closed form expression of the divergence). Consider the label generating process in Section 4. Let  $\pi(t, \phi)$  be the joint density of the label and features after mapping  $\Phi$ , i.e., the distribution of the  $(t, \phi)$  generated by the simulation process:

- 1. sample  $(\theta, y) \sim p(\theta, y)$ ;
- 2. sample  $(\theta_1, \ldots, \theta_M)$  from  $q(\theta|y)$ ;
- 3. generate  $((t_1, \phi_1), \ldots, (t_r, \phi_L))) = \Phi((y, \theta_1, \theta_1, \ldots, \theta_M));$
- 4. sample an index | from Uniform(1, ..., L)
- 5. return  $(t_1, \phi_1)$ .

We define

$$\pi_k(\phi) = \pi(\phi|t=k) = \sum_{y}^{Z} \pi(y)\pi(\phi|y, t=k)dy$$

to be the Y-averaged law of  $\phi | t = k$ . Note that Y has been averaged out.

Then classification divergence defended through (7) in Theorem 1 has a closed form expression, a Jensen–Shannon divergence in this projected  $\phi$  space.

$$D^{\text{opt}}(p, q, \Phi) = \bigvee_{k=0}^{k+1} W_k \mathsf{KL} \left( \left. \prod_{j=0}^{k+1} W_j \prod_{j=0}^{k+1} W_j \prod_{j=0}^{k+1} W_j \right) \right). \tag{13}$$

This Theorem 7 gives a connection between our divergence and the familiar Jensen–Shannon divergence, which is well known to be linked to IID sample classification error (note that the classification examples we generate are not IID for they share *Y*).

As a special case, when K=2,  $W_0=W_1=1/2$ , we recover the symmetric two-group Jensen–Shannon divergence, or JSD. In general, between two densities  $\pi_1(x)$  and  $\pi_2(x)$ , the JSD is

Jensen Shannon divergence 
$$(\pi \pi_2) = \frac{1}{2} \text{KL} \quad \pi_1 \parallel \frac{\pi_1 + \pi_2}{2} + \frac{1}{2} \text{KL} \quad \pi_2 \parallel \frac{\pi_1 + \pi_2}{2} \cdot (14)$$

As before, the standard notation [5] of conditional divergence is defined as taking expectations over conditional variables:

$$JSD(\pi_{1}(x|z), \pi_{2}(x|z)) = \frac{1}{2} KL \quad \pi_{1}(x|z) \mid \mid \frac{\pi_{1}(x|z) + \pi_{2}(x|z)}{2} + \frac{1}{2} KL \quad \pi_{2}(x|z) \mid \mid \frac{\pi_{1}(x|z) + \pi_{2}(x|z)}{2} \quad .$$
 (15)

**Theorem 8** (weighting). The binary label mapping generates M paris label-1 examples and one pair of label-0 examples per simulation. We can reweight the binary classifier by letting the ELPD be  $E[\frac{C}{M+1} - 1(t=1) \log p(t=1|c(\phi)) + \frac{CM}{M+1} - 1(t=0) \log p(t=0|c(\phi))]$ , where  $C = \frac{(M+1)^{-2}}{2M}$  is a normalizing constant.

That is if we modify the training utility function or the validation data LPD to be:

$$\frac{1}{n} \sum_{i=1}^{N} \frac{C}{M+1} 1(t_i = 1) \log \Pr(t = 1 | c(\phi_i)) + \frac{CM}{M+1} 1(t_i = 0) \log \Pr(t = 0 | c(\phi_i))$$

then the resulting binary classification divergence in Thm. 1, i.e.,

is the conditional Jensen Shannon divergence (15).

#### A.1 SBC test and generative classifier

In Example 3, we stated the intuitive connection between the SBC rank test and using rank as the feature in a "generative" classifier (in contrast to our "discriminate" approach). Now we make such comparison more precise: Assuming the parameter  $\theta$  is one-dimensional. This is what SBC does: From  $\rho$  we obtain  $\theta$ , y, and for each  $\theta$  we compute the rank of it in the paired q samples who shared the same y:  $r = \begin{bmatrix} M \\ m=1 \end{bmatrix}$   $1(\theta \le \tilde{\theta}_m)$ . Repeating this simulation many times, we obtain a sample of r, and SBC will test if such r is discrete-uniformly distributed. This test could be done by simulating a reference distribution that matches the null. We do so by generating the ranks in the q samples:  $r_n = \begin{bmatrix} M \\ m \ne n \end{bmatrix}$   $1(\tilde{\theta}_n \le \tilde{\theta}_m) + 1(\tilde{\theta}_n \le \theta)$ . To be clear, there are other ways to generate reference samples, but all it matters is that the reference distribution matches the distribution of r under the null. The uniformity test in SBC now becomes to test if r and r have the same distribution. For example, we can do a two-sample Kolmogorov–Smirnov test or a permutation test.

In comparison, this is what naive Bayes would do: the naive Bayes classifier first estimates the two conditional distributions: Pr(r|t=0) from the empirical distribution of the prior rank r, and Pr(r|t=1) from the empirical distribution of r. Then to test if Pr(r|t=0) = Pr(r|t=1) becomes the same to testing if the empirical distribution of r is the same as the empirical distribution of r—the same task as in SBC. The point is that given any two-sample based test,

is operationally the same as

only use rank + naive Bayes classifier + density estimated by histogram + sample-based test.

#### **B** Implementations

### **B.1** A cheat sheet of useful pre-learned features

Table 1 summarizes how to use the statistical feature (whenever available) in both the binary and multiclass classifiers.

	binary	multiclass
full feature	$(\theta, y)$	$(\theta_1, \ldots, \not \theta, y)$
MCMC	$p(\theta, y), r(\theta, y)$	$p(\theta_t, y) r(\theta_t, y), 1 \le t \le K$
VI	$p(\theta, y), q(\theta y)$	$p(\theta_t, y) q(\theta_t y), 1 \le t \le K$
likelihood-free	$p(\theta), q(\theta y)$	$p(\theta_t), \ q(\theta_t y), \ 1 \le t \le K$
Table 1: A sheet sheet of useful mrs learned features		

Table 1: A cheat sheet of useful pre-learned features

#### B.2 Code and data

We share the python and Jax implementation of our binary and multiclass calibration code in https://github.com/yao-yl/DiscCalibration

In the MLP training we include a standard  $L^2$  weight decay (i.e., training loss function = cross entropy loss + tuning weight \*  $L^2$  penalization). We tune the weight of the decay term by a 5 fold cross-validation in the training set on a fixed grid{0.1, 0.01, 0.001, 0.0001}

In the cosmology experiment in Section 5, we use one-hidden-layer MLP with 64 nodes to parameterize the classifier with the form (11), with additional pre-learned features such as  $\log q(\theta|y)$  added as linear features. The approximate log density  $\log q(\theta|y)$  is known (i.e., we can evaluate the for any  $\theta$  and Y, or at least least up to a constant) in either the normalizing flow or the ensemble of the normalizing flows. One classification run with roughly one million examples took roughly two hour

	model 1 ( $d_y = 38$ )	model 2 ( $d_y = 151$ )	$\mod 3 (d_y = 1031)$
1-flow	0.008 (0.003)	0.05(0.006)	0.0001 (0.0002)
5-ensemble	0.001 (0.002)	0.008(0.003)	0.001(0.001)

Figure 8: The estimated divergence of two simulation based inferences (either using one flow or an ensemble of five flows) in three cosmology models. These three models have the same parameter space and involve expensive simulations. Here we apply our classifier approach and estimate the Jensen–Shannon divergence (and standard deviation) using weighted binary classifiers from S = 2000, M = 500 simulations: that is  $2000 * 50 \pm 1$  million classification examples. In this table, we run this classifier by the black-box multilayer perceptron (MLP) on the full( $\theta$ , y) vector, while the table we have in the main paper further adds  $\log q(\theta|y)$  as a pre-learned feature since it is known in the normalizing flows. In this black box estimate table, we would pass at-test and not even detect miscalculation. Adding a statistically-meaningful feature greatly improves the classifier prediction, and tightens the bound of the divergence estimate.

cpu time on a local laptop. It would be more efficient to run the classification on GPU, thought the classification cost is overall negligible compared with the simulation step  $y|\theta$  which is pre-computed and stored.

In the closed-form Gaussian experiment in Section 5, we consider the easiest setting:  $\theta \in \mathbb{R}^d \sim \text{MVN}(\mathbf{0}, \text{Id}\,d)$  and a normal data model  $y|\theta \sim \text{MVN}(\theta, \text{Id}\,d)$ . Indeed, we have kept both the true distribution and the sampling corruption mean-field to make the calibration task easier for the traditional SBC rank-test. The power of the traditional SBC rank-test would be even worse if we add some interaction terms in the corrupted samples  $q(\theta|y)$ , while our method leans the joint space sampling quality by default. The exact posterior  $p(\theta|y)$  is known but we pretend we cannot evaluate it. We set d=16 to be comparable with the real data example. The right panel validates Thm. 3 fixed the simulation size S=5000 and vary M. We have tried other fixed M=1000 but it seems the resulting graph is very similiar. The confidence band in Figure 3 and Figure 4 is computed using 1.96 times standard error from repeated experiments.

Table 8 shows the naive classifier: using the bandbox MLP without the linear features. The estimate is rather loose compared with the table we obtained in Section 5 which used statistically-meaningful features.

#### C Proofs of the theorems

We prove the theorems in the following order: Thm. 7, Thm. 1, Thm. 8, Thm. 6, Thm. 3, Thm. 5, Thm. 4, and finally Thm. 2.

#### **Contents**

C.1	The classification divergence	15
C.2	Reweighting (Theorem 8) leads to symmetric divergence	18
C.3	MCMC without thinning (Theorem 6)	19
C.4	Large sample limit and rate as $M \to \infty$ (Theorem 3)	20
C.5	Valid hypothesis testing (Theorem 4)	22
C.6	Why does classifier help test? (Theorem 5)	23
C.7	The maximum discriminative generator (Theorem 2)	23

#### C.1 The classification divergence

To prove our Theorem 1, we first state an equivalent but a more general theorem. Note that we are allowing non-IID classification examples for the shared z variables in this theory.

**Theorem 9** (The general relation between divergence and classification error). *Define a simulation process*  $\pi(z, k, x)$  *through the following process, among the three random variables*, $z \in \mathbb{R}^n$ ,  $x \in \mathbb{R}^d$  can have any dimension, and  $k \in \{1, 2, ..., K\}$  is an integer.

- 1. Sample  $z \sim \pi(z)$
- 2. Sample  $k \sim \text{Categorical}(w)$ , where W is a simplex.
- 3. Sample  $X \sim \pi_k(x|z)$
- 4. Return (z, k, x)

Define  $U(\pi)$  to be the expected log probability or cross-entropy of an optimal classifier trained to predict k from (x, z), i.e.

$$U(\pi) = \max_{x} E_{\pi} U(k, c(x, z)),$$

where  $U(k, c) = \log c(k)$ . Then<sup>3</sup>

$$U(\pi) - \bigvee_{k=1}^{K} W_k \log w_k = \bigvee_{n=1}^{K} W_k KL (\pi_n(x|z) \parallel \pi(x|z))$$

where  $\pi(x|z) = \bigcap_{k=1}^{K} W_k \pi_k(x|z)$  is the mixture over all groups. We define  $D(\pi) = U(\pi) - \bigcap_{k=1}^{K} W_k \log w_k$ . Then this D has the following properties:

- 1. (lower bound).  $D(\pi) \ge 0$  . It achieves zero if and only if all  $P_k$  are the same, i.e.,  $\Pr(X^k \in A) = \Pr(X^j \in A)$  for any measurable set A.
- 2. (upper bound).  $D(\pi) \leq \bigcap_{n=1}^{P} W_k \log W_k \cdot This \ maximum \ is \ achieved \ is \ P_k \ are \ disjoint, \ i.e., Pr(X^k \in A) Pr(X^j \in A) = 0 \ for \ any \ measurable \ set \ A.$
- 3. (alternative upper bound when the difference is small).  $D(\pi) \leq \max_{j,k} KL(\pi_k(x|z), \pi_j(x|z))$ .
- 4. (reparametrization invariance). For any fixed bijective transformation  $X \to g(x)$ , let  $P_k^g$  be the law of  $g(x^k)$ , then  $D(p_k^g, \dots, p) = D(p_k, \dots, p)$ .
- 5. (transformation decreases divergence). For any fixed transformation  $X \to g(x)$  let  $P_k^g$  be the law of  $g(x^k)$ , then  $D(p_k^g, \dots, p) \leq D(p_k, \dots, p)$ .

*Proof.* The (unknown) optimal conditional classifier is the Bayes classifier:

$$\hat{\mathsf{Pr}}(t=k|x,\,z) = \quad \underbrace{\mathsf{P}^{W_k \pi_k(x|z)}}_{j \ W_k \pi_j(x|z)}.$$

Let  $r(x|z) = \prod_{j=1}^{K} W_k \pi_j(x|z)$  be mixture density marginalized over the group index K.

$$\mathsf{KL}\; (\pi_n\,(x|z) \parallel \pi(x|z)) \;\; := \;\; \begin{array}{c} \mathsf{Z} & \mathsf{Z} \\ \rho(z) & \pi_k\,(x|z) \log & \frac{\pi_k\,(x|z)}{r(x|z)} & dx & dz. \end{array}$$

<sup>&</sup>lt;sup>3</sup>We emphasize that we are using the notation of conditional KL divergence:

The expected log predictive density of the optimal classification is then

ELPD( 
$$\hat{\Pr}(t = k|x, z)$$
) = E  $z$   $r(x|z)$   $\frac{X}{k=1}$   $\frac{W_k \pi_k(x|z)}{r(x|z)} \log \frac{W_k \pi_k(x|z)}{r(x|z)}$   $dx$ 

$$= E z \frac{X}{k=1} \frac{X}{W_k} \frac{Z}{\pi_k(x|z) \log \frac{\pi_k(x|z)}{r(x|z)}} + \log w_k dx$$

$$= E z \frac{X}{k=1} \frac{X}{W_k \log(w_k)} + \frac{X}{k=1} \frac{Z}{m_k(x|z) \log \frac{\pi_k(x|z)}{r(x|z)}} dx$$

$$= \frac{X}{k=1} \frac{X}{k=1} \frac{X}{m_k \log(w_k)} + \frac{X}{m_k \log(w_k)} \frac{Z}{m_k \log(w_k)} \frac{Z}{m_k \log(w_k)} + \frac{Z}{m_k \log(w_k)} \frac{Z}{m_k \log(w_k)} \frac{Z}{m_k \log(w_k)} \frac{Z}{m_k \log(w_k)} + \frac{Z}{m_k \log(w_k)} \frac{Z}{m_k \log(w_k)$$

When ELPD –  $\Pr_{k=1}^{P} W_k \log(w_k) = 0$ ,  $\Pr_{k=1}^{P} W_k \pi_k(x|z) = 0$ , hence these K conditional densities  $\pi_k(x|z)$  are the same almost everywhere.

The reparametrization invariance is directly inherited from KL diverge.

The upper bound is the consequence of ELPD  $\leq 0$  for category classification and the lower bound is that KL divergence  $\geq 0$ .

Let's emphasize again that we are following the notation of conditional divergence in  $\S$ ], such that for any two joint density p(theta, y) and  $q(\theta, y)$ , the conditional KL divergence is  $KL(p(\theta|y) | q(\theta|y)) := E_{p(y)} E_{p(y|\theta)} \log \frac{p(y|\theta)}{q(y|\theta)}$ , an expectation is taken over Y.

The procedure in Theorem 9 can reflect many different generating processes. Let's recap the examples we show in Section 2.

In **Example 1** of the main paper, the binary case, K = 2, and we generate M samples from Q.  $\pi_1 = p(\theta|y)$ ,  $\pi_2 = q(\theta|y)$ ,  $W_1 = 1/M$ ,  $W_2 = (M-1)/M$ . Our simulation process in example 1 is equivalent to

- 1. Sample  $Y \sim \text{marginal } p(y)$ ,
- 2. Sample  $k \sim \text{Categorical}(w)$ , where  $W_1 = 1/M$ ,  $W_2 = (M-1)/M$ .
- 3. Sample  $X \sim \pi_k(x|z)$ , where  $\pi_1 = p(\theta|y)$ ,  $\pi_2 = q(\theta|y)$ .
- 4. Return (z, k, x)

As a direct consequence of Theorem 9, the resulting divergence in the binary case is

$$D_1(p, q) = w_1 KL \left( p(\theta|y) \parallel r(\theta|y) \right) + w_2 KL \left( q(\theta|y) \parallel r(\theta|y) \right),$$

where  $r = W_1 D + W_2 Q$ .

In **Example 2**: Binary label with no  $\mathcal{Y}$ , we generate M samples from Q. K = 2,  $\pi_1 = p(\theta)$ ,  $\pi_2 = q(\theta)$ ,  $W_1 = 1/(M+1)$ ,  $W_2 = M/(M+1)$ . Our simulation process in example 2 is equivalent to

- 1. Sample  $k \sim \text{Categorical}(w)$ , where  $W_1 = 1/M$ ,  $W_2 = (M-1)/M$ .
- 2. Sample  $X \sim \pi_k(x)$ , where  $\pi_1 = p(\theta)$ ,  $\pi_2 = q(\theta)$ .
- 3. Return (k, x)

From Theorem 9, the resulting divergence reads

$$D_3(p, q) = w_1 KL(p(\theta) \parallel r(\theta)) + w_2 KL(q(\theta) \parallel r(\theta))$$

In the multivariate Example 4, K = M + 1, and the individual density is

$$\pi_k(\theta_{1:K}) = p(\theta_k|y) \sum_{m = k}^{Y} q(\theta_m|y)$$

.

From Theorem 9, the resulting divergence reads

$$D_{4}(p, q) = \mathsf{KL} \left[ \begin{array}{c} \mathsf{Y} \\ p(\theta_{1}|y) \\ p(\theta_{m}|y) \end{array} \right] \left[ \begin{array}{c} \mathsf{X} \\ \mathsf{K} \\ \mathsf{k}=1 \end{array} \right] p(\theta_{k}|y) \left[ \begin{array}{c} \mathsf{Y} \\ \mathsf{k} \\ \mathsf{$$

We give more results on this  $D_4$  in Section C.4.

The discriminative calibration is designed to detect and gauge the difference between  $p(\theta|y)$  and  $q(\theta|y)$ . From now on we use a long vector  $\theta_{1:M+1}$  (the index starts by 1) to represent the simulation draw first generates one  $\theta_1$  from the prior, one data point y from the sampling distribution  $p(y|\theta_1)$ , and finally M draws from the approximate inference  $\theta_{1,...,M+1} \sim q(\theta|y)$ .

The "label generating process"  $\Phi$  takes a simulated draw as its input and returns a vector of features  $\phi$  and labels t.

$$\Phi: (y,\; \theta_{1,\ldots,M+1} \;\;)\; 7 \rightarrow \{(\phi_{1},\; t_{1}),\; \ldots\;,\; (\phi_{t}\;,\; t_{K}\;)\}.$$

The classifier will then see this feature and labels (the index k is disregarded) across all simulations. We only ask that the map  $\Phi$  be within in set F: satisfying that:

F:  $\phi | y$  is independent of t if  $p(\theta | y) = q(\theta | y)$ , a.e.

#### **Proof of Theorem 1.**

*Proof.* The classifier sees features  $\phi$  and labels t. Given a t,  $(\phi, y)$  are IID draws from  $\rho(y)\pi(\phi|y, t)$ . In case Y is not observable by the classifier,  $\phi | t$  are IID draws from  $(p, p(y))p(\phi | y, t)$ .

To use Theorem 9, let  $\mathbf{x}^k = \{(\phi_j, y)|t_j = k\}$  and  $W_k = \Pr(t = k)$ . The divergence from Theorem 9: reads

reads 
$$D^{\text{opt}}(p, q, \Phi) = \bigvee_{k=1}^{\infty} W_k \text{KL} \quad \pi_1(\phi) , \quad \bigvee_{k=1}^{\infty} W_k \pi_1(\phi) ,$$
 in which  $\pi_k(\phi) = p(\phi|t=k) = \bigvee_{k=1}^{\infty} P(y)\pi(\phi|y, t=k)dy$  is the *Y*-averaged  $\phi$  margin. This is

essentially Theorem 7.

If p = q, from F, the law of  $\Phi$ , Y is independent of f (Y is always independent of f marginally), so that  $\pi_k(\phi) = \int_{Y} p(y)\pi(\phi|y, t = k)dy$  is also independent of t, as it only depends on Y and  $\phi$ , hence the divergence remains 0. 

## C.2 Reweighting (Theorem 8) leads to symmetric divergence

We now prove Theorem 8. The binary label mapping generates pairs label-1 examples and one pair of label-0 examples per simulation. We can reweight the binary classifier by letting the ELPD be  $E[\frac{C}{M+1} \ 1(t=1) \log p(t=1|c(\phi)) + \frac{CM}{M+1} \ 1(t=0) \log p(t=0|c(\phi))]$ , where  $C = \frac{(M+1)^2}{2M}$  is a normalizing constant. We want to prove that after weighting, the classification divergence is symmetric Jensen shannon distance, as if we have balanced classification data.

*Proof.* After the label-reweighing, the optimal classifier is the balanced-class classifier

$$c(t=1|\phi) = \frac{\Pr(\phi|t=1)}{\Pr(\phi|t=1) + \Pr(\phi|t=0)}.$$

The re-weighted ELPD in the population is

$$\begin{split} \mathsf{ELPD} &= \mathsf{E}_{\,\,\phi} \big[ \mathsf{Pr}(t=0|\phi) \, \, \frac{CM}{M+1} \log \, c(t=0|\phi) \, + \mathsf{Pr}(t=0|\phi) \, \, \frac{C}{M+1} \log \, c(t=1|\phi) \big] \\ &= \mathsf{E}_{\,\,y} \, \, \frac{p(\theta|y) + M \, q(\theta|y)}{M+1} \, \, \frac{p(\theta|y)}{p(\theta|y) + M \, q(\theta|y)} \frac{CM}{M+1} \log \frac{p(\theta|y)}{2r(\theta|y)} \, + \, \frac{M \, q(\theta|y)}{p(\theta|y) + M \, q(\theta|y)} \frac{C}{M+1} \log \frac{q(\theta|y)}{2r(\theta|y)} \\ &= \mathsf{E}_{\,\,y} \, \, \, C \frac{2M}{(M+1)^{\,\,2}} \, \, \frac{p(\theta|y)}{2} \log \frac{p(\theta|y)}{2r(\theta|y)} \, + \, \frac{q(\theta|y)}{2} \log \frac{q(\theta|y)}{2r(\theta|y)} \, \, . \end{split}$$

With an appropriate normalizing constant  $C = \frac{(M+1)^{-2}}{2M}$ ,

ELPD = E<sub>y</sub> 
$$\frac{1}{2}p(\theta|y)\log\frac{p(\theta|y)}{r(\theta|y)} + \frac{1}{2}q(\theta|y)\log\frac{p(\theta|y)}{r(\theta|y)} - \log 2.$$

That is, the reweighted ELPD + log 2 is the conditional Jensen Shannon divergence (15) between  $p(\theta|y)$  and  $q(\theta|y)$ .

## C.3 MCMC without thinning (Theorem 6)

To diagnose MCMC, we sample  $(\theta, y) \sim p(\theta, y)$ , and sample (potentially autocorrelated) draws  $(\tilde{\theta}_1, \dots, \tilde{\theta}_M)$  from a MCMC sampler—whose stationary distribution we believe is  $q(\theta|y)$  (i.e., marginally  $\tilde{\theta}_i$  is from  $q(\theta|y)$ ), and generate examples  $((t_1, \phi_1), \dots, (t_{M+1}, \phi_{M+1}))$ , from the multiclass permutation (see definition in Example 4), such that

$$\phi = (\theta_0, \theta_1, \ldots, \theta_n).$$

Then we run an exchangeable classifier parameterized by

$$\Pr(t = k | (\theta_0, \theta_1, \dots, \theta_n, y)) = \Pr \frac{\exp(g(\theta_k, y))}{\prod_{k'=0}^{M} \exp(g(\theta_{k'}, y))},$$
(16)

where g is any  $\Theta \times Y \rightarrow R$  mapping to be learned.

Note that here  $\Pr(t = k | (\theta_0, \theta_1, \dots, \theta_n, y))$  is the classifier model we restrict it to be, not the true population  $\pi$ . In general  $\pi(t = k | (\theta_0, \theta_1, \dots, \theta_n, y)) \neq \Pr(t = k | (\theta_0, \theta_1, \dots, \theta_n, y))$  Roughly speaking, the optimal classifier should be the Bayes classifier projected to the restricted space (16), and we need to proof that this *projected restricted* solution turns out to be the same as the IID case 10, which is not trivial in general.

*Proof.* Intuitively, this separable network design(16) avoids the interaction between  $\tilde{\theta}_m$  with  $\tilde{\theta}_m$  and disallows the network to predict t based on the autocorrelation or clustering of  $\tilde{\theta}$ .

Because of the permutation design and we define q to be the marginal distribution of  $\theta$ , first, we know that

$$\pi(\theta_k | y, t = k) = p(\theta_k | y), \quad \pi(\theta_k | y, t \neq k) = q(\theta_k | y),$$

in the MCMC population  $\pi$  (there is no need to address the joint for now).

From this we obtain the conditionals in the population

$$\pi(t = k | \theta_k, y) = \frac{p(\theta_k | y)}{p(\theta_k | y) + M q(\theta | y)}$$

and further the ratios for any m, k withtin the index set  $\{0, 2, \ldots, M\}$  and  $m \neq k$ :

$$\pi(t=k|\theta_k,\,\theta_m,\,y)=\frac{p(\theta_k|y)q(\theta_k|y)}{p(\theta_k|y)q(\theta_m|y)+q(\theta_k|y)p(\theta_k|y)+(M-1)q_{12}(\theta_k,\,\theta_m)}$$

Here  $q_{12}(\theta_k, \theta_m)$  is the joint distribution of two out of M draws from the MCMC sampler  $q(\theta|y)$ , which would often not be the same as  $q(\theta_k|y)q(\theta_m|y)$  because of the autocorrelation of Markov chains. We do not need to specify the form of this joint. The key blessing is that when  $\theta_k$  is from  $p(\theta|y)$  and  $\theta_m$  is from  $q(\theta|y)$ , then they are independent (the numerator).

Next, from the line above we obtain the ratio estimate in the true population  $\pi$ :

$$\frac{\pi(t=k|\theta_{k},\theta_{m},y)}{\pi(t=m|\theta_{k},\theta_{m},y)} = \frac{p(\theta_{k}|y)q(\theta_{k}|y)}{q(\theta_{k}|y)p(\theta_{k}|y)} = \frac{p(\theta_{k}|y)/q(\theta_{k}|y)}{p(\theta_{m}|y)/q(\theta_{m}|y)}.$$
(17)

Now we project the restricted classifier (16). Intuitively, the separable classifier "almost" only depends on  $\theta_t$ , except for the normalizing constant. We can remove the dependence on the normalizing constant by specifying the ratio

$$\frac{\Pr(t = k | (\theta_0, \theta_1, \dots, \theta_n, y))}{\Pr(t = m | (\theta_0, \theta_1, \dots, \theta_n, y))} = \frac{\exp(g(\theta_k, y))}{\exp(g(\theta_m, y))}$$

Marginalizing out all other components, we obtain

$$\frac{\Pr(t=k|(\theta_{k},\theta_{m},y))}{\Pr(t=m|(\theta_{k},\theta_{m},y))} = \frac{\exp(g(\theta_{k},y))}{\exp(g(\theta_{m},y))}.$$
(18)

Matching the restriction (18) with the true population (17), the restricted projection is attainable if and only if

$$\exp(g(\theta_t, y)) = p(\theta_t, y)/q(\theta_t|y),$$

so that the optimal classifier needs to be

$$Pr(t|\theta_1,\ldots,\theta,y) = P \frac{p(\theta_t,y)/q(\theta_t|y)}{\sum_{k=1}^{K} p(\theta_k,y)/q(\theta_k|y)}.$$

It happens that this MCMC restricted optimal classifier matches the IID optimal classifier (10). It follows from the proof in Theorem 9 that the classification divergence using the restricted classifier is still  $D_4(p, q)$ , as if  $\{\partial_m\}_{m=1}^M$  are IID samples from  $q(\theta|y)$ .

## C.4 Large sample limit and rate as $M \to \infty$ (Theorem 3)

In the multivariate **Example 4**, from each draw  $\theta$ , y,  $\theta_{1:M}$ , we generate M+1 examples from K:=M+1 classes; one label from each. We will use the index starting from 1 in this subsection. For the k-th example,  $t_k = k$ , and the permutation for classifier reads  $\theta_k$  is a vector including y and K copies of  $\theta$ , where the k-th copy is the prior draw, and the remaining  $\theta$  are from q(|y). Slightly abused the notation, in this subsection, we call this long feature vector as  $\theta_{1:K}$ ; it is a vector in the  $\Theta^{M+1}$  space. We now prove Theorem 3.

*Proof.* First, we write the true conditional label probability in this process:

$$\pi(t|\theta_{1:K}, y) = \underbrace{P \frac{p(\theta_t|y)}{t} \underbrace{Q}_{j=t} \frac{p(\theta_t|y)}{p(\theta_t|y)}}_{t} \underbrace{P(\theta_t|y)}_{j=t} \underbrace{P \frac{p(\theta_t|y)}{p(\theta_t|y)}}_{j} = \underbrace{P \frac{p(\theta_t|y)}{p(\theta_t|y)}}_{j} \underbrace{P(\theta_t|y)}_{j} \underbrace{P(\theta_t|y)}_{j}.$$

Plug it as the classifier, we obtain the optimal ELPD or negative cross entropy: ELPD =  $E \log \pi(t|\theta_{1:K}, y)$ ,

ELPD = E 
$$P \frac{p(\theta_{t}|y)/q(\theta_{t}|y)}{\sum_{k=1}^{K} p(\theta_{k}|y)/q(\theta_{k}|y)}$$

$$= E \log \frac{p(\theta_{t}|y)}{q(\theta_{t}|y)} - E \log \frac{X}{k} p(\theta_{k}|y)/q(\theta_{k}).$$

The first term above is simply

$$\mathsf{E} \log \frac{p(\theta_t | y)}{q(\theta_t | y)} = \mathsf{KL} \left( p(\theta | y) || q(\theta | y) \right)$$

According to our definition (4), the divergence is ELPD offset by an entropy term,

$$D_A := ELPD + log K$$
.

We now derive the limit of  $D_4 - KL(p(\theta|y)||q(\theta|y))$  when  $M \to \infty$  (or equivalently  $K = M + 1 \to \infty$ )

$$\Delta := D_4 - KL (p(\theta|y)||q(\theta|y))$$

$$= ELPD + \log K - KL (p(\theta|y)||q(\theta|y))$$

$$= \log K - E \log \frac{X}{k} \frac{p(\theta_k|y)}{q(\theta_k|y)}$$

$$= - E \log \frac{1}{K} \frac{X}{k=1} \frac{p(\theta_k|y)}{q(\theta_k|y)}$$

Given any label value  $1 \le K$ ,  $\theta_t \sim p(\cdot|y)$ , and all the remaining  $\theta_i \sim q(\cdot|y)$  for j'=t.

Let

$$X_k = \frac{1}{K} \frac{X}{k} \frac{p(\theta_k | y)}{q(\theta_k | y)} = \frac{1}{K} \frac{p(\theta_t | y)}{q(\theta_t | y)} + \frac{1}{K} \frac{X}{k} \frac{p(\theta_t | y)}{q(\theta_t | y)}.$$

The first term

$$\frac{1}{K} \sum_{k=1}^{K} \frac{p(\theta_t | y)}{q(\theta_t | y)} \rightarrow 0.$$

The second term, call it  $\Delta_2$  is the mean of IID means as  $\theta_j \sim q(\cdot|y)$  for  $j \not= t$ , the law of large number yields

$$\Delta_2 := \frac{1}{K} \frac{X}{\sum_{i \neq j} \frac{p(\theta_i \mid y)}{q(\theta_i \mid y)}}, \quad \Delta_2 \to \mathsf{E}_{x \sim q(x \mid y)} \quad \frac{p(x \mid y)}{q(x \mid y)} = 1.$$

This proves that

$$\Delta \rightarrow 0$$
, as  $K \rightarrow \infty$ .

Hence,

$$D_4 - KL (p(\theta|y)||q(\theta|y)) \rightarrow 0$$

which finished the first part of the proof.

Now let's derive the rate. Under regularity conditions, for example, the variance of density ratio is bounded, i.e., there exists a constant  $C < \infty$ , such that for all Y,

$$\operatorname{Var}_{\theta_t \sim q(\theta^{-t}|y)} \quad \frac{p(\theta_t|y)}{q(\theta_t|y)} < C,$$

then CLT holds, such that the second term above has a normal limit,

$$\overline{K}(\Delta_2 - 1) \rightarrow \text{normal}(0, \sigma^2)$$
, in distribution,

where

$$\sigma^{2} = \operatorname{Var}_{q}\left(\frac{p(\theta_{t}|y)}{q(\theta_{t}|y)}\right)$$

$$= \operatorname{E}_{\theta, y \sim q(\theta, y)} \frac{p(\theta|y)}{q(\theta|y)} - 1$$

$$= \operatorname{E}_{y} \operatorname{E}_{\theta \sim q(\theta|y)} \frac{p(\theta|y)}{q(\theta|y)} - 1$$

$$= \chi^{2}\left(p(\theta|y) \mid\mid q(\theta|y)\right),$$

which is the definition of the conditional chi-squared divergence.

Consider a Taylor series expansion,  $log(1 + x) = x - \frac{1}{2}X^2 + o(x^2)$ . Using the Delta method to express the log function and the Slutsky theorem to ignore the zero term, we get

$$K \in \log(\Delta_2) \rightarrow -\frac{\sigma^2}{2}$$

Plug this in the definition of  $\Delta$  we obtain the desired convergence rate,

$$D_4 = KL (p(\theta|y)||q(\theta|y)) - \frac{1}{2M}X^2 (p(\theta|y)||q(\theta|y)) + o(1/M).$$

This proves the claims in Theorem 3.

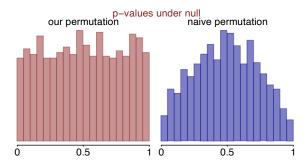


Figure 9: Our designed permutation vs the naive permutation.

Table 2: Batch permutation of labels

Original batch from one run label t | features **b** 

label t	features Ψ	
$(\theta, y)$	1	
$(\tilde{\theta}_1, y)$	0	
$(\tilde{\theta}_2, \mathbf{y})$	0	
$(\tilde{\theta}_M, y)$	0	

	permuted labels		atch of permutation
label t		label t	
$\overline{}$			

label t	label <sup>t</sup>
0	0
1	0
0	0
0	1

### C.5 Valid hypothesis testing (Theorem 4)

It is not trivial to have a valid permutation test since we do not have IID examples. One naive permutation is to permutate all the labels t (across batches t). The resulting permutation testing is not uniform under the null (Figure 9).

In contrast, when we permute the label, we only ask for within-batch permeation (permute the index  $\{t_1, \dots, t\}$ ) in each batch). See Table 2 for an illustration. Let's prove this permutation is valid.

*Proof.* Under the null,  $p(\theta|y) = q(\theta|y)$  almost everywhere. According to the label mapping  $\Phi \in F$ , label t is independent of  $\theta$  given y. We first show that label t is independent of  $\theta$ .

In general conditional independence does not lead to unconditional independence. But because here we design the generating process by  $\pi(y, t, \phi) = \pi_Y(y)\pi_t(t)\pi(\phi|y, \theta)$ . Under the null we have  $\pi(y, t, \phi) = \pi_t(t)(\pi_Y(y)\pi(\phi|y))$ . Hence t needs to be independent of  $\phi$ .

For any classifier c, because  $c(\phi)$  is a function of  $\phi$ ,  $c(\phi)$  and t are also independent. Let  $\pi_{\Phi}(\phi)$  be the marginal distribution of  $\phi$  in  $\pi$ .

Now we are computing the cross entropy (with respect to the population  $\pi$ ). LPD =  $\int_{j=1}^{N} U(c(\phi_n), t_n)$  where U is the log score  $U(P, x) = \log P(x)$ . It is a random variable because we have a finite validation-set. Now we are conducting a permutation of label<sup>t</sup>, the permuted label  $\tilde{t}$  is an independent draw from the same marginal distribution of t,  $\pi_t(t)$ . Because of the independence,  $E_{\phi,t}$   $U(c(\phi), t) = E_t E_{\phi} U(c(\phi), t)$ , hence the permuted LPD  $_b = LPD$ , where LPD  $_b$  is the computed LPD in the L-th permutation. Therefore L-th permutation (with an infinitely amount of repeated random permutation), this L-value will be uniform on L-th permutation of L-th permutation).

In practice, the computed p-value from finite permutation can only take values on a finite set  $\{0, 1/B, \ldots, (B-1)/B, 1\}$ . More precisely, under the null hypothesis, this permutation test p-value is *discreetly-uniform* on this set. This is because for any  $0 \le m \le B$ 

$$Pr(p = m/B) = \begin{pmatrix} Z_1 & B \\ & m \end{pmatrix} p^m (1-p)^{B-m} dp = 1/(B+1), \quad \forall m.$$

Hence, the permutation P-value based on B random permutations is uniformly distributed on the set  $\{0, 1/B, \ldots, (B-1)/B, 1\}$ .

#### C.6 Why does classifier help test? (Theorem 5)

Let's recap the binary label-generating process:

- 1. Sample  $Y \sim \text{marginal } p(y)$ ,
- 2. Sample  $k \sim \text{Bernoulli}(w)$ , where  $W_1 = 1/M$ ,  $W_2 = (M-1)/M$ .
- 3. Sample  $\theta \sim \pi_k(\theta|z)$ , where  $\pi_0 = p(\theta|y)$ ,  $\pi_1 = q(\theta|y)$ .
- 4. Return (*z*, *k*, *x*)

Theorem 5 state that the optimal classifier has a sufficiency propriety in that (a) le  $\mathcal{C}$  be the probability of label 1 in the optimal classifier as per (7), and let  $\mathcal{T}_c^0$  and  $\mathcal{T}_c^q$  be the one-dimensional distribution of this  $\hat{c}(\phi)$  when  $(\theta, y)$  is sampled from  $p(\theta, y)$  or from  $p(y)q(\theta|y)$  respectively, then (i) Conditional on the summary statistic  $\mathcal{C}$ , the label  $\mathcal{C}$  is independent of all features  $\phi = (\theta, y)$ . (ii) There is no loss of information in divergence as the joint divergence is the same as the projected divergence,  $D_1(p, q) = D_1(\mathcal{T}_c^0, \mathcal{T}_c^q)$ .

*Proof.* There are three random variables,  $\theta$ , y, and t. The optimal classifier c is the probability in this true joint:

$$\hat{c}(\theta,\,y)=\Pr(t=1|(\theta,\,y)).$$

To show sufficiency or conditional independence, all we need to show is that, conditional on any given value of  $^{C}$ ,  $Pr(t = 1 | (\theta, y), c)$  does not depend on  $(\theta, y)$  (in other words,  $Pr(t = 1 | (\theta, y), c)$  is a function of  $^{C}$  only). This becomes obvious as

$$Pr(t = 1 | (\theta, y), c(\theta, y)) = Pr(t = 1 | (\theta, y)) = ^c c(\theta, y).$$

Now we prove that there is "no loss" in divergence.

$$D_1(p, q) = w \mathsf{KL}(p(\theta, y)||r(\theta, y)) + (1 - w) \mathsf{KL}(q(\theta, y)||r(\theta, y))$$

We express the first term in C

$$KL(p(\theta, y) || wp(\theta, y) + (1 - w)q(\theta, y))$$
=  $KL(\pi(\theta, y|t = 0) || \pi(\theta, y))$   
=  $KL(\pi(\hat{c}|t = 0) || \pi(\hat{c})) - KL(\pi(\theta, y|\hat{c}, t = 0) || \pi(\theta, y|\hat{c}))$ 

This steps uses the chain rule of KL divergence:  $KL[p(x, y) \mid q(x, y)] = KL[p(x) \mid q(x)] + KL[p(y \mid x) \mid q(y \mid x)]$ 

Using conditional independence:

$$\pi(\theta, y|\hat{c}, t = 0) = \pi(\theta, y|\hat{c}, t = 1)$$

Hence  $KL(\pi(\theta, y|\hat{c}, t = 0) || \pi(\theta, y|\hat{c})) = 0$ . Therefore,

$$\mathsf{KL}(p(\theta, y)||r(\theta, y)) = \mathsf{KL}(\pi(\hat{c}|t=0)||\pi(\hat{c}))$$

where  $\pi(c) = \pi(\hat{c}|t=0) + (1-w)\pi(\hat{c}|t=0)$ 

Similarly,

$$\mathsf{KL}(q(\theta, y)||r(\theta, y)) = \mathsf{KL}(\pi(\hat{c}|t=1)||\pi(\hat{c}))$$

This proves  $D_1(p, q) = D_1(\pi_c^p, \mathcal{T}_c^q)$ .

#### C.7 The maximum discriminative generator (Theorem 2)

We save the proof of Theorem 2 in the end for its length.

The generator  $\phi$  contains a few degrees of freedom: the number of classes K, the number of examples L, and how to design and the label-feature pairs. In the binary labeling:  $t_1 = 1$ ,  $\phi_1 = (\theta_1, y)$  and the remaining  $t_k = 0$ ,  $\phi_k = (\theta_k, y)$  for  $2 \le k \le M + 1$ . The multi-class  $\Phi^*$  assigns labels 1:K as

$$\Phi^* : t_k = k, \ \phi_k = (\text{Perm}_{1 \to k} (\theta_1, \dots, \theta_n), \ y), \ 1 \le k \le K = M + 1.$$
 (19)

Before the main proof that the multiclass permutation creates the largest divergence, let's first convince that the multi-class classification produced higher divergence than the binary one.

In the binary classification with M = 1 (one draw from P and one draw from Q)

$$D_1(p, q) = \frac{1}{2} \text{KL} \quad p(\theta|y), \frac{p(\theta|y) + q(\theta|y)}{2} \quad + \frac{1}{2} \text{KL} \quad q(\theta|y), \frac{p(\theta|y) + q(\theta|y)}{2}$$

In the multi-class classification with M = 1,

$$D_4(p, q) = \mathsf{KL} \quad p(\theta^1 | y) q(\theta^2 | y), \ \frac{p(\theta^1 | y) q(\theta^2 | y) + q(\theta^1 | y) p(\theta^2 | y)}{2}$$

Use the joint KL > marginal KL, we have

$$D_4(p, q) \ge KL \quad p(\theta^1|y), \frac{p(\theta^1|y) + q(\theta^1|y)}{2}$$

Likewise,

$$D_4(p, q) \ge KL \quad q(\theta^2|y), \frac{p(\theta^2|y) + q(\theta^2|y)}{2}$$

Add these two lines we obtain

$$D_4(p, q) \ge D_1(p, q)$$
.

To prove that this multi-class permutation produced the uniformly largest divergence (across M, K, P, Q) optimal, we organize the proof into lemmas 5 to 9. For notation brevity, we denote  $\hat{M} := M + 1$  in these lemmas to avoid using index M + 1.

**Lemma 10.** For an arbitrary integer  $^L$ , any given output space  $^Y$ , and any input space  $^X$  that has at least  $^L$  elements, if there are two functions mapping  $^X$  to  $^Y$ 

$$f_1, f_2 : (x_1, \ldots, x) \not\rightarrow y \in Y.$$

satisfying the following propriety:

• for any probability distribution  $\pi$  on X, when  $x_1, \dots, X$  are L iid random variables with law  $\pi$ ,  $f_1(x_1, \dots, X)$  has the same distribution as  $f_2(x_1, \dots, X)$ ,

then there must exist a permutation of 1: L, denoted by  $\sigma(1:L)$ , such that

$$f_2(x_1,\ldots,X) = f_1(x_{\sigma(1)},\ldots,X_{\sigma(L)}).$$

*Proof.* For any L distinct values  $a_1, \ldots, a_n, a_n \in X$ , let  $\pi$  be a mixture distribution of L delta functions:

$$\pi = \sum_{m=1}^{X} \delta(a_m) p_m,$$

where the m-th mixture probability is

$$p_m = C \quad \frac{1}{2}^{L^{m-1}}, \quad C^{-1} = \overset{\times}{\sum} \quad \frac{1}{2}^{L^{m-1}}.$$

C is chosen such that  $P_{m=1}^{L} p_m = 1$ .

Now that  $X_1, \dots, X$  are L IID random variables from this  $\pi, f_1(x_1, \dots, X)$  is also a mixture of delta functions, For any sequence of input indices  $(u_1, \dots, U) \in \{1 : L\}^{-L}$ ,

$$\Pr(f_{1}(x_{1}, \ldots, x_{l}) = f_{1}(au_{1}, \ldots, a_{l})) = \bigvee_{m=1}^{k} ((C/2^{L^{m-1}})^{P_{j=1}^{L}} (u_{j} = m))$$

$$= C^{L}(\frac{1}{2})^{P_{m=1}^{L}} (P_{j=1}^{L} (u_{j} = m)L^{m-1}), \qquad (20)$$

in which the power index can be written as

$$\begin{cases} x & x \\ x & x \\ y & 1 \\ y$$

as an L-decimal-integer.

Next, we study the law of  $f_2(x_1, \ldots, x_n)$ :

$$Pr(f_2(x_1,\ldots,x_n)=f_2(a_1,\ldots,x_n))=C^{-1}(\frac{1}{2})^{(1,1,\ldots,1)-1}$$

Because  $f_2(x_1, \dots, X)$  and  $f_1(x_1, \dots, X)$  have the same distribution,  $f_2(a_1, \dots, A)$ ) needs to match the value at which  $f_1(x_1, \dots, X)$  has probability  $C^L(\frac{1}{2})^{(1,1,\dots,1)}$  to attain. Comparing with (20), this probability is only attained when  $\sum_{j=1}^{L} 1(u_j = m) = 1$ ,  $\forall m$ . That is, there exists a  $\sigma$ , a permutation of 1:L, such that  $u_1, \dots, u_n = \sigma(1, 2, \dots, L)$ 

Matching the value of  $f_2(x_1, \dots, X)$  we obtain

$$f_2(a_1, \ldots, a) = f_1(a_{\sigma(1)}, \ldots, a_{(L)}).$$

Because the choice of vector  $(a_1, \dots, a)$  is arbitrary, we have

$$f_2(x_1,\ldots,X) = f_1(x_{\sigma(1)},\ldots,X_{\sigma(L)}).$$

Because augmentation increases divergence, for the purpose of finding the largest divergence, to produce the largest divergence, we only need to consider an augmented generator that includes all

$$\Phi^{aug}:(y, \theta_{1,\ldots,L}) \nearrow \{((\phi_1, y), t_1), \ldots, ((\phi_r, y), t_K)\}.$$

It is enough to consider the generator of which  $\phi_k = \phi_k(\theta_1, \dots, \theta)$  are K functions of  $(\theta_1, \dots, \theta)$ .

**Lemma 11.** For any augmented generator  $\Phi^{arg}$  satisfies F, the null, there must exists K permutations  $\sigma_1(1:L),\ldots,\sigma_1(1:L)$ , with the convention  $\sigma_1(1:L)=1:L$ , such that

$$\phi_k(\theta_1,\ldots,\theta) = \phi_1(\theta_{\sigma_k(1:(L))}).$$

*Proof.* Use Lemma (10) for (K - 1) times.

**Lemma 12.** For any augmented generator  $\Phi^{arg}$  satisfies F, all feature-label generator can be replaced by a permutation,  $\phi_k(\theta_1, \dots, \theta_k) = (\theta_{\sigma_k(1)}, \dots, \theta_k)$ , while the divergence does not decrease.

*Proof.* From the previous lemma,

$$\phi_k(\theta_1,\ldots,\theta)=\phi_1(\theta_{\sigma_k(1:L)}).$$

The augmented feature is now  $(\phi_1(\theta_{\sigma_k(1:L)}), y)$ , a transformation of  $(\theta_{\sigma_k(1:L)}), y)$ . Using the raw feature  $(\theta_{\sigma_k(1:L)}), y$  keeps the divergence non-decreasing.

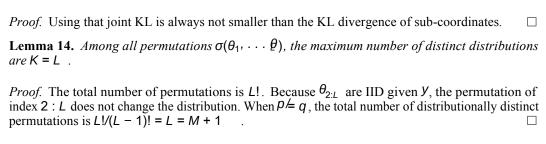
Now that we only want to consider permutation-based generators: there exists  $a^K$ , and K permutations  $\sigma_k(1:L)$ .

$$\Phi^{aug}: (y, \theta_{1,\dots,L}) 7 \rightarrow \{((\phi_1, y), t_1), \dots, ((\phi_t, y), t_k)\}.$$

$$\phi_k(\theta_1, \dots, \theta) = \theta_{\sigma_k(1:L)}.$$

Given a  $p \not\models q$ , any permutations  $\theta_{\sigma_i(1:L)}$  contains one copy from P and (L-1) copies from Q. It suffices to only consider those permutation  $\theta_{\sigma_i(1:L)}$  whose distributions are unique.

**Lemma 13.** Given any  $p \not\models q$  and L fixed, assuming all  $\theta_{\sigma_i(1:L)}$  has different distributions, then  $D(q, p, \Phi^{arg})$  is an increasing function on K.



It is clear that the proposed multi-class permutation  $\Phi^*$  attains this maximum number of distinct distributions, which proves its optimality among all generators from all previous lemmas, thereby the proof of our Theorem 2.