Sampling Permutations for Shapley Value Estimation

Rory Mitchell

Nvidia Corporation Santa Clara CA 95051, USA

Joshua Cooper

Department of Mathematics University of South Carolina 1523 Greene St. Columbia, SC 29223, USA

Eibe Frank

Department of Computer Science University of Waikato Hamilton, New Zealand

Geoffrey Holmes

Department of Computer Science University of Waikato Hamilton, New Zealand

RAMITCHELLNZ@GMAIL.COM

COOPER@MATH.SC.EDU

EIBE@CS.WAIKATO.AC.NZ

GEOFF@CS.WAIKATO.AC.NZ

Abstract

Game-theoretic attribution techniques based on Shapley values are used extensively to interpret black-box machine learning models, but their exact calculation is generally NP-hard, requiring approximation methods for non-trivial models. As the computation of Shapley values can be expressed as a summation over a set of permutations, a common approach is to sample a subset of these permutations for approximation. Unfortunately, standard Monte Carlo sampling methods can exhibit slow convergence, and more sophisticated quasi Monte Carlo methods are not well defined on the space of permutations. To address this, we investigate new approaches based on two classes of approximation methods and compare them empirically. First, we demonstrate quadrature techniques in a RKHS containing functions of permutations, using the Mallows kernel to obtain explicit convergence rates of O(1/n), improving on $O(1/\sqrt{n})$ for plain Monte Carlo. The RKHS perspective also leads to quasi Monte Carlo type error bounds, with a tractable discrepancy measure defined on permutations. Second, we exploit connections between the hypersphere \mathbb{S}^{d-2} and permutations to create practical algorithms for generating permutation samples with good properties. Experiments show the above techniques provide significant improvements for Shapley value estimates over existing methods, converging to a smaller RMSE in the same number of model evaluations.

Keywords: Interpretability, Quasi Monte Carlo, Shapley values

1. Introduction

The seminal work of Shapley (1953) introduces an axiomatic attribution of collaborative game outcomes among coalitions of participating players. Shapley values are popular in machine learning (Cohen et al. (2007); Strumbelj and Kononenko (2010); Štrumbelj and

Kononenko (2014); Lundberg and Lee (2017)) because the assignment of feature relevance to model outputs is structured according to axioms consistent with human notions of attribution. In the machine learning context, each feature is treated as a player participating in the prediction provided by a machine learning model and the prediction is considered the outcome of the game. Feature attributions via Shapley values provide valuable insight into the output of complex models that are otherwise difficult to interpret.

Exact computation of Shapley values is known to be NP-hard in general (Deng and Papadimitriou (1994)) and approximations based on sampling have been proposed by Mann and Shapley (1960); Owen (1972); Castro et al. (2009); Maleki (2015); Castro et al. (2017). In particular, a simple Monte Carlo estimate for the Shapley value is obtained by sampling from a uniform distribution of permutations. An extensively developed Quasi Monte Carlo theory for integration on the unit cube shows that careful selection of samples can improve convergence significantly over random sampling, however, these results do not extend to the space of permutations. Our goal is to better characterise 'good' sample sets for this unique approximation problem, and to develop tractable methods of obtaining these samples, reducing computation time for high-quality approximations of Shapley values. Crucially, we observe that sample evaluations, in this context corresponding to evaluations of machine learning models, dominate the execution time of approximations. Due to the high cost of each sample evaluation, considerable computational effort can be justified in finding such sample sets.

In Section 3, exploiting the direct connection between Shapley values and permutations, we define a reproducing kernel Hilbert space (RKHS) with several possible kernels over permutations. Using these kernels, we apply kernel herding, and sequential Bayesian quadrature algorithms to estimate Shapley values. In particular, we observe that kernel herding, in conjunction with the universal Mallows kernel, leads to an explicit convergence rate of $O(\frac{1}{n})$ as compared to $O(\frac{1}{\sqrt{n}})$ for ordinary Monte Carlo. An outcome of our investigation into kernels is a quasi Monte Carlo type error bound, with a tractable discrepancy formula.

In Section 4, we describe another family of methods for efficiently sampling Shapley values, utilising a convenient isomorphism between the symmetric group \mathfrak{S}_d and points on the hypersphere \mathbb{S}^{d-2} . These methods are motivated by the relative ease of selecting well-spaced points on the sphere, as compared to the discrete space of permutations. We develop two new sampling methods, termed orthogonal spherical codes and Sobol permutations, that select high-quality samples by choosing points well-distributed on \mathbb{S}^{d-2} .

Our empirical evaluation in Section 5 examines the performance of the above methods compared to existing methods on a range of practical machine learning models, tracking the reduction in root mean squared error against exactly calculated Shapley values. Additionally, we evaluate explicit measures of discrepancy (in the quasi Monte Carlo sense) for the sample sets generated by our algorithms. This evaluation of discrepancy for the generated samples of permutations may be of broader interest, as quasi Monte Carlo error bounds based on discrepancy apply to any statistics of functions of permutations and not just Shapley values.

2. Background and Related Work

We first introduce some common notation for permutations and provide the formal definition of Shapley values. Then, we briefly review the literature for existing techniques for approximating Shapley values.

2.1 Notation

We refer to the symmetric group of permutations of d elements as \mathfrak{S}_d . We reserve the use of n to refer to the number of samples. The permutation $\sigma \in \mathfrak{S}_d$ assigns rank j to element i by $\sigma(i) = j$. For example, given the permutation written in one-line notation:

$$\sigma = \begin{pmatrix} 1 & 4 & 2 & 3 \end{pmatrix}$$

and the list of items

$$(x_1, x_2, x_3, x_4)$$

the items are reordered such that x_i occupies the $\sigma(i)$ coordinate

$$(x_1, x_3, x_4, x_2)$$

and the inverse $\sigma^{-1}(j) = i$ is

 $\sigma^{-1} = \begin{pmatrix} 1 & 3 & 4 & 2 \end{pmatrix}.$

2.2 Shapley Values

Shapley values (Shapley (1953)) provide a mechanism to distribute the proceeds of a cooperative game among the members of the winning coalition by measuring marginal contribution to the final outcome. The Shapley value Sh_i for coalition member *i* is defined as

$$\operatorname{Sh}_{i}(v) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|! \, (|N| - |S| - 1)!}{|N|!} (v(S \cup \{i\}) - v(S)) \tag{1}$$

where S is a partial coalition, N is the grand coalition (consisting of all members), and v is the so-called "characteristic function" that is assumed to return the proceeds (i.e., value) obtained by any coalition.

The Shapley value function may also be conveniently expressed in terms of permutations

$$Sh_{i}(v) = \frac{1}{|N|!} \sum_{\sigma \in \mathfrak{S}_{d}} \left[v([\sigma]_{i-1} \cup \{i\}) - v([\sigma]_{i-1}) \right]$$
(2)

where $[\sigma]_{i-1}$ represents the set of players ranked lower than *i* in the ordering σ . The Shapley value is unique and has the following desirable properties:

- 1. Efficiency: $\sum_{i=1}^{n} \operatorname{Sh}_{i}(v) = v(N)$. The sum of Shapley values for each coalition member is the value of the grand coalition N.
- 2. Symmetry: If, $\forall S \subseteq N \setminus \{i, j\}, v(S \cup \{i\}) = v(S \cup \{j\})$, then $\text{Sh}_i = \text{Sh}_j$. If two players have the same marginal effect on each coalition, their Shapley values are the same.

- 3. Linearity: $\text{Sh}_i(v+w) = \text{Sh}_i(v) + \text{Sh}_i(w)$. The Shapley values of sums of games are the sum of the Shapley values of the respective games.
- 4. Dummy: If, $\forall S \subseteq N \setminus \{i\}, v(S \cup \{i\}) = v(S)$, then $\text{Sh}_i = 0$. The coalition member whose marginal impact is always zero has a Shapley value of zero.

Evaluation of the Shapley value is known to be NP-hard in general (Deng and Papadimitriou (1994)) but may be approximated by sampling terms from the sum of either Equation 1 or Equation 2. This paper focuses on techniques for approximating Equation 2 via carefully chosen samples of permutations. We discuss characteristic functions v that arise in the context of machine learning models, with the goal of attributing predictions to input features.

Shapley values have been used as a feature attribution method for machine learning in many prior works (Cohen et al. (2007); Strumbelj and Kononenko (2010); Štrumbelj and Kononenko (2014); Lundberg and Lee (2017)). In the terminology of supervised learning, we have some learned model f(x) = y that maps a vector of features x to a prediction y. The value of the characteristic function is assumed to be given by y, and the grand coalition is given by the full set of features. In a partial coalition, only some of the features are considered "active" and their values made available to the model to obtain a prediction. Applying the characteristic function for partial coalitions requires the definition of $f(x_S)$, where the input features x are perturbed in some way according to the active subset S. A taxonomy of possible approaches is given in Covert et al. (2020). In this paper, we marginalise out features using the joint marginal distribution with some background data set $p(X_S)$, which is the default behaviour of the SHAP python package of Lundberg and Lee (2017).

2.3 Monte Carlo

An obvious Shapley value approximation is the simple Monte Carlo estimator,

$$\bar{Sh}_{i}(v) = \frac{1}{n} \sum_{\sigma \in \Pi} \left[v([\sigma]_{i-1} \cup \{i\}) - v([\sigma]_{i-1}) \right],$$
(3)

for a uniform sample of permutations $\Pi \subset \mathfrak{S}_d$ of size n. Monte Carlo techniques were used to solve electoral college voting games in Mann and Shapley (1960), and a more general analysis is given in Castro et al. (2009). Equation 3 is an unbiased estimator that converges asymptotically to the Shapley value at a rate of $O(1/\sqrt{n})$ according to the Central Limit Theorem.

From a practical implementation perspective, note that a single sample of permutations Π can be used to evaluate Sh_i for all features *i*. For each permutation $\sigma \in \Pi$ of length *d*, first evaluate the empty set $v(\{\})$, then walk through the permutation, incrementing *i* and evaluating $v([\sigma]_i)$, yielding d + 1 function evaluations for σ that are used to construct marginal contributions for each feature. $v([\sigma]_{i-1})$ is not evaluated, but reused from the previous function evaluation, providing a factor of two improvement over the naive approach.

2.4 Antithetic Sampling

Antithetic sampling is a variance reduction technique for Monte Carlo integration where samples are taken as correlated pairs instead of standard i.i.d. samples. The antithetic Monte Carlo estimate (see Rubinstein and Kroese (2016)) is

$$\hat{\mu}_{anti} = \frac{1}{n} \sum_{i=1}^{n/2} f(X_i) + f(Y_i)$$

with variance given by

$$\operatorname{Var}(\hat{\mu}_{anti}) = \frac{\sigma}{n} (1 + \operatorname{Corr}(f(X), f(Y))), \tag{4}$$

such that if f(X) and f(Y) are negatively correlated, the variance is reduced. A common choice for sampling on the unit cube is $X \sim U(0,1)^d$ with $Y_i = 1 - X_i$. Antithetic sampling for functions of permutations is discussed in Lomeli et al. (2019), with a simple strategy being to take permutations and their reverse. We implement this sampling strategy in our experiments with antithetic sampling.

2.5 Multilinear Extension

Another Shapley values approximation is the multilinear extension of Owen (1972). The sum over feature subsets from (1) can be represented equivalently as an integral by introducing a random variable for feature subsets. The Shapley value is calculated as

$$\operatorname{Sh}_{i}(v) = \int_{0}^{1} e_{i}(q) dq \tag{5}$$

where

$$e_i(q) = \mathbb{E}[v(E_i \cup i) - v(E_i)]$$

and E_i is a random subset of features, excluding *i*, where each feature has probability *q* of being selected. $e_i(q)$ is estimated with samples. In our experiments, we implement a version of the multilinear extension algorithm using the trapezoid rule to sample *q* at fixed intervals. A form of this algorithm incorporating antithetic sampling is also presented in Okhrati and Lipani (2020), by rewriting Equation 5 as

$$\operatorname{Sh}_{i}(v) = \int_{0}^{\frac{1}{2}} e_{i}(q) + e_{i}(1-q)dq$$

where the sample set E_i is used to estimate $e_i(q)$ and the 'inverse set', $\{N \setminus \{E_i, i\}\}$, is used to estimate $e_i(1-q)$. In Section 5, we include experiments for the multilinear extension method, both with and without antithetic sampling.

2.6 Stratified Sampling

Another common variance reduction technique is stratified sampling, where the domain of interest is divided into mutually exclusive subregions, an estimate is obtained for each subregion independently, and the estimates are combined to obtain the final estimate. For integral $\mu = \int_{\mathcal{D}} f(x)p(x)dx$ in domain \mathcal{D} , separable into J non-overlapping regions $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_J$ where $w_j = P(X \in \mathcal{D}_j)$ and $p_j(x) = w_j^{-1}p(x)\mathbb{1}_{x \in \mathcal{D}_j}$, the basic stratified sampling estimator is

$$\hat{\mu}_{strat} = \sum_{j=1}^{J} \frac{w_j}{n_j} \sum_{i=1}^{n_j} f(X_{ij}),$$

where $X_{ij} \sim p_j$ for $i = 1, \dots, n_j$ and $j = 1, \dots, J$ (see Owen (2003)). The stratum size n_j can be chosen according to the Neyman allocation if estimates of the variance in each region are known. The stratified sampling method was first applied to Shapley value estimation by Maleki (2015), then improved by Castro et al. (2017). We implement the version in Castro et al. (2017), where strata \mathcal{D}_i^{ℓ} are considered for all $i = 1, \dots, d$ and $\ell = 1, \dots, d$, where \mathcal{D}_i^{ℓ} is the subset of marginal contributions with feature i at position ℓ .

This concludes discussion of existing works; the next sections introduce the primary contributions of this paper.

3. Kernel Methods

A majority of Monte Carlo integration works deal with continuous functions on \mathbb{R}^d where the distribution of samples is well defined. In the space of permutations, distances between samples are implicitly defined, so we impose a similarity metric samples via a kernel and select samples with good distributions relative to these kernels.

Given a positive definite kernel $K : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ over some input space \mathcal{X} , there is an embedding $\phi : \mathcal{X} \to \mathcal{F}$ of elements of \mathcal{X} into a Hilbert space \mathcal{F} , where the kernel computes an inner product $K(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{F}}$ given $x, y \in \mathcal{X}$. Hilbert spaces associated with a kernel are known as reproducing kernel Hilbert spaces (RKHS). Kernels are used extensively in machine learning for learning relations between arbitrary structured data. In this paper, we use kernels over permutations to develop a notion of the quality of finite point sets for the Shapley value estimation problem, and for the optimisation of such point sets. For this task, we investigate three established kernels over permutations: the Kendall, Mallows, and Spearman kernels.

The Kendall and Mallows kernels are defined in Jiao and Vert (2015). Given two permutations σ and σ' of the same length, both kernels are based on the number of concordant and discordant pairs between the permutations:

$$n_{\rm con}(\sigma,\sigma') = \sum_{i< j} [\mathbbm{1}_{\sigma(i)<\sigma(j)} \mathbbm{1}_{\sigma'(i)<\sigma'(j)} + \mathbbm{1}_{\sigma(i)>\sigma(j)} \mathbbm{1}_{\sigma'(i)>\sigma'(j)}]$$
$$n_{\rm dis}(\sigma,\sigma') = \sum_{i< j} [\mathbbm{1}_{\sigma(i)<\sigma(j)} \mathbbm{1}_{\sigma'(i)>\sigma'(j)} + \mathbbm{1}_{\sigma(i)>\sigma(j)} \mathbbm{1}_{\sigma'(i)<\sigma'(j)}]$$

Assuming the length of the permutation is d, the Kendall kernel, corresponding to the well-known Kendall tau correlation coefficient (Kendall (1938)), is

$$K_{\tau}(\sigma, \sigma') = \frac{n_{\rm con}(\sigma, \sigma') - n_{\rm dis}(\sigma, \sigma')}{\binom{d}{2}}.$$

The Mallows kernel, for $\lambda \geq 0$, is defined as

$$K_M^{\lambda}(\sigma, \sigma') = e^{-\lambda n_{\rm dis}(\sigma, \sigma')/{d \choose 2}}$$

Here, the Mallows kernel differs slightly from that of Jiao and Vert (2015). We normalise the $n_{dis(\sigma,\sigma')}$ term relative to d, allowing a consistent selection of the λ parameter across permutations of different length.

While the straightforward implementation of Kendall and Mallows kernels is of order $O(d^2)$, a $O(d \log d)$ variant based on merge-sort is given by Knight (1966).

Note that K_{τ} can also be expressed in terms of a feature map of $\binom{d}{2}$ elements,

$$\Phi_{\tau}(\sigma) = \left(\frac{1}{\sqrt{\binom{d}{2}}} (\mathbbm{1}_{\sigma(i) > \sigma(j)} - \mathbbm{1}_{\sigma(i) < \sigma(j)})\right)_{1 \le i < j \le d}$$

so that

$$K_{\tau}(\sigma, \sigma') = \Phi(\sigma)^T \Phi(\sigma').$$

The Mallows kernel corresponds to a more complicated feature map, although still finite dimensional, given in Mania et al. (2018).

We also define a third kernel based on Spearman's ρ . The (unnormalised) Spearman rank distance

$$d_{\rho}(\sigma, \sigma') = \sum_{i=1}^{d} (\sigma(i) - \sigma'(i))^2 = ||\sigma - \sigma'||_2^2$$

is a semimetric of negative type (Diaconis (1988)), therefore we can exploit the relationship between semimetrics of negative type and kernels from Sejdinovic et al. (2013) to obtain a valid kernel. Writing the product $\sum_{i=0}^{d} \sigma(i)\sigma(i)'$ using vector notation as $\sigma^{T}\sigma'$, we have

$$d(\sigma, \sigma') = K(\sigma, \sigma) + K(\sigma', \sigma') - 2K(\sigma, \sigma')$$
$$d_{\rho}(\sigma, \sigma') = \sigma^{T} \sigma + \sigma'^{T} \sigma' - 2\sigma^{T} \sigma'$$
$$\implies K_{\rho}(\sigma, \sigma') = \sigma^{T} \sigma'$$

and the kernel's feature map is trivially

 $\Phi_{\rho}(\sigma) = \sigma.$

Before introducing sampling algorithms, we derive an additional property for the above kernels: analytic formulas for the expected value evaluated at some fixed point σ and values drawn from a given probability distribution $\sigma' \sim p$. The distribution of interest for approximating (2) is the uniform distribution U. The expected value is straightforward to obtain for the Spearman and Kendall kernel:

$$\forall \sigma \in \Pi, \quad \mathbb{E}_{\sigma' \sim U}[K_{\rho}(\sigma, \sigma')] = \frac{d(d+1)^2}{4}$$
$$\forall \sigma \in \Pi, \quad \mathbb{E}_{\sigma' \sim U}[K_{\tau}(\sigma, \sigma')] = 0$$

The Mallows kernel is more difficult. Let X be a random variable representing the number of inversions over all permutations of length d. Its distribution is studied in Muir (1898), with probability generating function given as

$$\phi_d(x) = \prod_{j=1}^d \frac{1-x^j}{j(1-x)}.$$

There is no convenient form in terms of standard functions for its associated density function. From the probability generating function of X, we obtain the moment generating function:

$$M_d(t) = \phi_d(e^t)$$
$$= \prod_{j=1}^d \frac{1 - e^{tj}}{j(1 - e^t)}$$
$$= \mathbb{E}[e^{tX}]$$

The quantity $n_{\rm dis}(I, \sigma)$, where I is the identity permutation, returns exactly the number of inversions in σ . Therefore, we have

$$M_d(-\lambda/\binom{d}{2}) = \mathbb{E}[e^{-\lambda X/\binom{d}{2}}]$$

= $\mathbb{E}_{\sigma' \sim U}[K_M(I, \sigma')].$

The quantity n_{dis} is right-invariant in the sense that $n_{\text{dis}}(\sigma, \sigma') = n_{\text{dis}}(\tau\sigma, \tau\sigma')$ for $\tau \in \mathfrak{S}_d$ (Diaconis (1988)), so

$$\begin{aligned} \forall \tau \in \mathfrak{S}_d, \quad \mathbb{E}_{\sigma' \sim U}[K_M(I, \sigma')] &= \mathbb{E}_{\sigma' \sim U}[K_M(\tau I, \tau \sigma')] \\ &= \mathbb{E}_{\sigma' \sim U}[K_M(\tau I, \sigma')] \\ \forall \sigma \in \mathfrak{S}_d, \quad \mathbb{E}_{\sigma' \sim U}[K_M(I, \sigma')] &= \mathbb{E}_{\sigma' \sim U}[K_M(\sigma, \sigma')] \\ &= \prod_{j=1}^d \frac{1 - e^{-\lambda j/\binom{d}{2}}}{j(1 - e^{-\lambda j/\binom{d}{2}})}, \end{aligned}$$

We now describe two greedy algorithms for generating point sets improving on simple Monte Carlo—kernel herding and sequential Bayesian quadrature.

3.1 Kernel Herding

A greedy process called "kernel herding" for selecting (unweighted) quadrature samples in a reproducing kernel Hilbert space is proposed in Chen et al. (2010). The sample n + 1 in kernel herding is given by

$$x_{n+1} = \arg\max_{x} \left[\mathbb{E}_{x' \sim p}[K(x, x')] - \frac{1}{n+1} \sum_{i=1}^{n} K(x, x_i) \right]$$
(6)

which can be interpreted as a greedy optimisation process selecting points for maximum separation, while also converging on the expected distribution p. In the case of Shapley value estimation, the samples are permutations $\sigma \in \mathfrak{S}_d$ and p is a uniform distribution with $p(\sigma) = \frac{1}{\sigma!}, \forall \sigma \in \mathfrak{S}_d$.

Kernel herding has time complexity $O(n^2)$ for n samples, assuming the argmax can be computed in O(1) time and $\mathbb{E}_{x'\sim p}[K(x,x')]$ is available. We have analytic formulas for $\mathbb{E}_{x'\sim p}[K(x,x')]$ from the previous section for the Spearman, Kendall, and Mallows kernels, and they give constant values depending only on the size of the permutation d. We compute an approximation to the argmax in constant time by taking a fixed number of random samples at each iteration and retaining the one yielding the maximum.

If certain conditions are met, kernel herding converges at the rate $O(\frac{1}{n})$, an improvement over $O(\frac{1}{\sqrt{n}})$ for standard Monte Carlo sampling. According to Chen et al. (2010), this improved convergence rate is achieved if the RKHS is universal, and mild assumptions are satisfied by the argmax (it need not be exact). Of the Spearman, Kendall and Mallows kernels, only the Mallows kernel has the universal property (Mania et al. (2018)).

Next, we describe a more sophisticated kernel-based algorithm generating weighted samples.

3.2 Sequential Bayesian Quadrature

Bayesian Quadrature (O'Hagan (1991); Rasmussen and Ghahramani (2003)) (BQ) formulates the integration problem,

$$Z_{f,p} = \int f(x)p(x)dx$$

as a Bayesian inference problem. Standard BQ imposes a Gaussian process prior on f with zero mean and kernel function K. A posterior distribution is inferred over f conditioned on a set of points (x_0, x_1, \dots, x_n) . This implies a distribution on $Z_{f,p}$ with expected value

$$\mathbb{E}_{GP}[Z] = z^T K^{-1} f(X)$$

where f(X) is the vector of function evaluations at points (x_0, x_1, \dots, x_n) , K^{-1} is the inverse of the kernel covariance matrix, and $z_i = \mathbb{E}_{x' \sim p}[K(x_i, x')]$. Effectively, for an arbitrary set of points, Bayesian quadrature solves the linear system Kw = z to obtain a reweighting of the sample evaluations, yielding the estimate

$$Z \simeq w^T f(X).$$

An advantage of the Bayesian approach is that uncertainty is propagated through to the final estimate. Its variance is given by

$$\mathbb{V}[Z_{f,p}|f(X)] = \mathbb{E}_{x,x' \sim p}[K(x,x')] - z^T K^{-1} z.$$
(7)

This variance estimate is used in Huszár and Duvenaud (2012) to develop sequential Bayesian quadrature (SBQ), a greedy algorithm selecting samples to minimise Equation 7. This procedure, summarised in Algorithm 1, is shown by Huszár and Duvenaud (2012) to be related to optimally weighted kernel herding.

SBQ has time complexity $O(n^3)$ for n samples if the argmin takes constant time, and an $O(n^2)$ Cholesky update algorithm is used to form K^{-1} , adding one sample at a time. In general, exact minimisation of Equation 7 is not tractable, so as with kernel herding, we approximate the argmin by drawing a fixed number of random samples and choosing the one yielding the minimum variance.

3.3 Error Analysis in RKHS

Canonical error analysis of quasi Monte-Carlo quadrature is performed using the Koksma-Hlawka inequality (Hlawka (1961); Niederreiter (1992)), decomposing error into a product Algorithm 1: Sequential Bayesian Quadrature

Input: n, kernel K, sampling distribution p, integrand f1 $X_0 \leftarrow RandomSample(p)$ 2 $K^{-1} = I$ // Inverse of covariance matrix $\mathbf{s} \ z_0 \leftarrow \mathbb{E}_{x' \sim p}[K(X_0, x')]$ 4 for $i \leftarrow 2$ to n do $X_i \leftarrow \arg\min \mathbb{E}_{x,x' \sim p}[K(x,x')] - z^T K^{-1} z$ $y \leftarrow \vec{0}$ 6 for $j \leftarrow 1$ to i do 7 $y_j = K(X_i, X_j)$ 8 $K^{-1} \leftarrow CholeskyUpdate(K^{-1}, y)$ 9 $z_i \leftarrow \mathbb{E}_{x' \sim p}[K(X_i, x')]$ 10 11 $w = z^T K^{-1}$ 12 return $w^T f(X)$

of function variation and discrepancy of the sample set. We derive a version of this inequality for Shapley value approximation in terms of reproducing kernel Hilbert spaces. Our derivation mostly follows Hickernell (2000), with modification of standard integrals to weighted sums of functions on \mathfrak{S}_d , allowing us to calculate discrepancies for point sets generated by kernel herding and SBQ with permutation kernels. The analysis is performed for the Mallows kernel, which is known to be a universal kernel (Mania et al. (2018)).

Given a symmetric, positive definite kernel K, we have a unique RKHS \mathcal{F} with inner product $\langle \cdot, \cdot \rangle_K$ and norm $|| \cdot ||_K$, where the kernel reproduces functions $f \in \mathcal{F}$ by

$$f(\sigma) = \langle f, K(\cdot, \sigma) \rangle_K.$$

Define error functional

$$\operatorname{Err}(f,\Pi,w) = \frac{1}{d!} \sum_{\sigma \in \mathfrak{S}_d} f(\sigma) - \sum_{\tau \in \Pi} w_{\tau} f(\tau),$$

where Π is a sample set of permutations and w_{τ} is the associated weight of sample τ . Because the Mallows kernel is a universal kernel, the Shapley value component functions $f(\sigma)$ belong to \mathcal{F} . Given that $\operatorname{Err}(f, \Pi, w)$ is a continuous linear functional on \mathcal{F} and assuming that it is bounded, by the Riesz Representation Theorem, there is a function $\xi \in \mathcal{F}$ that is its representer: $\operatorname{Err}(f, \Pi, w) = \langle \xi, f \rangle_K$. Using the Cauchy-Schwarz inequality, the quadrature error is bounded by

$$|\operatorname{Err}(f,\Pi,w)| = |\langle \xi, f \rangle_K| \le ||\xi||_K ||f||_K = D(\Pi,w)V(f)$$

where $D(\Pi, w) = ||\xi||_K$ is the discrepancy of point set Π with weights w and $V(f) = ||f||_K$ is the function variation. The quantity $D(\Pi, w)$ has an explicit formula. As the function ξ is reproduced by the kernel, we have:

$$\begin{aligned} \xi(\sigma') &= \langle \xi, K(\cdot, \sigma') \rangle_K = \operatorname{Err}(K(\cdot, \sigma'), \Pi, w) \\ &= \frac{1}{d!} \sum_{\sigma \in \mathfrak{S}_d} K(\sigma, \sigma') - \sum_{\tau \in \Pi} w_\tau K(\tau, \sigma'). \end{aligned}$$

Then the discrepancy can be obtained, using the fact that $\operatorname{Err}(f,\Pi,w) = \langle \xi, f \rangle_K$, by

$$D(\Pi, w) = ||\xi||_{k} = \sqrt{\langle \xi, \xi \rangle_{K}} = \sqrt{\operatorname{Err}(\xi, \Pi, w)}$$

$$= \left(\frac{1}{d!} \sum_{\sigma \in \mathfrak{S}_{d}} \xi(\sigma) - \sum_{\tau \in \Pi} w_{\tau} \xi(\tau)\right)^{\frac{1}{2}}$$

$$= \left(\frac{1}{d!} \sum_{\sigma \in \mathfrak{S}_{d}} \left[\frac{1}{d!} \sum_{\sigma' \in \mathfrak{S}_{d}} K(\sigma, \sigma') - \sum_{\tau \in \Pi} w_{\tau} K(\tau, \sigma)\right] - \sum_{\tau \in \Pi} w_{\tau} \left[\frac{1}{d!} \sum_{\sigma \in \mathfrak{S}_{d}} K(\sigma, \tau) - \sum_{\tau' \in \Pi} w_{\tau'} K(\tau, \tau')\right]\right)^{\frac{1}{2}}$$

$$= \left(\frac{1}{(d!)^{2}} \sum_{\sigma, \sigma' \in \mathfrak{S}_{d}} K(\sigma, \sigma') - \frac{2}{d!} \sum_{\sigma \in \mathfrak{S}_{d}} \sum_{\tau \in \Pi} w_{\tau} K(\tau, \sigma) + \sum_{\tau, \tau' \in \Pi} w_{\tau} w_{\tau'} K(\tau, \tau')\right)^{\frac{1}{2}}$$

$$= \left(\mathbb{E}_{\sigma, \sigma' \sim U} [K(\sigma, \sigma')] - 2 \sum_{\tau \in \Pi} w_{\tau} \mathbb{E}_{\sigma \sim U} [K(\tau, \sigma)] + \sum_{\tau, \tau' \in \Pi} w_{\tau} w_{\tau'} K(\tau, \tau')\right)^{\frac{1}{2}}$$
(8)

It can be seen that kernel herding (Equation 6) greedily minimises $D(\Pi, w)^2$ with constant weights $\frac{1}{n}$, by examining the reduction in $D(\Pi, \frac{1}{n})^2$ by addition of sample π to Π . The kernel herding algorithm for sample $\sigma_{n+1} \in \Pi$ is

$$\sigma_{n+1} = \underset{\sigma}{\arg\max} \left[\mathbb{E}_{\sigma' \sim U}[K(\sigma, \sigma')] - \frac{1}{n+1} \sum_{i=1}^{n} K(\sigma, \sigma_i) \right]$$

Note that, since $K(\cdot, \cdot)$ is right-invariant, the quantity $\mathbb{E}_{\sigma' \sim U}[K(\sigma, \sigma')]$ does not depend on σ , so the arg max above is simply minimizing $\sum_{i=1}^{n} K(\sigma, \sigma_i)$. On the other hand, denoting the identity permutation by I,

$$\begin{split} D(\Pi, \frac{1}{n})^2 - D(\Pi \cup \{\pi\}, \frac{1}{n+1})^2 &= 2 \sum_{\tau \in \Pi \cup \{\pi\}} \frac{1}{n+1} \mathbb{E}_{\sigma \sim U}[K(\tau, \sigma)] - 2 \sum_{\tau \in \Pi} \frac{1}{n} \mathbb{E}_{\sigma \sim U}[K(\tau, \sigma)] \\ &+ \sum_{\tau, \tau' \in \Pi} \frac{1}{n^2} K(\tau, \tau') - \sum_{\tau, \tau' \in \Pi \cup \{\pi\}} \frac{1}{(n+1)^2} K(\tau, \tau') \\ &= 2 \frac{n+1}{n+1} \mathbb{E}_{\sigma \sim U}[K(I, \sigma)] - 2 \frac{n}{n} \mathbb{E}_{\sigma \sim U}[K(I, \sigma)] \\ &+ \sum_{\tau, \tau' \in \Pi} \frac{2n+1}{n^2(n+1)^2} K(\tau, \tau') - 2 \sum_{\tau \in \Pi} \frac{1}{(n+1)^2} K(\tau, \pi) \\ &= \frac{K(I, I)}{(n+1)^2} + \sum_{\tau, \tau' \in \Pi} \frac{2n+1}{n^2(n+1)^2} K(\tau, \tau') \\ &- \frac{2}{(n+1)^2} \sum_{\tau \in \Pi} K(\tau, \pi) \end{split}$$

where both equalities use right-invariance. Note that the first two summands in the last expression are constants (i.e., do not depend on the choice of π), so maximizing this quantity is the same as minimizing $\sum_{\tau \in \Pi} K(\tau, \pi)$, i.e., the same as the kernel herding optimization subproblem.

Furthermore, we can show Bayesian quadrature minimises squared discrepancy via optimisation of weights. Writing $z_i = \mathbb{E}_{\sigma' \sim p}[K(\sigma_i, \sigma')]$ and switching to vector notation we have

$$D(\Pi, w)^2 = c - 2w^T z + w^T K w$$

where the first term is a constant not depending on w. Taking the gradient with respect to w, setting it to 0, and solving for w, we obtain:

$$\nabla D(\Pi, w)^2 = -2z + 2w^T K = 0$$

$$w^* = z^T K^{-1},$$
 (9)

where (9) is exactly line 11 of Algorithm 1.

We use the discrepancy measure in (8) for numerical experiments in Section 5.2 to determine the quality of a set of sampled permutations in a way that is independent of the integrand f.

4. Sampling Permutations on \mathbb{S}^{d-1}

Kernel herding and sequential Bayesian quadrature directly reduce the discrepancy of the sampled permutations via greedy optimisation procedures. We now describe two approaches to sampling permutations of length d based on a relaxation to the Euclidean sphere $\mathbb{S}^{d-2} = \{x \in \mathbb{R}^{d-1} : ||x|| = 1\}$, where the problem of selecting well distributed samples is simplified. We describe a simple procedure for mapping points on the surface of this hypersphere to the nearest permutation, where the candidate nearest neighbours form the vertices of a Cayley graph inscribing the sphere. This representation provides a natural connection between distance metrics over permutations, such as Kendall's tau and Spearman's rho, and Euclidean space. We show that samples taken uniformly on the sphere result in a uniform distribution over permutations, and evaluate two unbiased sampling algorithms. Our approach here is closely related to Plis et al. (2010), where an angular view of permutations is used to solve inference problems.

4.1 Spheres, Permutahedrons, and the Cayley Graph

Consider the projection of permutations $\sigma \in \mathfrak{S}_d$ as points in \mathbb{R}^d , where the *i*-th coordinate is given by $\sigma^{-1}(i)$. These points form the vertices of a polytope known as the permutohedron (Guilbaud and Rosenstiehl (1963)). The permutohedron is a d-1 dimensional object embedded in d dimensional space, lying on the hyperplane given by

$$\sum_{i=1}^{d} \sigma^{-1}(i) = \frac{d(d+1)}{2},$$



[4 12 3] [4 12 3] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [4 13 2] [3 1 12] [3 1 12] [3 1 12] [3 1 12] [3 1 12] [3 1 12] [3 1 12] [3 1 12] [3 1 12] [3 1 12] [3 1 12] [4 13 2] [4 13 2] [3 1 12] [4 13 2] [3 1 12] [4 13 2] [3 1 12] [4 13 2] [3 1 12] [4 13 2] [3 1 12] [4 13 1] [4 13 2] [3 1 12] [4 13 1] [4 13 1] [4 13 2] [4 13 1] [4 13 2] [4 13 1] [4 13 2] [3 1 12] [3 1 12] [3 1 12] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 1 12] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 12 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1] [4 1]

Figure 1: Cayley Graph of d = 3



with normal vector

$$\vec{n} = \begin{bmatrix} \frac{1}{\sqrt{d}} \\ \frac{1}{\sqrt{d}} \\ \vdots \\ \frac{1}{\sqrt{d}} \end{bmatrix}, \qquad (10)$$

and inscribing the hypersphere \mathbb{S}^{d-2} lying on the hyperplane, defined by

$$\sum_{i=1}^{d} \sigma^{-1}(i)^2 = \frac{d(d+1)(2d+1)}{6}.$$

Inverting the permutations at the vertices of the permutohedron gives a Cayley graph of the symmetric group with adjacent transpositions as the generating set. Figure 1 shows the Cayley graph for \mathfrak{S}_3 , whose vertices form a hexagon inscribing a circle on a hyperplane, and Figure 2 shows the Cayley graph of \mathfrak{S}_4 projected into 3 dimensions (its vertices lie on a hyperplane in four dimensions). Each vertex σ^{-1} in the Cayley graph has d-1neighbours, where each neighbour differs by exactly one adjacent transposition (one bubblesort operation). Critically for our application, this graph has an interpretation in terms of distance metrics on permutations. The Kendall-tau distance is the graph distance in the 1skeleton of this polytope, and Spearman distance is the squared Euclidean distance between two vertices (Thompson (1993)). Additionally, the antipode of a permutation is its reverse permutation. With this intuition, we use the hypersphere as a continuous relaxation of the space of permutations, where selecting samples far apart on the hypersphere corresponds to sampling permutations far apart in the distance metrics of interest.

We now describe a process for sampling from the set of permutations inscribing \mathbb{S}^{d-2} . First, shift and scale the permutohedron to lie around the origin with radius r = 1. The transformation on vertex σ^{-1} is given by

$$\hat{\sigma}^{-1} = \frac{\sigma^{-1} - \mu}{||\sigma^{-1}||},\tag{11}$$

where $\mu = (\frac{d+1}{2}, \frac{d+1}{2}, \cdots)$ is the mean vector of all permutations.

Now select some vector x of dimension d-1, say, uniformly at random from the surface of \mathbb{S}^{d-2} . Project x onto the hyperplane in \mathbb{R}^d using the following matrix $d-1 \times d$ matrix:

$$U = \begin{bmatrix} 1 & -1 & 0 & \dots & 0 \\ 1 & 1 & -2 & \dots & 0 \\ \vdots & & \ddots & \\ 1 & 1 & 1 & \dots & -(d-1) \end{bmatrix}$$

It is easily verifiable that this basis of row vectors is orthogonal to hyperplane normal \vec{n} . Normalising the row vectors of U gives a transformation matrix \hat{U} used to project vector x to the hyperplane by

 $\tilde{x} = \hat{U}^T x.$

so that

$$\tilde{x}^T \vec{n} = 0$$

Given \tilde{x} , find the closest permutation $\hat{\sigma}^{-1}$ by maximising the inner product

$$\hat{y} = \operatorname*{arg\,max}_{\hat{\sigma}^{-1}} \tilde{x}^T \hat{\sigma}^{-1}.$$
(12)

This maximisation is simplified by noting that $\hat{\sigma}^{-1}$ is always a reordering of the same constants ($\hat{\sigma}^{-1}$ is a scaled and shifted permutation). The inner product is therefore maximised by matching the largest element in $\hat{\sigma}^{-1}$ against the largest element in \tilde{x} , then proceeding to the second-largest, and so on. Thus the argmax is performed by finding the permutation corresponding to the order type of \tilde{x} , which is order-isomorphic to the coordinates of \tilde{x} . The output \hat{y} is a vertex on a scaled permutohedron - to get the corresponding point on the Cayley graph, undo the scale/shift of Eq. 11 to get a true permutation, then invert that permutation:

$$y = \operatorname{inverse}(\hat{y}||\sigma^{-1}|| + \mu).$$
(13)

In fact, both Eq. 12 and 13 can be simplified via a routine *argsort*, defined by

$$\operatorname{argsort}(a) = b$$

such that

$$a_{b_0} \le a_{b_1} \le \dots \le a_{b_n}$$

In other words, b contains the indices of the elements of a in sorted position.

Algorithm 2 describes the end-to-end process of sampling. We use the algorithm of Knuth (1997) for generating points uniformly at random on \mathbb{S}^{d-2} : sample from d-1 independent Gaussian random variables and normalise the resulting vector to have unit length. We now make the claim that Algorithm 2 is unbiased.

Theorem 1 Algorithm 2 generates permutations uniformly at random, i.e., $Pr(\sigma) = \frac{1}{d!}, \forall \sigma \in \mathfrak{S}_d$, from a uniform random sample on \mathbb{S}^{d-2} .

Algorithm 2: Sample permutation from \mathbb{S}^{d-2}					
Output: σ , a permutation of length d					
$1 \ x \leftarrow N(0, 1)$	// x is a vector of $d-1$ i.i.d. normal samples				
2 $x \leftarrow \frac{x}{ x }$	// x lies uniformly on \mathbb{S}^{d-2}				
з $ ilde{x} = \hat{U}^T x$					
4 $\sigma \leftarrow \operatorname{argsort}(\tilde{x})$	// σ is a uniform random permutation				

Proof The point $x \in \mathbb{S}^{d-2}$ from Algorithm 2, Line 2, has multivariate normal distribution with mean 0 and covariance $\Sigma = aI$ for some scalar a and I as the identity matrix. $\tilde{x} = \hat{U}^T x$ is an affine transformation of a multivariate normal and so has covariance

$$Cov(\tilde{\mathbf{x}}) = \hat{U}^T \Sigma \hat{U}$$
$$= a \hat{U}^T I \hat{U}$$
$$= a \hat{U}^T \hat{U}$$

The $d \times d$ matrix $\hat{U}^T \hat{U}$ has the form

$$\hat{U}^T \hat{U} = \begin{bmatrix} \frac{d-1}{d} & \frac{-1}{d} & \dots & \frac{-1}{d} \\ \frac{-1}{d} & \frac{d-1}{d} & \dots & \frac{-1}{d} \\ & \vdots & \ddots & \\ \frac{-1}{d} & \frac{-1}{d} & \dots & \frac{d-1}{d} \end{bmatrix}$$

with all diagonal elements $\frac{d-1}{d}$ and off diagonal elements $\frac{-1}{d}$, and so \tilde{x} is equicorrelated. Due to equicorrelation, \tilde{x} has order type such that $\forall \tilde{x}_i, \tilde{x}_j \in x, i \neq j : Pr(\tilde{x}_i < \tilde{x}_j) = \frac{1}{2}$. In other words, all orderings of \tilde{x} are equally likely. The function *argsort* implies an order-isomorphic bijection, that is, *argsort* returns a unique permutation for every unique ordering over its input. As every ordering of \tilde{x} is equally likely, Algorithm 2 outputs permutations $\sigma \in \mathfrak{S}_d$ with $p(\sigma) = \frac{1}{d!}, \forall \sigma \in \mathfrak{S}_d$.

Furthermore, Equation 12 associates a point on the surface of \mathbb{S}^{d-2} to the nearest permutation. This implies that there is a Voronoi cell on the same surface associated with each permutation σ_i , and a sample \tilde{x} is associated with σ_i if it lands in its cell. Figure 3 shows the Voronoi cells on the hypersphere surface for d = 4, where the green points are equidistant from nearby permutations. A corollary of Theorem 1 is that these Voronoi cells must have equal measure, which is easily verified for d = 4.

4.2 Orthogonal Spherical Codes

Having established an order isomorphism $\mathbb{S}^{d-2} \to \mathfrak{S}_d$, we consider selecting well-distributed points on \mathbb{S}^{d-2} . Our first approach, described in Algorithm 3, is to select 2(d-1) correlated samples on \mathbb{S}^{d-2} from a basis of orthogonal vectors. Algorithm 3 uses the Gram-Schmidt process to incrementally generate a random basis, then converts each component and its reverse into permutations by the same mechanism as Algorithm 2. The cost of each additional sample is proportional to $O(d^2)$. This sampling method is related to orthogonal Monte Carlo

Algorithm 3: Sample k = 2(d-1) permutations from d-21 $X \sim N(0,1)_{k/2,d}$ // iid. normal random Matrix2 $Y \leftarrow 0_{k,d}$ // Matrix storing output permutations3 for $i \leftarrow 1$ to k/2 do// Matrix storing output permutations4for $j \leftarrow 1$ to i do5 $\lfloor X_i \leftarrow X_i - X_j X_i^T \cdot X_j$ 6 $X_i \leftarrow \frac{X_i}{||X_i||}$

techniques discussed in Choromanski et al. (2019). Writing $v([\sigma]_{i-1} \cup \{i\}) - v([\sigma]_{i-1}) = g_i(\sigma)$, the Shapley value estimate for samples given by Algorithm 3 is

$$\bar{\mathrm{Sh}}_{i}^{\mathrm{orth}}(v) = \frac{1}{n} \sum_{\ell=1}^{n/k} \sum_{j=1}^{k} g_{i}(\sigma_{\ell j}),$$
 (14)

where $(\sigma_{\ell 1}, \sigma_{\ell 2}, \cdots, \sigma_{\ell k})$ are a set of correlated samples and n is a multiple of k.

Proposition 1 $\bar{Sh}_i^{orth}(v)$ is an unbiased estimator of $Sh_i(v)$.

 $Y_{2i} \leftarrow \operatorname{argsort}(\hat{U}^T X_i)$

 $Y_{2i+1} \leftarrow \operatorname{argsort}(\hat{U}^T(-X_i))$

7

8

9 return Y

Proof The Shapley value $Sh_i(v)$ is equivalently expressed as an expectation over uniformly distributed permutations:

$$\operatorname{Sh}_{i}(v) = \frac{1}{|N|!} \sum_{\sigma \in \mathfrak{S}_{d}} \left[v([\sigma]_{i-1} \cup \{i\}) - v([\sigma]_{i-1}) \right]$$

$$\operatorname{Sh}_{i}(v) = \mathbb{E}_{\sigma \sim U}[g_{i}(\sigma)]$$

The distribution of permutations drawn as orthogonal samples is clearly symmetric, so $p(\sigma_{\ell,j}) = p(\sigma_{\ell,m})$ for any two indices j, m in a set of k samples, and $\mathbb{E}[g_i(\sigma_{\ell,j})] = \mathbb{E}[g_i(\sigma_{\ell,m}))] = \mathbb{E}[g_i(\sigma^{ortho})]$. As the estimator (14) is a sum, by the linearity of expectation

$$\mathbb{E}[\bar{\mathrm{Sh}}_{i}^{\mathrm{orth}}(v)] = \frac{1}{n} \sum_{\ell=1}^{n/k} \sum_{j=1}^{k} \mathbb{E}[g_{i}(\sigma_{\ell j})] = \mathbb{E}[g_{i}(\sigma^{ortho})].$$

By Theorem 1, the random variable σ^{ortho} has a uniform distribution if its associated sample $x \in \mathbb{S}_{d-2}$ is drawn with uniform distribution. Let x be a component of a random orthogonal basis. If the random basis is drawn with equal probability from the set of orthogonal matrices of order d - 1 (i.e. with Haar distribution for the orthogonal group), then it follows that $\mathbb{E}[g_i(\sigma^{ortho})] = \mathbb{E}_{\sigma \sim U}[g_i(\sigma)]$. The Gram-Schmidt process applied to a square matrix with elements as i.i.d. standard normal random variables yields a random orthogonal matrix with Haar distribution (Mezzadri (2006)). Therefore

$$Sh_i(v) = \mathbb{E}_{\sigma \sim U}[g_i(\sigma)] = \mathbb{E}_{\sigma \sim U}[g_i(\sigma)]$$
$$= \mathbb{E}[Sh_i^{orth}(v)]$$

The variance of the estimator (14) can be analysed similarly to the antithetic sampling of Section 2.4 to k correlated random variables. By extension of the antithetic variance in Equation 4, we have

$$\operatorname{Var}(\bar{\operatorname{Sh}}_{i}^{\operatorname{orth}}(v)) = \frac{1}{n} \sum_{\ell=1}^{n/k} \sum_{j,m=1}^{k} \operatorname{Cov}(g(\sigma_{\ell j}), g(\sigma_{\ell m})).$$

The variance is therefore minimised by selecting k negatively correlated samples. Our experimental evaluation in Section 5 suggests that, for the domain of interest, orthogonal samples on the sphere are indeed strongly negatively correlated, and the resulting estimators are more accurate than standard Monte Carlo and antithetic sampling in all evaluations.

Samples from Algorithm 3 can also be considered as a type of spherical code. Spherical codes describe configurations of points on the unit sphere maximising the angle between any two points (see Conway et al. (1987)). A spherical code $A(n, \phi)$ gives the maximum number of points in dimension n with minimum angle ϕ . The orthonormal basis and its antipodes trivially yield the optimal code $A(d-1, \frac{\pi}{2}) = 2(d-1)$.

From their relative positions on the Cayley graph we obtain bounds on the Kendall tau kernel $K_{\tau}(\sigma, \sigma')$ from Section 3 for the samples of Algorithm 3. The angle between vertices of the Cayley graph is related to $K_{\tau}(\sigma, \sigma')$ in that the maximum kernel value of 1 occurs for two permutations at angle 0 and the minimum kernel value of -1 occurs for a permutation and its reverse, separated by angle π . As the angle between two points (x, x')on \mathbb{S}_{d-2} increases from 0 to π , the kernel $K_{\tau}(\sigma, \sigma')$ for the nearest permutations (σ, σ') increases monotonically and linearly with the angle, aside from quantisation error. If the angle between two distinct points (x, x') in our spherical codes is $\frac{\pi}{2}$, we obtain via the map, $\mathbb{S}^{d-2} \to \mathfrak{S}_d$, the permutations (σ, σ') such that

$$|K_{\tau}(\sigma, \sigma')| \le 1/2 + \epsilon,$$

with some small constant quantisation error ϵ . Figure 4 shows k = 6 samples for the d = 4case. This is made precise in the following result. Note that the statement and its proof are in terms of σ and σ' instead of their inverses (which label the vertices of the permutohedron in our convention), for simplicity; without this change, the meaning is the same, since $n_{\rm dis}(\sigma,\sigma') = n_{\rm dis}(\sigma^{-1},\sigma'^{-1})$ and $A(\sigma)^T A(\sigma') = A(\sigma^{-1})^T A(\sigma'^{-1})$ for any permutations σ , σ'. First, let $\rho = \sqrt{d(d^2 - 1)/12}$, so that the map $A(y) = (y - \mu)/\rho$ maps the permutohedron to an isometric copy of \mathbb{S}^{d-2} centered at the origin in \mathbb{R}^d , the intersection of the unit sphere \mathbb{S}^{d-1} with the hyperplane orthogonal to \vec{n} .

Theorem 2 Suppose $\sigma, \sigma' \in \mathfrak{S}_d$. Then

$$-2+4\left(\frac{1-K_{\tau}(\sigma,\sigma')}{2}\right)^{3/2} \le A(\sigma)^{T}A(\sigma') - 3K_{\tau}(\sigma,\sigma') + O(d^{-1}) \le 2-4\left(\frac{1+K_{\tau}(\sigma,\sigma')}{2}\right)^{3/2}$$

and if $A(\sigma)^{T}A(\sigma') = o(1)$ then

and, if $A(\sigma)^{*} A(\sigma') = o(1)$, then

$$|K_{\tau}(\sigma, \sigma')| \le 1/2 + o(1).$$



Figure 3: Voronoi cells on n-sphere

Figure 4: Orthogonal spherical codes

The above result is a kind of converse to the so-called Rearrangement Inequality, which states that the maximum dot product between a vector and a vector consisting of any permutation of its coordinates is maximized when the permutation is the identity and minimized when it is the reverse identity. Here, we show what happens in between: as one varies from the identity to its reverse one adjacent transposition at a time, the dot product smoothly transitions from maximal to minimal, with some variability across permutations having the same number of inversions. Interestingly, we do not know if the above bound is best possible. A quick calculation shows that, letting $k \approx d2^{-1/3}$ be an integer, the permutation

$$\pi = (k, k - 1, \dots, 2, 1, k + 1, k + 2, \dots, d - 1, d)$$

has $\nu(\pi) = I^T \pi = d^3(1/4 + o(1))$, i.e, $A(I)^T A(\pi) \approx 0$. However, π admits $d^2(2^{-5/3} + o(1))$ inversions, whence $K_{\tau}(I, \pi) \approx 1 - 2^{-2/3} \approx 0.37 < 1/2$.

Figure 5 shows the distribution of pairs of unique samples taken from random vectors, versus unique samples from an orthogonal basis, at d = 10. Samples corresponding to orthogonal vectors are tightly distributed around $K_{\tau}(\sigma, \sigma') = 0$, and pairs corresponding to a vector and its antipodes are clustered at $K_{\tau}(\sigma, \sigma') = -1$. Figure 6 plots the bounds from Theorem 2 relating the dot product of vectors on \mathbb{S}^{d-2} to the Kendall tau kernel at d = 15.

4.3 Sobol Sequences on the Sphere

We now describe another approach to sampling permutations via \mathbb{S}_{d-2} , based on standard quasi Monte Carlo techniques. Low discrepancy point sets on the unit cube \mathbb{R}^{d-2} may be projected to \mathbb{S}^{d-2} via area preserving transformations. Such projections are discussed in depth in Brauchart and Dick (2012); Hardin et al. (2016), where they are observed to have good properties for numerical integration. Below we define transformations in terms of the inverse cumulative distribution of the generalised polar coordinate system and use transformed high-dimensional Sobol sequences to obtain well-distributed permutations.

In the generalised polar coordinate system of Blumenson (1960), a point on \mathbb{S}^{d-2} is defined by radius r and d-2 angular coordinates $(r, \varphi_1, \varphi_2, \cdots, \varphi_{d-2})$, where $(\varphi_1, \cdots, \varphi_{d-3})$ range from $[0, \pi]$ and φ_{d-2} ranges from $[0, 2\pi]$.



Figure 5: Sampling random pairs

Figure 6: Bounds on $K_{\tau}(I, \sigma)$

The polar coordinates on the sphere are independent and have probability density functions

$$f(\varphi_{d-2}) = \frac{1}{2\pi}$$

and for $1 \leq j < d-2$:

$$f(\varphi_j) = \frac{1}{B(\frac{d-j-1}{2}, \frac{1}{2})} \sin^{(d-j-2)}(\varphi_j).$$

where B is the beta function. The above density function is obtained by normalising the formula for the surface area element of a hypersphere to integrate to 1 (Blumenson (1960)). The cumulative distribution function for the polar coordinates is then

$$F_j(\varphi_j) = \int_0^{\varphi_j} f_j(u) du.$$

As per standard inverse transform sampling, we draw samples $x \in [0, 1]^{d-2}$ uniformly from the unit cube and project them to polar coordinates uniformly distributed on the sphere as $\varphi_j = F_j^{-1}(x_j)$. F_j^{-1} can be obtained quickly via a root finding algorithm, such as the bracketing method described in Press et al. (2007).

The points $x \in [0, 1]^{d-2}$ are generated using the Sobol sequence (Sobol' (1967)), also referred to as (t, s)-sequences in base 2. Analogously to our discrepancy for functions of permutations in Equation 8, derived with the Mallows kernel, Sobol points can be shown to minimise a discrepancy for the kernel

$$K(x, x') = \prod_{i=1}^{d} \min(1 - x_j, 1 - x'_j)$$

with $x, x' \in [0, 1]^d$, where the discrepancy decreases at the rate $O(\frac{(\log n)^d}{n})$ (see Dick and Pillichshammer (2010)). Sobol points are relatively inexpensive to generate compared with other algorithms discussed in this paper, although explicit convergence rates for discrepancy on the cube do not translate to \mathbb{S}^{d-2} or \mathfrak{S}_d .

Algorithm 4: Sobol Permutations

```
1 Function PolarToCartesian((r, \varphi_1, \varphi_2, \cdots, \varphi_{d-2})):
          Output: \vec{x}
          for i \leftarrow 1 to d - 1 do
 \mathbf{2}
               x_i \leftarrow r
 з
               for j \leftarrow 1 to i - 1 do
 4
                    x_i \leftarrow x_i \sin \varphi_j
 5
               if i < d - 2 then
 6
 7
                     x_i \leftarrow x_i \cos \varphi_i
 8
          return x
 9
    Function SobolPermutations (n, d):
10
          Output: \Pi
          for i \leftarrow 1 to n do
11
               x \leftarrow \text{SobolPoint}(i, n, d)
                                                                                                // x has d-2 elements
12
               \varphi \leftarrow \vec{0}
13
               for j \leftarrow 1 to d - 2 do
14
                \varphi_j \leftarrow F_i^{-1}(x_j)
                                                                                     // Inverse CDF transformation
15
               y \leftarrow \texttt{PolarToCartesian}(1, \varphi)
                                                                                                // y has d-1 elements
16
               z \leftarrow \hat{U}^T u
                                                                                                      // z has d elements
17
               \Pi_i \leftarrow \operatorname{argsort}(z)
18
          return Π
19
20
```

Combining Sobol points with inverse transform sampling yields uniformly distributed points on \mathbb{S}^{d-2} . To map these points to permutations, we project from \mathbb{R}^{d-1} to the hyperplane in \mathbb{R}^d containing the permutahedron (such that points are orthogonal to the normal in Eq. 10) using the matrix \hat{U} , and apply argsort to obtain permutations.

Combining all of the above, Algorithm 4 describes the process of generating permutation samples from a Sobol sequence. Figure 7 shows 200 Sobol points distributed on the surface of the sphere. As our Sobol sequence and inverse CDF sampling generate points uniformly distributed on the n-sphere, Theorem 1 applies, and Algorithm 4 samples permutations from a uniform distribution in an unbiased way. Figure 8 shows the distribution of 1000 permutations sampled with d = 4, which is clearly uniform.

5. Evaluation

Sampling techniques are evaluated using two strategies: the root mean squared error (RMSE) of Shapley value estimates as compared to exact Shapley values, and the discrepancies of point sets with respect to the Mallows kernel. Machine learning models are generated for the datasets listed in Table 1. The XGBoost algorithm by Chen and Guestrin (2016) is used to generate gradient boosted decision tree models (GBDT). GBDT models are used as they are nonlinear, yielding nontrivial Shapley values, but we also have access to exact Shapley values for comparison. Exact values for GBDT models are computed via



Figure 7: Sobol sphere

Figure 8: Sobol permutations

Table 1: Datasets

NAME	ROWS	COLS	TASK	REF
ADULT	48842	14	CLASS	Конауі (1996)
BREAST_CANCER	699	30	CLASS	Mangasarian and Wolberg (1990)
CAL_HOUSING	20640	8	REGR	Pace and Barry (1997)
MAKE_REGRESSION	1000	10	REGR	Pedregosa et al. (2011)

the domain-specific, polynomial-time, TreeShap algorithm of Lundberg et al. (2020), which otherwise can be intractable to obtain. For each dataset/algorithm combination, models are trained on the entire dataset, then Shapley values are evaluated for all features of 10 randomly chosen instances, using a background dataset of 100 instances to marginalise out features. The root mean squared error is obtained by comparison to exact Shapley values. This process is repeated 25 times to obtain error bars. In addition to the permutation-based algorithms proposed in this paper, we evaluate and compare to algorithms sampling from binary sets. For a fair comparison, all algorithms are evaluated according to number of evaluations of $v(S \cup i) - v(S)$, written as 'marginal_evals' on the x-axis of figures. If the algorithm samples permutations, the number of marginal evaluations is proportional to nd, where n is the number of permutations sampled.

5.1 Sampling Methods

Given three possible kernels for kernel herding, Bayesian quadrature, and sequential Bayesian quadrature, as well as a free parameter for the (normalised) Mallows kernel, we first make a comparison among kernels. Figure 9 shows the RMSE for Shapley value estimation using various kernels with the kernel herding algorithm. The Mallows kernel with $\lambda = 5$ either has the lowest error, or contains the lowest error inside its 95% confidence interval, on all 4 datasets. For this reason, as well as its superior theoretical properties (it is a universal kernel), we use this configuration for subsequent experiments. Figure 10 shows the behaviour of the Mallows kernel with different values for the λ parameter—it controls the rate of change in the similarity measure as the number of discordant pairs between two permutations is increased.



Figure 9: Kernel herding with various kernels



Figure 10: Mallows kernel falloff

Figure 11: Argmax trials

The kernel herding algorithm relies on an argmax procedure to select each new sample. Exact solutions are not easily available, so we rely on random sampling, where the best sample of k trials is selected in each iteration. As the number of argmax samples k is increased, the execution time increases linearly and the accuracy of the optimisation procedure increases. Figure 11 shows the improvement in RMSE for kernel herding on the cal_housing dataset, with the Mallows kernel and $\lambda = 5$, by varying the number of argmax trials from 5 to 50. We settle on 25 trials as a compromise between accuracy and runtime.

Due to the large number of algorithms under evaluation, we first compare the existing algorithms from the literature, described in Section 2: plain Monte Carlo (MC), antithetic sampling (MC-antithetic), multilinear extension (Owen), multilinear extension with anti-thetic sampling (Owen-Halved), and stratified sampling (Stratified). Results are shown in Figure 12.

As MC-antithetic has the lowest error in all cases, we compare it to the algorithms introduced in this work: kernel herding (Herding), sequential Bayesian quadrature (SBQ), orthogonal spherical codes (Orthogonal) and Sobol permutations (Sobol). As discussed above, kernel methods use the Mallows kernel with $\lambda = 5$ and 25 argmax trials. Results are shown in Figure 13. Sequential Bayesian quadrature is highly effective on the *adult*, *cal_housing* and *make_regression* datasets, but it is also the most expensive method and impractical to evaluate for n > 100. Herding is effective on the same datasets as SBQ, but can take larger numbers of samples and so achieves the highest accuracies of any algorithm for *adult*, *cal_housing* and *make_regression*. The *breast_cancer* dataset exhibits different properties, with Herding and SBQ both being outperformed by simple antithetic sampling. Sampling with orthogonal spherical codes shows similar empirical performance to antithetic sampling, but with better convergence on all datasets, and the lowest error of any algorithm for *breast_cancer*. The Sobol method is less effective for small numbers of samples, but achieves better results than orthogonal sampling given sufficiently many samples, and is cheaper to evaluate than herding or SBQ.



Figure 12: Existing algorithms



Figure 13: New algorithms

	Argopress	Uppppyg		Operation	CDO	Capar
	ALGORITHM	HERDING	MC-ANTI	ORTHO	SBQ	SOBOL
D	Ν					
	10	0.259652	0.280529	0.262992	0.258927	0.274002
10	100	0.066303	0.090138	0.076668	0.064028	0.073511
	1000	0.015502	0.028335	0.024012	-	0.01909
	10	0.285612	0.289526	0.28641	0.285547	0.286821
50	100	0.08563	0.091634	0.07941	0.08531	0.084681
	1000	0.024938	0.028939	0.025099	-	0.024717
	10	0.29287	0.290015	0.289343	0.292625	0.288284
200	100	0.089084	0.091779	0.088869	0.089019	0.089375
	1000	0.027349	0.029016	0.025562	-	0.025572

Table 2: Discrepancy (lower is better)

Table 3: Complexity in n

Algorithm	Complexity		
MC-Antithetic Herding SBQ Orthogonal Sobol	$ \begin{array}{c} O(n) \\ O(n^2) \\ O(n^3) \\ O(n) \\ O(n) \\ O(n) \end{array} $		

5.2 Discrepancy

Table 2 shows mean discrepancies over 25 trials for the various permutation sampling algorithms, calculated as per Equation 8 using the Mallows kernel with $\lambda = 5$. We omit results for SBQ at n = 1000 due to runtime. Computational complexity with respect to nfor each algorithm is listed in Table 3. At low dimension the methods directly optimising discrepancy, herding and SBQ, achieve significantly lower discrepancies than other methods. For d = 10, n = 1000, herding achieves almost a twofold reduction in discrepancy over antithetic sampling, directly corresponding to an almost twofold lower error bound under the Koksma-Hlawka inequality. Antithetic sampling has a higher discrepancy than all other methods here, except in one case, d = 200, n = 10, where it achieves lower discrepancy than herding and SBQ. In general we see the orthogonal and Sobol methods are the most effective at higher dimensions, collectively accounting for the lowest discrepancies at d = 200. These results show that no single approach is best for all problems but all of the newly introduced methods provide significant improvements over the status quo.

The discrepancies computed above are applicable beyond the applications discussed in this paper. Tables 2 and 3 provide a reference for how to select samples of permutations at a given computational budget and dimension, not just for Shapley value estimation, but for any bounded function $f : \mathfrak{S}_d \to \mathbb{R}$.

6. Conclusion

We show new techniques for the approximation of Shapley values in machine learning applications based on careful selection of samples from the symmetric group \mathfrak{S}_d . One set of techniques draws on theory of reproducing kernel Hilbert spaces and the optimisation of discrepancies for functions of permutations, and another exploits connections between permutations and the hypersphere \mathbb{S}^{d-2} . We perform empirical analysis of approximation error against exact Shapley values, and also calculate discrepancies for generated sets of permutations. The introduced sampling methods show improved convergence over existing state-of-the-art methods in almost all cases. Our results show that kernel-based methods are more effective for lower dimensional problems, and methods sampling from \mathbb{S}^{d-2} are more effective for higher dimensional problems.

References

- LE Blumenson. A derivation of n-dimensional spherical coordinates. *The American Mathematical Monthly*, 67(1):63–66, 1960.
- Johann S Brauchart and Josef Dick. Quasi-monte carlo rules for numerical integration over the unit sphere S². Numerische Mathematik, 121(3):473-502, 2012.
- Javier Castro, Daniel Gómez, and Juan Tejada. Polynomial calculation of the shapley value based on sampling. *Computers & Operations Research*, 36(5):1726-1730, 2009. ISSN 0305-0548. doi: https://doi.org/10.1016/j.cor.2008.04.004. URL https: //www.sciencedirect.com/science/article/pii/S0305054808000804. Selected papers presented at the Tenth International Symposium on Locational Decisions (ISOLDE X).
- Javier Castro, Daniel Gómez, Elisenda Molina, and Juan Tejada. Improving polynomial estimation of the shapley value by stratified random sampling with optimum allocation. Computers & Operations Research, 82:180–188, 2017. ISSN 0305-0548. doi: https: //doi.org/10.1016/j.cor.2017.01.019. URL https://www.sciencedirect.com/science/ article/pii/S030505481730028X.
- Tianqi Chen and Carlos Guestrin. XGBoost: A scalable tree boosting system. In KDD, pages 785–794. ACM, 2016.
- Yutian Chen, Max Welling, and Alex Smola. Super-samples from kernel herding. In Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence, UAI'10, page 109–116, Arlington, Virginia, USA, 2010. AUAI Press. ISBN 9780974903965.
- Krzysztof Choromanski, Mark Rowland, Wenyu Chen, and Adrian Weller. Unifying orthogonal monte carlo methods. In *International Conference on Machine Learning*, pages 1203–1212. PMLR, 2019.
- Shay Cohen, Gideon Dror, and Eytan Ruppin. Feature selection via coalitional game theory. Neural Computation, 19(7):1939–1961, 2007.

- J. H. Conway, N. J. A. Sloane, and E. Bannai. *Sphere-Packings, Lattices, and Groups.* Springer-Verlag, Berlin, Heidelberg, 1987. ISBN 038796617X.
- Ian Covert, Scott Lundberg, and Su-In Lee. Explaining by removing: A unified framework for model explanation. arXiv preprint arXiv:2011.14878, 2020.
- Xiaotie Deng and Christos H Papadimitriou. On the complexity of cooperative solution concepts. *Mathematics of operations research*, 19(2):257–266, 1994.
- Persi Diaconis. Group representations in probability and statistics. Institute of Mathematical Statistics Lecture Notes—Monograph Series, 11. Institute of Mathematical Statistics, Hayward, CA, 1988. ISBN 0-940600-14-5. URL http://projecteuclid.org/euclid. lnms/1215467407.
- Josef Dick and Friedrich Pillichshammer. Digital nets and sequences: discrepancy theory and quasi-Monte Carlo integration. Cambridge University Press, 2010.
- G. Th. Guilbaud and P. Rosenstiehl. Analyse algébrique d'un scrutin. Mathématiques et Sciences humaines, 4:9-33, 1963. URL www.numdam.org/item/MSH_1963_4_9_0/.
- Doug P Hardin, TJ Michaels, and Edward B Saff. A comparison of popular point configurations on \mathbb{S}^2 . Dolomites Research Notes on Approximation, 9:16–49, 2016.
- Fred J. Hickernell. What affects the accuracy of quasi-monte carlo quadrature? In Harald Niederreiter and Jerome Spanier, editors, *Monte-Carlo and Quasi-Monte Carlo Methods* 1998, pages 16–55, Berlin, Heidelberg, 2000. Springer Berlin Heidelberg. ISBN 978-3-642-59657-5.
- Edmund Hlawka. Funktionen von beschränkter variatiou in der theorie der gleichverteilung. Annali di Matematica Pura ed Applicata, 54(1):325–333, 1961.
- Ferenc Huszár and David Duvenaud. Optimally-weighted herding is bayesian quadrature. In Proceedings of the Twenty-Eighth Conference on Uncertainty in Artificial Intelligence, UAI'12, page 377–386, Arlington, Virginia, USA, 2012. AUAI Press. ISBN 9780974903989.
- Yunlong Jiao and Jean-Philippe Vert. The kendall and mallows kernels for permutations. In Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37, ICML'15, page 1935–1944. JMLR.org, 2015.
- Maurice G Kendall. A new measure of rank correlation. *Biometrika*, 30(1/2):81–93, 1938.
- William R. Knight. A computer method for calculating kendall's tau with ungrouped data. Journal of the American Statistical Association, 61(314):436-439, 1966. ISSN 01621459. URL http://www.jstor.org/stable/2282833.
- Donald E. Knuth. The Art of Computer Programming, Volume 2 (3rd Ed.): Seminumerical Algorithms. Addison-Wesley Longman Publishing Co., Inc., USA, 1997. ISBN 0201896842.

- Ron Kohavi. Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid. In KDD, pages 202–207. AAAI Press, 1996.
- Maria Lomeli, Mark Rowland, Arthur Gretton, and Zoubin Ghahramani. Antithetic and monte carlo kernel estimators for partial rankings. *Statistics and Computing*, 29(5):1127– 1147, 2019.
- Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, Advances in Neural Information Processing Systems 30, pages 4765-4774. Curran Associates, Inc., 2017. URL http://papers.nips.cc/ paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf.
- Scott M. Lundberg, Gabriel Erion, Hugh Chen, Alex DeGrave, Jordan M. Prutkin, Bala Nair, Ronit Katz, Jonathan Himmelfarb, Nisha Bansal, and Su-In Lee. From local explanations to global understanding with explainable ai for trees. *Nature Machine Intelligence*, 2(1):2522–5839, 2020.
- Sasan Maleki. Addressing the computational issues of the Shapley value with applications in the smart grid. PhD thesis, University of Southampton, 2015.
- Olvi L Mangasarian and William H Wolberg. Cancer diagnosis via linear programming. Technical report, University of Wisconsin-Madison Department of Computer Sciences, 1990.
- Horia Mania, Aaditya Ramdas, Martin J Wainwright, Michael I Jordan, and Benjamin Recht. On kernel methods for covariates that are rankings. *Electronic Journal of Statis*tics, 12:2537–2577, 2018.
- Irwin Mann and Lloyd S Shapley. Values of large games, IV: Evaluating the electoral college by Montecarlo techniques. Rand Corporation, 1960.
- Francesco Mezzadri. How to generate random matrices from the classical compact groups. arXiv preprint math-ph/0609050, 2006.
- Thomas Muir. On a simple term of a determinant. In *Proc. Royal Society Edinburg*, volume 21, pages 441–477, 1898.
- Harald Niederreiter. Random Number Generation and Quasi-Monte Carlo Methods. Society for Industrial and Applied Mathematics, USA, 1992. ISBN 0898712955.
- A. O'Hagan. Bayes-hermite quadrature. Journal of Statistical Planning and Inference, 29 (3):245-260, 1991. ISSN 0378-3758. doi: https://doi.org/10.1016/0378-3758(91)90002-V. URL https://www.sciencedirect.com/science/article/pii/037837589190002V.
- Ramin Okhrati and Aldo Lipani. A multilinear sampling algorithm to estimate shapley values. In *Proc. of ICPR*, ICPR, 2020.
- Art B Owen. Quasi-monte carlo sampling. Monte Carlo Ray Tracing: Siggraph, 1:69–88, 2003.

- Guillermo Owen. Multilinear extensions of games. Management Science, 18(5):P64-P79, 1972. ISSN 00251909, 15265501. URL http://www.jstor.org/stable/2661445.
- R Kelley Pace and Ronald Barry. Sparse spatial autoregressions. *Statistics & Probability Letters*, 33(3):291–297, 1997.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- S. M. Plis, T. Lane, and V. D. Calhoun. Permutations as angular data: Efficient inference in factorial spaces. In 2010 IEEE International Conference on Data Mining, pages 403–410, 2010. doi: 10.1109/ICDM.2010.122.
- William H. Press, Saul A. Teukolsky, William T. Vetterling, and Brian P. Flannery. Numerical Recipes 3rd Edition: The Art of Scientific Computing. Cambridge University Press, USA, 3 edition, 2007. ISBN 0521880688.
- Carl Edward Rasmussen and Zoubin Ghahramani. Bayesian monte carlo. Advances in neural information processing systems, pages 505–512, 2003.
- Reuven Y Rubinstein and Dirk P Kroese. Simulation and the Monte Carlo method, volume 10. John Wiley & Sons, 2016.
- Dino Sejdinovic, Bharath Sriperumbudur, Arthur Gretton, and Kenji Fukumizu. Equivalence of distance-based and rkhs-based statistics in hypothesis testing. *The Annals of Statistics*, pages 2263–2291, 2013.
- Lloyd S Shapley. A value for n-person games. Contributions to the Theory of Games, 2(28): 307–317, 1953.
- Il'ya Meerovich Sobol'. On the distribution of points in a cube and the approximate evaluation of integrals. *Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki*, 7(4): 784–802, 1967.
- Erik Strumbelj and Igor Kononenko. An efficient explanation of individual classifications using game theory. J. Mach. Learn. Res., 11:1–18, March 2010. ISSN 1532-4435.
- G. L. Thompson. Generalized permutation polytopes and exploratory graphical methods for ranked data. *The Annals of Statistics*, 21(3):1401–1430, 1993. ISSN 00905364. URL http://www.jstor.org/stable/2242202.
- Erik Štrumbelj and Igor Kononenko. Explaining prediction models and individual predictions with feature contributions. *Knowl. Inf. Syst.*, 41(3):647–665, December 2014. ISSN 0219-1377. doi: 10.1007/s10115-013-0679-x. URL https://doi.org/10.1007/ s10115-013-0679-x.

Appendix A. Proof of Theorem 2 (See page 17)

Theorem 2 Suppose $\sigma, \sigma' \in \mathfrak{S}_d$. Then

$$-2 + 4\left(\frac{1 - K_{\tau}(\sigma, \sigma')}{2}\right)^{3/2} \le A(\sigma)^T A(\sigma') - 3K_{\tau}(\sigma, \sigma') + O(d^{-1}) \le 2 - 4\left(\frac{1 + K_{\tau}(\sigma, \sigma')}{2}\right)^{3/2}$$

and, if $A(\sigma)^T A(\sigma') = o(1)$, then

$$|K_{\tau}(\sigma, \sigma')| \le 1/2 + o(1).$$

Proof For $1 \leq a \leq d-1$, write $t_a \in \mathfrak{S}_d$ for the adjacent transposition of a and a+1, i.e., the permutation so that $t_a(j) = j$ for $j \neq a, a+1, t_a(a) = a+1$ and $t_a(a+1) = a$. We interpret a product of permutations to be their composition as functions. For a permutation $\pi \in \mathfrak{S}_d$, write $\nu(\pi)$ for the quantity $\sum_{j=1}^d j\pi(j)$, and note that $\nu(I) = \sum_{j=1}^d j^2 = d(d+1)(2d+1)/6$. It is well-known that the number of inversions $n_{\text{dis}}(I,\pi) = |\{(i,j) : i < j \text{ and } \pi(i) > i\}$

It is well-known that the number of inversions $n_{\text{dis}}(I,\pi) = |\{(i,j) : i < j \text{ and } \pi(i) : \pi(j)\}|$ in a permutation π equals the least k so that there exist a_1, \ldots, a_k with

$$\pi = \prod_{i=1}^{k} t_{a_i}.$$
(15)

This quantity k is known as the "length" of π and is exactly the distance in the 1-skeleton of the permutohedron representation of \mathfrak{S}_d . Furthermore, the a_i can be obtained via bubble sort, i.e., the product (15) begins with

$$t_{\pi(1)-1}t_{\pi(1)-2}\cdots t_1$$

and proceeds recursively on $\pi|_{\{2,\dots,d\}}$. Write π_j for the product of the first j terms in (15) for $1 \leq j \leq k$, i.e., $\pi_j = \prod_{i=1}^j t_{a_i}$, with $\pi_0 = I$. Then the pairs $e_j = \{\pi_j(a_j), \pi_j(a_j+1)\}$ are all distinct, because entries of π in one-line notation switch places at most once when applying the adjacent transpositions, i.e., a larger value a, once it switches places with a smaller value b immediately to its left, never switches place with b again. Furthermore, note that

$$\nu(\pi_{j+1}) - \nu(\pi_j) = (j\pi_{j+1}(a_j) + (j+1)\pi_{j+1}(a_j+1)) - (j\pi_j(a_j) + (j+1)\pi_j(a_j+1)) = (j\pi_j(a_j+1) + (j+1)\pi_j(a_j)) - (j\pi_j(a_j) + (j+1)\pi_j(a_j+1)) = \pi_j(a_j+1) - \pi_j(a_j),$$

a quantity which is always negative because the sequence of transpositions obtained above only ever increases the number of inversions. Therefore, the collection $\{e_j\}_{j=1}^k$ consists of k distinct edges of a complete graph on $\{1, \ldots, d\}$ and

$$\nu(\pi) = \nu(\pi_k) = \nu(\pi_k) - \nu(I) + \frac{d(d+1)(2d+1)}{6}$$
$$= \frac{d(d+1)(2d+1)}{6} + \sum_{j=1}^k \pi_j(a_j+1) - \pi_j(a_j)$$
$$= \frac{d(d+1)(2d+1)}{6} - \sum_{j=1}^k \operatorname{wt}(e_j)$$

where wt($\{a, b\}$) = |b - a|. By greedily selecting the highest-weight or lowest-weight edges of the complete graph K_d weighted by wt(·), the quantity $\sum_{j=1}^k \text{wt}(e_j)$ is always at least

$$1 \cdot (d-1) + 2 \cdot (d-2) + \dots + (d-m) \cdot m = \frac{(d+2m-1)(d-m+1)(d-m)}{6}$$

where m is the smallest integer so that $\sum_{j=1}^{d-m} (d-j) = (d+m-1)(d-m)/2 \le k$, because the summands correspond to d-1 edges of weight 1, d-2 edges of weight 2, and so on up to m edges of weight d-m. Similarly, $\sum_{j=1}^{k} \operatorname{wt}(e_j)$ is at most

$$(d-1)\cdot 1 + (d-2)\cdot 2 + \dots + M\cdot (d-M) = \frac{(d+2M-1)(d-M+1)(d-M)}{6}$$

where M is the largest integer so that $\sum_{j=1}^{d-M} j = (d-M)(d-M+1)/2 \ge k$, since in this case we bound the total edge weight via 1 edge of weight d-1, 2 edges of weight d-2, and so on up to d-M edges of weight M. Then, letting $\alpha = k/\binom{d}{2}$ (so that $\alpha \in [0,1]$),

$$m = \left\lfloor \frac{\sqrt{4d^2 - 4d - 8k + 1} + 1}{2} \right\rfloor = d\sqrt{1 - \alpha} \pm 1$$
$$M = \left\lceil \frac{2d - \sqrt{8k + 1} + 1}{2} \right\rceil = d(1 - \sqrt{\alpha}) \pm 1$$

It is straightforward to verify that, if f(s) = (d+2s-1)(d-s+1)(d-s)/6, then s = O(d) implies $f(s \pm 1) = f(s) + O(d^2)$. So, letting $\alpha = k/\binom{d}{2}$ (so that $\alpha \in [0, 1]$)

$$\begin{split} \nu(\pi) &\leq \frac{d(d+1)(2d+1)}{6} - f(M) \\ &= \frac{d(d+1)(2d+1)}{6} - f(d\sqrt{1-\alpha}) + O(d^2) \\ &= \frac{d^3}{3} - \frac{d^3(1+2\sqrt{1-\alpha})(1-\sqrt{1-\alpha})^2}{6} + O(d^2) \\ &= d^3 \left(\frac{2}{3} - \frac{\alpha}{2} - \frac{(1-\alpha)^{3/2}}{3}\right) + O(d^2) \end{split}$$

and

$$\begin{split} \nu(\pi) &\geq \frac{d(d+1)(2d+1)}{6} - f(m) \\ &= \frac{d^3}{3} - f(d(1-\sqrt{\alpha})) + O(d^2) \\ &= \frac{d^3}{3} - \frac{d^3(1+2(1-\sqrt{\alpha}))(1-(1-\sqrt{\alpha}))^2}{6} + O(d^2) \\ &= d^3 \left(\frac{1}{3} - \frac{\alpha}{2} + \frac{\alpha^{3/2}}{3}\right) + O(d^2). \end{split}$$

(Note that the functions in parentheses meet for $\alpha = 0, 1$.) Thus, applying the fact that $\nu(\sigma' \circ \sigma^{-1}) = I^T(\sigma' \circ \sigma^{-1}) = \sigma^T \sigma'$, where we regard permutations both as functions π of $\{1, \ldots, d\}$ and as vectors $(\pi(1), \ldots, \pi(d))$,

$$2 + 2\alpha^{3/2} \le \frac{6\sigma^T \sigma'}{d^3} + O(d^{-1}) + 3\alpha \le 4 - 2(1-\alpha)^{3/2}$$

Then, since

$$K_{\tau}(\sigma, \sigma') = 1 - \frac{2n_{\text{dis}}(I, \sigma'\sigma^{-1})}{\binom{d}{2}} = 1 - 2\alpha$$

we have

$$\frac{1}{4} + \left(\frac{1 - K_{\tau}(\sigma, \sigma')}{2}\right)^{3/2} \le \frac{3\sigma^T \sigma'}{d^3} + O(d^{-1}) - \frac{3K_{\tau}(\sigma, \sigma')}{4} \le \frac{5}{4} - \left(\frac{1 + K_{\tau}(\sigma, \sigma')}{2}\right)^{3/2}.$$

Writing $\sigma = \rho x + \mu$ and $\sigma' = \rho x' + \mu$ yields the first claim of the result, since then

$$\sigma^T \sigma' = \frac{d(d^2 - 1)}{12} A(\sigma)^T A(\sigma') + \frac{d(d + 1)^2}{4}.$$

For the second claim, note that, if $\sigma^T \sigma' = d^3(1/4 + o(1))$ (the expected value for random permutations, corresponding to $A(\sigma)^T A(\sigma') \approx 0$),

$$-2 + 4\left(\frac{1 - K_{\tau}(\sigma, \sigma')}{2}\right)^{3/2} \le -3K_{\tau}(\sigma, \sigma') + O(d^{-1}) \le 2 - 4\left(\frac{1 + K_{\tau}(\sigma, \sigma')}{2}\right)^{3/2},$$

i.e.,

$$|K_{\tau}(\sigma, \sigma')| \le 1/2 + o(1).$$