
COMPUTING EXACT SHAPLEY VALUES IN POLYNOMIAL TIME FOR PRODUCT-KERNEL METHODS

Anonymous authors

Paper under double-blind review

ABSTRACT

Kernel methods are widely used in machine learning due to their flexibility and expressiveness. However, their black-box nature poses significant challenges to interpretability, limiting their adoption in high-stakes applications. Shapley value-based feature attribution techniques, such as SHAP and kernel method-specific adaptation like RKHS-SHAP, offer a promising path toward explainability. Yet, computing exact Shapley values is generally intractable, leading existing methods to rely on approximations and thereby incur unavoidable error. In this work, we introduce PKeX-Shapley, a novel algorithm that utilizes the multiplicative structure of product kernels to enable the exact computation of Shapley values in polynomial time. The core of our approach is a new value function, the *functional baseline value function*, specifically designed for product-kernel models. This value function removes the influence of a feature subset by setting its functional component to the least informative state. Crucially, it allows a recursive thus efficient computation of Shapley values in polynomial time. As an important additional contribution, we show that our framework extends beyond predictive modeling to statistical inference. In particular, it generalizes to popular kernel-based discrepancy measures such as the Maximum Mean Discrepancy (MMD) and the Hilbert–Schmidt Independence Criterion (HSIC), thereby providing new tools for interpretable statistical inference.

1 INTRODUCTION

Shapley values (Shapley, 1953), a solution concept originating from cooperative game theory, offer a principled, axiomatic framework for feature attribution in machine learning (ML) (Lundberg and Lee, 2017; Covert et al., 2020). Thanks to their rigorous axiomatic foundation, there has been a widespread adoption within the explainable AI community (Sundararajan and Najmi, 2020). They allow a model’s output—such as predictions or losses—to be *fairly* distributed across input features based on their individual and joint contributions. Consequently, several algorithms have been proposed to estimate Shapley values under different modelling scenarios, ranging from model-agnostic methods like Kernel SHAP (Lundberg and Lee, 2017) to model-specific methods that leverage structural properties to improve statistical or computational efficiency. Well-known examples of the latter include Tree SHAP for tree-based models (Lundberg et al., 2020), GPSHAP for Gaussian processes (Chau et al., 2024), Deep SHAP for deep networks (Lundberg and Lee, 2017), and, most relevant to our work, RKHS-SHAP for kernel methods (Chau et al., 2022). Kernel methods are particularly notable—not only for their use in prediction tasks but also in a broad range of statistical inference problems, including measuring distributional closeness (Gretton et al., 2006; Naslidnyk et al., 2025), two-sample testing (Schrab, 2025; Chau et al., 2025), goodness-of-fit testing (Chwialkowski et al., 2016; Liu et al., 2016), independence testing (Albert et al., 2022), causal inference and discovery (Mitrovic et al., 2018; Chau et al., 2021; Sejdinovic, 2024), among others. Hence, as kernel methods gain widespread adoption in high-stakes applications for their flexibility and expressive power, the need for interpretability has become increasingly vital.

Despite their game-theoretic foundation, Shapley values face two major challenges when applied to machine learning. The first is defining and estimating a suitable value function that quantifies the contribution of a coalition of features. Ideally, this function should capture the model’s behavior when the complementary features are absent. A natural approach is to retrain the model on each subset and use the resulting prediction as the value (Lipovetsky and Conklin, 2001), but this is

054 computationally infeasible due to the exponential number of retrainings required. A more common
 055 alternative simulates feature absence through marginal or conditional expectations of the model’s
 056 output, given the fixed subset (see Sundararajan and Najmi (2020) for other formulations). However,
 057 estimating these expectations reliably is difficult, as it often involves density estimation or simplifying
 058 assumptions such as feature independence—assumptions that can lead to misleading and unfaithful
 059 explanations (Kumar et al., 2020). Chau et al. (2022) tackle this problem in the context of kernel
 060 methods by exploiting the structure of reproducing kernel Hilbert spaces (RKHS) and using kernel
 061 distributional embeddings (Muandet et al., 2017) to estimate value functions nonparametrically,
 062 thereby avoiding density estimation and independence assumptions.

063 Once a value function is fixed, the second challenge lies in efficiently estimating the Shapley
 064 values themselves. Exact computation requires evaluating the value function over all 2^d possible
 065 feature subsets for d features, which is computationally demanding. To alleviate this, approximation
 066 techniques such as Monte Carlo sampling and regression-based methods are widely used, relying
 067 on a smaller set of evaluations. While these approaches reduce computational costs, they inevitably
 068 introduce estimation errors that grow with both sample size and feature dimensionality (Kumar
 069 et al., 2020). In some cases, model-specific structures can be exploited for efficiency—for example,
 070 Tree SHAP leverages tree decompositions for exact polynomial-time computation. By contrast,
 071 RKHS-SHAP improves the statistical estimation of the value function but still relies on regression-
 072 based approximations for the Shapley values themselves, and thus inherits the same computational
 073 limitations.

074 To design an efficient Shapley value-based attribution algorithm for kernel methods, we focus on a subclass that
 075 employs product kernels—referred to as *product-kernel methods*—and introduce *PKeX-Shapley* (Product-**K**ernel-
 076 based **eX**act **S**hapley attribution). Product kernels are particularly attractive as they are simple to implement
 077 while retaining strong theoretical guarantees. For instance, the product of universal kernels remains universal, so the
 078 resulting RKHS can approximate any continuous function provided that each component kernel is sufficiently expressive
 079 (Szabó and Sriperumbudur, 2018). Although this focus narrows the scope compared with Chau et al. (2022),
 080 which considers generic kernels, the additional structure enables us to directly address the two challenges outlined
 081 above. Specifically, we introduce a new value function, called the *functional baseline value function*,
 082 which leverages the multiplicative structure of product kernels. Conceptually, it parallels baseline
 083 value functions: rather than setting features to fixed baseline inputs, we set the corresponding kernel
 084 components to their least informative state, thereby eliminating their influence. This design admits a
 085 natural interpretation as the optimal orthogonal projection of the predictive function onto a function
 086 class restricted to a given feature subset. Crucially, this formulation enables a recursive algorithm
 087 for computing exact Shapley values in polynomial time (see Figure 1 for a runtime comparison
 088 with naïve computation). This yields substantial efficiency gains without relying on approximation
 089 techniques such as sampling or regression-based estimation.

095 As an additional, independent contribution, we demonstrate that our framework extends beyond
 096 kernel-based predictive models. In particular, we show how our method applies to kernel-based
 097 statistical discrepancies, including the Maximum Mean Discrepancy (MMD) and Hilbert–Schmidt
 098 Independence Criterion (HSIC), thereby enabling interpretable statistical inference. A Python
 099 implementation of our method is publicly available (pke, 2025).

101 2 BACKGROUND ON SHAPLEY VALUES

102 **Notation.** Let \mathcal{D} denote the set of d features, and $2^{\mathcal{D}}$ its power set. The training set $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
 103 consists of n samples, where $\mathbf{x}^{(i)} \in \mathbb{R}^d$ and $y^{(i)} \in \mathbb{R}$ (or a discrete label set for classification tasks).
 104 Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the matrix of features, and $\mathbf{X}_{\mathcal{S}} \in \mathbb{R}^{n \times s}$ the submatrix restricted to features
 105 in subset $\mathcal{S} \subseteq \mathcal{D}$, and we write $\mathbf{X}_j := \mathbf{X}_{\{j\}}$. We use capital letters for random variables, bold capital
 106 letters for matrices, calligraphic letters for sets, and bold lowercase letters for vectors. The restriction
 107

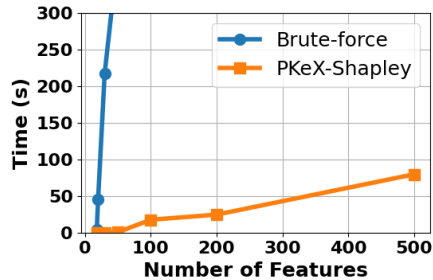


Figure 1: Execution time of brute-force and our PKeX-Shapley with a 300s budget (see details in Appendix H.2).

of a vector \mathbf{x} to features in \mathcal{S} is denoted by $\mathbf{x}_{\mathcal{S}}$. The elementwise product is denoted by \odot , and expectation by \mathbb{E} . A symmetric positive (semi-)definite kernel function over \mathcal{D} is denoted by k , and its restriction to subset \mathcal{S} by $k_{\mathcal{S}}$. The corresponding kernel matrices are denoted by \mathbf{K} and $\mathbf{K}_{\mathcal{S}}$, respectively. All proofs are provided in Appendix D.

Shapley value. The Shapley value (Shapley, 1953) is a widely used game-theoretic method for feature attribution in predictive models. It assigns importance to each input feature by averaging its marginal contribution across all possible feature subsets, with weights determined by coalition size. This construction uniquely satisfies several fairness axioms, making it a principled approach to distributing a model’s output among its features. Specifically, given a value function $v : 2^{\mathcal{D}} \rightarrow \mathbb{R}$, which quantifies the contribution of a feature subset, the Shapley value for a feature j is defined as

$$\phi_j = \sum_{\mathcal{S} \subseteq \mathcal{D} \setminus \{j\}} \mu(|\mathcal{S}|) (v(\mathcal{S} \cup \{j\}) - v(\mathcal{S})), \quad (1)$$

where $\mu(s) = s!(d-s-1)!/d!$ and $v(\mathcal{S} \cup \{j\}) - v(\mathcal{S})$ measures the marginal contribution of feature j when added to subset \mathcal{S} . The Shapley value is the unique solution satisfying four axioms of efficiency, null player, symmetry, and linearity for any given value function (Shapley, 1953). However, in the context of explainability, defining this value function v —which determines how subsets of features contribute to the model’s output—is a non-trivial design choice. Several formulations have been proposed (Sundararajan and Najmi, 2020), but two are widely adopted: the interventional value function $v_{\mathbf{x}}(\mathcal{S}) = \mathbb{E}_{X_{\mathcal{D} \setminus \mathcal{S}} | X_{\mathcal{S}} = \mathbf{x}_{\mathcal{S}}} [f(\mathbf{x}_{\mathcal{S}}, X_{\mathcal{D} \setminus \mathcal{S}})]$, which replaces missing features with samples from their marginal distributions (Janzing et al., 2020), and the observational value function $v_{\mathbf{x}}(\mathcal{S}) = \mathbb{E}_{X_{\mathcal{D} \setminus \mathcal{S}}} [f(X) | X_{\mathcal{S}} = \mathbf{x}_{\mathcal{S}}]$, which imputes them using the conditional distribution given observed features (Lundberg and Lee, 2017). While the observational approach preserves feature dependencies, it requires estimating complex conditional expectations, making it computationally demanding. In practice, both approaches are often implemented by sampling from marginal or conditional distributions, requiring repeated model evaluations that become costly for large models or high-dimensional inputs (Sundararajan and Najmi, 2020). Moreover, these estimates typically rely on a finite set of background samples (commonly drawn from the training set), which can substantially influence the resulting explanations (Molnar, 2023, Chapter 21).

3 EXACT SHAPLEY VALUE FOR PRODUCT-KERNEL LEARNING METHODS

In this section, we demonstrate how to compute *exact* Shapley values for local explanation in polynomial time when the predictive model is constructed through product-kernel methods, e.g., an SVM or kernel ridge regressor built using product kernels.

3.1 A NEW FUNCTIONAL BASELINE VALUE FUNCTION FOR PRODUCT-KERNEL METHODS

Product-kernel methods rely on kernel functions to capture complex relationships between input features and output. A kernel-based decision function is generally expressed as

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}^{(i)}) = \boldsymbol{\alpha}^{\top} k(\mathbf{X}, \mathbf{x})$$

where k is a product kernel function, and α_i are the learned coefficients associated with the model. Recall that the product kernel k can be expressed as products of base kernels k_j , $j \in \mathcal{D}$:

$$k(\mathbf{x}, \mathbf{x}^{(i)}) = \prod_{j \in \mathcal{D}} k_j(x_j, x_j^{(i)}). \quad (2)$$

Product-kernel methods are widely used in machine learning for their simplicity and effectiveness in modeling similarities in high-dimensional data—by designing a kernel for each feature and then multiplying them together (Gardner et al., 2018). They also come with strong theoretical guarantees: if the base kernels are universal—i.e., capable of approximating any continuous function defined on the marginal input—then the product kernel inherits this property (Szabó and Sriperumbudur, 2018). Some well-known kernels, such as the radial basis function (RBF) with an isotropic or anisotropic bandwidth, belong to the family of product kernels, i.e., $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2/2\sigma^2) = \prod_{j=1}^d \exp(-(x_j - x_j')^2/2\sigma^2)$ (see Appendix A for more details on product kernels).

We now introduce our *functional baseline value function*, motivated by the multiplicative structure of product kernels. In standard settings, interventional value functions remove the effect of a feature by replacing it with its average contribution, thereby isolating its marginal impact. More generally, this can be viewed as a baseline value function (Sundararajan and Najmi, 2020, Sec. 2.3), where features are set to a non-informative state, such as zero vectors in image data. In our setting, however, the functional form of the predictive model is directly accessible. We therefore remove a feature’s influence by factoring out its corresponding functional component and setting it to one—the least informative contribution in the multiplicative structure. This design principle yields the following definition for product-kernel methods.

Definition 1. Let α be the coefficients of the function f constructed using the product-kernel method. The functional baseline value function for an input \mathbf{x} , a feature coalition S , and coefficients α is

$$\nu_{\mathbf{x}}(S) = \alpha^\top k_S(\mathbf{X}_S, \mathbf{x}_S). \quad (3)$$

Notably, computing this value function requires only linear time, making it highly efficient. A detailed comparison with baseline and interventional value functions is provided in Appendix B.

We emphasize that the choice of a value function determines how feature subset importance is measured, and is therefore inherently subjective. Our introduction of the functional baseline value function is motivated not only by its conceptual alignment with existing formulations, but also by the substantial computational benefits it enables, as we will demonstrate in the following subsection. We do not claim that this value function is intrinsically superior in capturing feature importance; however, we demonstrate in Experiment #2 (Section 5.1) that its empirical performance in recovering influential features is promising and on par with other established baselines.

3.2 A RECURSIVE FORMULATION FOR COMPUTING SHAPLEY VALUES

Having motivated our functional baseline value function, we now demonstrate the computational gain we obtain from this formulation.

Theorem 2. Let $\mathcal{Z} := \{\mathbf{z}_1, \dots, \mathbf{z}_d\}$ with $\mathbf{z}_i = k_i(\mathbf{X}_i, x_i)$ and $e_q(\mathcal{Z})$ the elementary symmetric polynomials (ESPs) of order q over \mathcal{Z} , defined recursively as

$$e_q(\mathcal{Z}_{-j}) = \frac{1}{q} \sum_{r=1}^q (-1)^{r-1} e_{q-r}(\mathcal{Z}_{-j}) \odot p_r(\mathcal{Z}_{-j}),$$

where $p_r(\mathcal{Z}) = \sum_{\mathbf{z}_i \in \mathcal{Z}} \mathbf{z}_i^r$ is the element-wise degree- r power sum. For product-kernel learning methods with coefficients α , the Shapley value ϕ_j^α for feature j of instance \mathbf{x} can then be expressed as

$$\phi_j^\alpha := \sum_{S \subseteq \mathcal{D} \setminus \{j\}} \mu(|S|) (v_{\mathbf{x}}(S \cup \{j\}) - v_{\mathbf{x}}(S)) = \alpha^\top \left((k_j(\mathbf{X}_j, x_j) - 1) \odot \sum_{q=0}^{d-1} \mu(q) e_q(\mathcal{Z}_{-j}) \right). \quad (4)$$

The recursion follows from Newton’s identities for ESPs (Egge, 2019), which due to space constraints, is explained further in Appendix C. Here we provide the key intuition behind equation 4. The multiplicative structure of the product kernel allows us to factorize the marginal contribution $v_{\mathbf{x}}(S \cup \{j\}) - v_{\mathbf{x}}(S)$ as $\alpha^\top ((k_j(\mathbf{X}_j, x_j) - 1) \odot k_S(\mathbf{X}_S, \mathbf{x}_S))$. This structure allows us then to push the summation inside the inner product between the coefficients α and the kernel evaluations, and express the total sum over all subsets $S \subseteq \mathcal{D} \setminus \{j\}$ in terms of weighted ESPs $e_q(\mathcal{Z}_{-j})$, which can then be computed recursively in $O(d^2)$. A similar trick has been used in the context of additive Gaussian processes (Duvenaud et al., 2011), but to our knowledge, it has not been previously leveraged for computing Shapley values in polynomial time. The full algorithm is given in Appendix E, together with a numerically stable modification.

Next, we show the additivity of the explanation for the value function in equation 3 (see Appendix F for discussion on the additivity of explanations for different values for the null game).

Lemma 3. For any instance \mathbf{x} , the sum of Shapley values satisfies: $\sum_{j=1}^d \phi_j^\alpha = f(\mathbf{x}) - f_\emptyset(\mathbf{x})$ where $f_\emptyset(\mathbf{x}) = \sum_{i=1}^n \alpha_i$ represents the baseline contribution with no features.

4 EXPLAINING KERNEL-BASED STATISTICAL DISCREPANCIES

While most explainability research has focused on predictive models¹, far less attention has been given to statistical inference tasks, where kernel methods are equally powerful. To both demonstrate the flexibility of our framework and advocate for interpretable statistical inference, we show how it applies to the Maximum Mean Discrepancy (MMD) (Gretton et al., 2006) and the Hilbert–Schmidt Independence Criterion (HSIC) (Gretton et al., 2007), two widely used measures of distributional discrepancy and independence.

4.1 DISTRIBUTING THE DISCREPANCY: EXPLAINING THE MMD

The *Maximum Mean Discrepancy* (MMD) quantifies the difference between two probability distributions in terms of their kernel mean embeddings: $\text{MMD}^2(\mathbb{P}, \mathbb{Q}) := \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2$, where \mathcal{H}_k is the RKHS associated with the kernel k and $\mu_{\mathbb{P}} := \int k(\mathbf{x}, \cdot) d\mathbb{P}(\mathbf{x}) \in \mathcal{H}_k$ the kernel mean embedding of \mathbb{P} (Muandet et al., 2017). The embedding $\mu_{\mathbb{Q}}$ is defined analogously. Based on the samples $\{\mathbf{x}^{(i)}\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathbb{P}$ and $\{\mathbf{z}^{(i)}\}_{i=1}^m \stackrel{i.i.d.}{\sim} \mathbb{Q}$, the empirical estimate of MMD, expressed entirely in terms of k , is given by $\widehat{\text{MMD}}^2(\mathbb{P}, \mathbb{Q}) = \frac{1}{n(n-1)} \sum_{i \neq j} k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) + \frac{1}{m(m-1)} \sum_{i \neq j} k(\mathbf{z}^{(i)}, \mathbf{z}^{(j)}) - \frac{2}{nm} \sum_{i,j} k(\mathbf{x}^{(i)}, \mathbf{z}^{(j)})$.

In this section, we propose an attribution method to allocate the overall discrepancy measured by MMD among the involved variables. Such an attribution is useful in many problems, including explaining MMD-based statistics for hypothesis testing (e.g., Gretton et al. (2012)) and determining the contribution of variables to covariance shift (e.g., Zhang et al. (2020)). Similar to the motivation presented in Section 3, as inspired by the multiplicative structure of product kernels, we define the value function of a coalition $\mathcal{S} \subseteq \mathcal{D}$ as the MMD computed only on the kernel components corresponding to \mathcal{S} , with the remaining components set to their least informative state (equal to 1 since we are doing multiplication). This isolates the discrepancy attributable to \mathcal{S} alone.

Definition 4. For a coalition $\mathcal{S} \subseteq \mathcal{D}$, the functional baseline value function for MMD is defined as $v_{\text{MMD}}(\mathcal{S}) = \frac{1}{n(n-1)} \sum_{i \neq j} k_{\mathcal{S}}(\mathbf{x}_S^{(i)}, \mathbf{x}_S^{(j)}) + \frac{1}{m(m-1)} \sum_{i \neq j} k_{\mathcal{S}}(\mathbf{z}_S^{(i)}, \mathbf{z}_S^{(j)}) - \frac{2}{nm} \sum_{i,j} k_{\mathcal{S}}(\mathbf{x}_S^{(i)}, \mathbf{z}_S^{(j)})$.

By applying the same trick as in Theorem 2 to the value function v_{MMD} , we can obtain a similar recursive formulation of Shapley values for MMD. That is, for the first term in v_{MMD} , we use the multiplicate structure of product kernels and write $v_{\text{MMD}}(\mathcal{S} \cup \{q\}) - v_{\text{MMD}}(\mathcal{S})$ as $k_{\mathcal{S} \cup \{q\}}(\mathbf{x}_{\mathcal{S} \cup \{q\}}^{(i)}, \mathbf{x}_{\mathcal{S} \cup \{q\}}^{(j)}) - k_{\mathcal{S}}(\mathbf{x}_S^{(i)}, \mathbf{x}_S^{(j)}) = (k_q(x_q^{(i)}, x_q^{(j)}) - 1)k_{\mathcal{S}}(\mathbf{x}_S^{(i)}, \mathbf{x}_S^{(j)})$. By pushing out $k_q(x_q^{(i)}, x_q^{(j)}) - 1$ from the summation in the Shapley value, we can express the total sum over $\mathcal{S} \subseteq \mathcal{D} \setminus \{q\}$ as the weighted ESPs. The following proposition summarizes this.

Proposition 5. Let $\mathcal{Z}(\mathbf{x}, \mathbf{x}') = \{k_1(x_1, x'_1), \dots, k_d(x_d, x'_d)\}$, and $e_r(\mathcal{Z}(\mathbf{x}, \mathbf{x}'))$ determined as

$$e_r(\mathcal{Z}_{-q}(\mathbf{x}, \mathbf{x}')) = \frac{1}{r} \sum_{s=1}^r (-1)^{s-1} e_{r-s}(\mathcal{Z}_{-q}(\mathbf{x}, \mathbf{x}')) p_s(\mathcal{Z}_{-q}(\mathbf{x}, \mathbf{x}')),$$

where $p_s(\mathcal{Z}) = \sum_{z \in \mathcal{Z}} z^s$ represents the degree- s power sum. Further, let $\gamma_q(\mathbf{x}, \mathbf{x}')$ be defined as $\gamma_q(\mathbf{x}, \mathbf{x}') := (k_q(x_q, x'_q) - 1) \sum_{r=0}^{d-1} \mu(r) e_r(\mathcal{Z}_{-q}(\mathbf{x}, \mathbf{x}'))$. Then, for product kernels, the Shapley value for the MMD can be recursively computed as

$$\phi_q^{\text{MMD}} = \frac{1}{n(n-1)} \sum_{i \neq j} \gamma_q(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) + \frac{1}{m(m-1)} \sum_{i \neq j} \gamma_q(\mathbf{z}^{(i)}, \mathbf{z}^{(j)}) - \frac{2}{nm} \sum_{i,j} \gamma_q(\mathbf{x}^{(i)}, \mathbf{z}^{(j)}).$$

The Shapley values ϕ_q^{MMD} allow us to allocate the overall distributional discrepancy between \mathbb{P} and \mathbb{Q} across the variables, identifying the most influential ones in distinguishing the two distributions.

4.2 DISTRIBUTING THE DEPENDENCE: EXPLAINING THE HSIC

The *Hilbert-Schmidt Independence Criterion* (HSIC) is a kernel-based dependence measure that quantifies the statistical dependence between two random variables. Let X and Y be two random variables with $k(\cdot, \cdot)$ and $l(\cdot, \cdot)$ as reproducing (product) kernels defined on them. Then, $\text{HSIC}(X, Y) :=$

¹Fleissner et al. (2024) designed an explanation algorithm for kernel-based unsupervised learning.

270 $\|\mathcal{C}_{XY}\|_{\text{HS}}^2$, where \mathcal{C}_{XY} is the cross-covariance operator and $\|\cdot\|_{\text{HS}}$ is the Hilbert-Schmidt (HS) norm;
 271 see, e.g., (Muandet et al., 2017, Sec. 3.6) for technical details. Given a sample $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathbb{P}(X, Y)$,
 272 HSIC(X, Y) can be estimated as $\widehat{\text{HSIC}}(X, Y) = (n-1)^{-2} \text{tr}(\mathbf{KHLH})$ where $\mathbf{K} \in \mathbb{R}^{n \times n}$
 273 is the kernel matrix computed using the kernel k , i.e., $\mathbf{K}_{ij} = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$, $\mathbf{L} \in \mathbb{R}^{n \times n}$ is the kernel
 274 matrix computed using the kernel l , i.e., $\mathbf{L}_{ij} = l(\mathbf{y}^{(i)}, \mathbf{y}^{(j)})$, and $\mathbf{H} = I - \frac{1}{n} \mathbf{1}\mathbf{1}^\top$ is the centering
 275 matrix ensuring zero mean in the feature space.
 276

277 Computing HSIC(X, Y) gives us the overall dependence between X and Y . We now focus on
 278 scenarios where X represents a random variable of the feature vector \mathbf{x} and Y represents the scalar
 279 prediction outcome y . Our interest is then to distribute the overall dependence over the features in
 280 \mathbf{x} . This is particularly useful for feature selection and global sensitivity analysis, where the target is
 281 usually univariate. Analogous to the predictive function and MMD cases, we define the functional
 282 baseline value function of a coalition $\mathcal{S} \subseteq \mathcal{D}$ as the HSIC computed only with the kernel restricted to
 283 \mathcal{S} , with all other kernel components set to one. This isolates the dependence captured by \mathcal{S} alone:

284 **Definition 6.** For a coalition $\mathcal{S} \subseteq \mathcal{D}$, the functional baseline value function for HSIC is defined
 285 as $v_{\text{HSIC}}(\mathcal{S}) = \frac{1}{(n-1)^2} \text{tr}(\mathbf{HLH} \mathbf{K}_{\mathcal{S}})$, where $\mathbf{K}_{\mathcal{S}} = \odot_{j \in \mathcal{S}} \mathbf{K}_j$ is the Hadamard product of kernel
 286 matrices restricted to \mathcal{S} .
 287

288 This construction follows the same principle as before: setting kernel contributions of features outside
 289 \mathcal{S} to their least informative value (equal to 1 as we are doing multiplication). Thus $v_{\text{HSIC}}(\mathcal{S})$ quantifies
 290 the dependency between $X_{\mathcal{S}}$ and Y , independently of the remaining features. Applying the same
 291 trick as in Proposition 5 to the value function v_{HSIC} yields a similar recursive formulation of Shapley
 292 values for HSIC.

293 **Proposition 7.** For the product kernel, the Shapley value for HSIC can be recursively computed as:

$$294 \phi_j^{\text{HSIC}} = \frac{1}{(n-1)^2} \text{tr} \left(\mathbf{HLH} \left((\mathbf{K}_j - \mathbf{1}\mathbf{1}^\top) \odot \sum_{q=0}^{d-1} \mu(q) E_q(\mathcal{K}_{-j}) \right) \right),$$

295 where $\mathcal{K} = \{\mathbf{K}_1, \dots, \mathbf{K}_d\}$ with \mathbf{K}_i being the kernel matrix for feature i only, and $E_q(\mathcal{K}_{-j}) =$
 296 $\frac{1}{q} \sum_{r=1}^q (-1)^{r-1} E_{q-r}(\mathcal{K}_{-j}) \odot P_r(\mathcal{K}_{-j})$, is the ESPs with $P_r(\mathcal{K}) = \sum_{\mathbf{k}_i \in \mathcal{K}} \mathbf{K}_i^r$ being the element-
 297 wise degree- r power sum.
 298

299 The Shapley values ϕ_j^{HSIC} allow us to allocate the overall dependence between X and Y across
 300 individual features, identifying the most influential ones for prediction. Our results can be generalized
 301 to scenarios when both X and Y are multivariate; see Appendix G for more details. The attribution of
 302 overall dependence to the involved multivariate variables has other applications in statistical inference,
 303 e.g., kernel-based (conditional) independence testing (Gretton et al., 2007; Albert et al., 2022).
 304
 305
 306

307 5 EXPERIMENTS

308 We evaluate the effectiveness of PKeX-Shapley for product-kernel methods through a series of
 309 experiments conducted on a 24-core machine with 16GB RAM. We further provide the experimental
 310 setups, training procedure, and extra experiments in Appendix H.
 311
 312

313 5.1 EFFECTIVENESS OF RECURSION AND VALUE FUNCTION IN LOCAL EXPLANATIONS

314 **Experiment 1: Recursion vs. Regression Formulation.** We empirically demonstrate the advantage
 315 of PKeX-Shapley, compared to sampling-based approximations such as Kernel SHAP. We generate
 316 four synthetic datasets, each with 1000 samples and $d \in \{10, 20, 30, 50\}$ features. Features are
 317 independently drawn from a standard normal distribution, and labels are generated using a linear
 318 model with additive Gaussian noise ($\sigma = 0.1$). For each dataset, we train a support vector regression
 319 model with an RBF kernel and use the trained model to compare explanation methods.
 320

321 For a fair comparison between the recursive and regression-based methods, we adopt the same
 322 functional baseline value function as in equation 3. Specifically, we first compute the exact Shapley
 323 values ϕ_j^x using PKeX-Shapley, and then estimate the Shapley values using a regression-based
 approach analogous to Kernel SHAP, but modified to employ our value function. The estimator uses

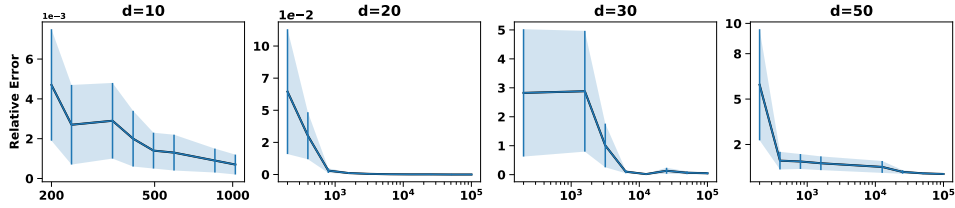


Figure 2: Relative estimation error of regression-based Shapley values versus the exact recursive values, shown across coalition sample sizes for feature dimensions $d = 10, 20, 30, 50$.

the paired coalition sampling scheme of Covert and Lee (2021), with sample sizes ranging from 200 to 10^5 , and computes $\hat{\phi}_j^{\mathbf{x}}$ by solving the corresponding weighted linear regression problem. To quantify the approximation error introduced by the regression formulation, we report the relative deviation error across 100 randomly selected instances, defined as $\sum_{j=1}^d |\phi_j^{\mathbf{x}} - \hat{\phi}_j^{\mathbf{x}}| / |\phi_j^{\mathbf{x}}|$.

Figure 2 plots the relative error as a function of sample size for each dataset over the 100 selected instances. When $d = 10$, the relative error is already around 0.005 with 200 coalition samples ($\approx 20\%$ of all possible coalitions). However, when $d = 50$, the approximation error remains above 9.0 with 200 samples, stays above 1.0 with 10^4 , and only approaches 0.05 near 10^5 coalition samples. These results highlight that as the number of features increases, a substantially larger number of samples is required to obtain reliable estimates—rendering sampling-based methods unreliable in high-dimensional settings. This, in turn, supports the importance of having a fast and exact computation algorithm.

Experiment 2: Effectiveness of the Functional Baseline Value Function. Recall that PKeX-Shapley employs a functional baseline value function, in contrast to standard expectation-based value functions. It is therefore natural to ask whether this formulation meaningfully captures the importance of feature subsets. To evaluate this, we assess the quality of the resulting explanations in recovering the most informative feature on synthetic datasets, and compare PKeX-Shapley against alternative explanation methods. We generate three regression tasks of $n = 1000$ samples in \mathbb{R}^{50} , where only the first one-third of features (denoted by \mathcal{S} , $|\mathcal{S}| = 17$) drive the target, and the remaining 33 features are redundant. The three target functions over \mathcal{S} are: a degree-5 polynomial, a degree-10 polynomial, and a squared-exponential response $y = \exp(\sum_{i \in \mathcal{S}} x_i^2)$. We train a support vector regressor with an RBF kernel on each dataset, and produce explanations using our exact recursive method alongside three baselines: RKHS-SHAP (Chau et al., 2022), GEMFIX (Mohammadi et al., 2025a), BiSHAP (Masoomi et al., 2021), and Sampling SHAP (Štrumbelj and Kononenko, 2014), each configured with 500 and 1000 coalition samples. All methods employ a fixed background set of 100 points to estimate their value functions.

Attribution accuracy is computed over 100 independent test instances by selecting the top-17 features returned by each method and measuring the fraction of true active features recovered. Figure 3 shows the average accuracy rate for each method across the three tasks. PKeX-Shapley achieves a competitive or superior performance in all cases, whereas Kernel SHAP, GEMFIX, and Sampling SHAP suffer accuracy degradation as the target function’s complexity increases (most notably for the degree-10 polynomial and the exponential model).

We also measure per-instance explanation runtime. Figure 4 presents error-bar plots (mean \pm standard deviation) of execution times in seconds. With 500 coalition samples, all methods incur comparable execution times. When the sample size increases to 1000, PKeX-Shapley remains significantly faster than the baselines, despite using the same background-sample budget of 100 for other methods. This demonstrates that PKeX-Shapley not only provides exact attributions but also outperforms sampling-based estimators in computational efficiency as coalition sample counts grow.

5.2 EXPLAINING DISTRIBUTION DISCREPANCY USING MMD WITH PKE-X-SHAPLEY

To illustrate how PKeX-Shapley can explain distributional discrepancies measured by MMD, we conduct two synthetic experiments following the standard two-sample testing setup (Schrab, 2025). Our goal is to attribute the observed MMD between two distributions to individual input variables. We present one of the synthetic experiments below, with the remaining experiments provided in Appendix H.3. In all MMD experiments, we use the RBF kernel with the bandwidth selected via the

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

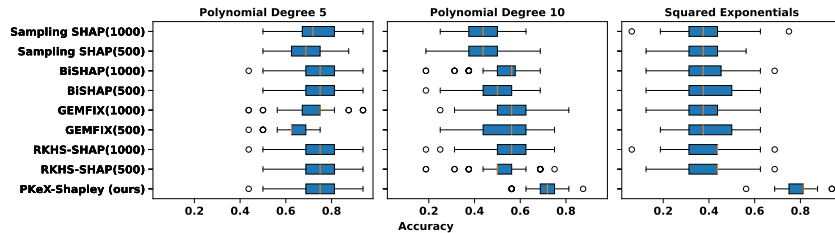


Figure 3: Recovery rate of true active features by each method on the three synthetic tasks.

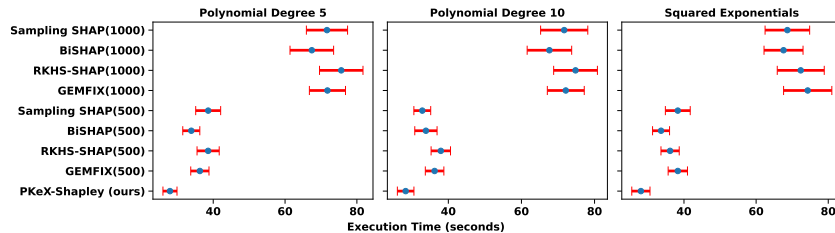


Figure 4: Per-instance explanation time (mean \pm standard deviation) for each method with 500 and 1000 coalition samples.

median heuristic method. We generate datasets X and Z , each comprising 20 variables. The first ten variables X_1, \dots, X_{10} and Z_1, \dots, Z_{10} are sampled from the same multivariate normal distribution, ensuring identical distributions. The remaining 10 variables X_{11}, \dots, X_{20} differ, with variables in X sampled from a multivariate normal distribution and those in Z from a Student’s t-distribution of the same mean. This introduces differences in higher-order moments and covariance structure, resulting in a measurable discrepancy.

We generate 1,000 samples from each distribution and compute the MMD. Shapley values are computed to quantify the contribution of each variable to the overall MMD. To ensure robustness, the experiment was repeated 1,000 times, and kernel density estimates (KDE) of the Shapley values are presented in Figure 5.

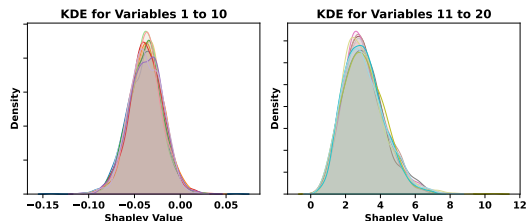


Figure 5: KDE plots of Shapley values for the synthetic dataset. Variables 1–10 reduce MMD via negative contributions, while 11–20 increase it via positive contributions.

The left subplot in Figure 5 displays the KDE plots for variables 1–10, while the right subplot corresponds to variables 11–20. For the first 10 variables, the consistently negative Shapley values indicate a reduction in the overall MMD, effectively “pulling” the discrepancy closer to zero. This aligns with our intuition: since these variables are identically distributed across both datasets, they should not contribute positively to the observed distributional difference. In contrast, variables 11–20 exhibit positive Shapley values, reflecting their contribution to the increase in MMD and, thus, their role in capturing the divergence between the distributions. Moreover, the KDE plots of the Shapley values within each group (variables 1–10 and 11–20, respectively) are nearly identical, which is consistent with the symmetric construction of the synthetic data. This experiment illustrates that Shapley values provide meaningful insights into the contribution of individual variables to distributional discrepancies measured by MMD.

5.3 EXPLAINING HSIC WITH PKE X-SHAPLEY: FEATURE SELECTION CASE STUDY

Lastly, we demonstrate how attributing the HSIC between input features and the target variable to individual features can support feature selection by quantifying their respective contributions to the overall statistical dependence. As a comparison, we also compare our approach with five feature importance methods: HSICLasso (Climente-González et al., 2019), Mutual Information (MI), Lasso, K-Best, and Tree Ensemble (with feature permutations). Our experiments were conducted on seven

Table 1: Performance (mean±standard deviation) when training on the top 20% of features. Datasets *breast cancer*, *skillcraft*, *sml*, and *parkinson* are regression (MAPE, lower is better); *sonar*, *Wisconsin*, *ionosphere* are classification (accuracy, higher is better).

| Method | sonar | Wisconsin | ionosphere | breast cancer | skillcraft | sml | parkinson |
|--------------|--------------------|--------------------|------------------------------------|--------------------|--------------------|--------------------|--------------------|
| accuracy (↑) | | | mean absolute percentage error (↓) | | | | |
| PKeX-Shapley | 0.808±0.030 | 0.909±0.015 | 0.878±0.036 | 0.850±0.010 | 1.000±0.020 | 0.999±0.001 | 0.125±0.022 |
| HSICLasso | 0.808±0.044 | 0.884±0.010 | 0.912±0.035 | 1.000±0.080 | 2.175±0.429 | 1.354±0.726 | 1.170±0.269 |
| MI | 0.875±0.053 | 0.900±0.015 | 0.937±0.023 | 1.000±0.080 | 1.134±0.099 | 0.196±0.071 | 0.214±0.004 |
| Lasso | 0.842±0.038 | 0.900±0.024 | 0.932±0.021 | 1.000±0.080 | 1.821±0.680 | 1.000±0.000 | 1.000±0.000 |
| K-Best | 0.779±0.040 | 0.909±0.015 | 0.869±0.018 | 1.000±0.080 | 1.134±0.099 | 0.257±0.086 | 1.018±0.076 |
| Tree Ens. | 0.837±0.059 | 0.887±0.033 | 0.926±0.031 | 1.000±0.080 | 2.175±0.429 | 1.000±0.000 | 0.214±0.004 |

datasets, where we trained Gaussian process (GP) models with an RBF kernel, using only the top 20% of features ranked by each selection method.

Table 1 presents the results, reporting the five-fold cross-validated mean and standard deviation for the GP models trained on the selected features. We use mean absolute percentage error (MAPE) as the performance metric for regression tasks and accuracy for classification tasks. For kernel computation in HSIC, we use an RBF kernel for features and regression targets, and a categorical kernel for classification targets. The bandwidth for the RBF kernel is selected using the median heuristic. PKeX-Shapley consistently delivers strong results across all datasets. On regression problems, it yields better MAPE on datasets *breast cancer*, *skillcraft*, and *parkinson*, while other methods often incur higher error (e.g., HSICLasso on *skillcraft* and *parkinson*) or greater variance (e.g., HSICLasso on *sml*). For classification, PKeX-Shapley maintains accuracies above 80%, matching or exceeding the baselines. It achieves the best-performing results on the *Wisconsin* datasets, and remains competitive on the other two datasets. Interestingly, when compared to HSICLasso—a method specifically tailored for feature selection—PKeX-Shapley demonstrates superior performance across the majority of datasets. Notably, we achieve better results in 5 out of the 7 datasets, tie in one, and only fall short in the *ionosphere* dataset. This is particularly noteworthy, as PKeX-Shapley is not explicitly designed for feature selection, yet it consistently outperforms a specialized method like HSICLasso.

6 CONCLUSION, LIMITATION, AND DISCUSSION

This work introduces PKeX-Shapley, a polynomial-time algorithm for computing exact Shapley values for product-kernel methods. We introduce the functional baseline value function for product-kernel methods, and show that it naturally induces a functional decomposition, which we exploit to develop a recursive algorithm that bypass the exponential complexity of naïve Shapley value computation. This approach enables exact, efficient feature attribution for both predictive models and kernel-based statistical discrepancies, including Maximum Mean Discrepancy (MMD) and Hilbert-Schmidt Independence Criterion (HSIC), providing interpretability of distributional differences and dependence structures. Our method achieves quadratic-time complexity and eliminates the approximation errors inherent in sampling-based estimators (and expectation-based value functions), as demonstrated through experiments on both synthetic and real-world datasets.

While our algorithm reduces computational complexity from exponential to quadratic time, its main limitation is that it applies only to product kernels. This trade-off is largely unavoidable: achieving tractable computation requires imposing structural constraints. As part of future work, we plan to explore whether further relaxation is possible. Another promising direction is to extend our approach to higher-order attribution methods, such as Shapley interaction indices (Sundararajan et al., 2020; Muschalik et al., 2024). It is also of interest to investigate how our computational techniques can speed up other explanation techniques, such as partial dependence plots and integrated gradients.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

REPRODUCIBILITY STATEMENT

We have made all code and implementation details publicly available (pke, 2025). The experimental setup, including dataset preprocessing, model architectures, training procedures, and hyperparameter settings, are fully described in either the main paper or the appendix. This ensures that all reported results can be independently verified.

REFERENCES

- Lloyd S Shapley. A value for n-person games. 1953.
- Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017.
- Ian Covert, Scott M Lundberg, and Su-In Lee. Understanding global feature contributions with additive importance measures. *Advances in Neural Information Processing Systems*, 33:17212–17223, 2020.
- Mukund Sundararajan and Amir Najmi. The many shapley values for model explanation. In *International conference on machine learning*, pages 9269–9278. PMLR, 2020.
- 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):56–67, 2020.
- Siu Lun Chau, Krikamol Muandet, and Dino Sejdinovic. Explaining the uncertain: Stochastic shapley values for gaussian process models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Siu Lun Chau, Robert Hu, Javier Gonzalez, and Dino Sejdinovic. Rkhs-shap: Shapley values for kernel methods. *Advances in neural information processing systems*, 35:13050–13063, 2022.
- Arthur Gretton, Karsten Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alex Smola. A kernel method for the two-sample-problem. *Advances in neural information processing systems*, 19, 2006.
- Masha Naslidnyk, Siu Lun Chau, Francois-Xavier Briol, and Krikamol Muandet. Kernel quantile embeddings and associated probability metrics. In *Forty-second International Conference on Machine Learning*, 2025.
- Antonin Schrab. A unified view of optimal kernel hypothesis testing. *arXiv preprint arXiv:2503.07084*, 2025.
- Siu Lun Chau, Antonin Schrab, Arthur Gretton, Dino Sejdinovic, and Krikamol Muandet. Credal two-sample tests of epistemic uncertainty. In *International Conference on Artificial Intelligence and Statistics*, pages 127–135. PMLR, 2025.
- Kacper Chwialkowski, Heiko Strathmann, and Arthur Gretton. A kernel test of goodness of fit. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 2606–2615, 2016.
- Qiang Liu, Jason Lee, and Michael Jordan. A kernelized stein discrepancy for goodness-of-fit tests. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48, pages 276–284, 2016.
- Mélisande Albert, Béatrice Laurent, Amandine Marrel, and Anouar Meynaoui. Adaptive test of independence based on hsc measures. *The Annals of Statistics*, 50(2):858–879, 2022.
- Jovana Mitrovic, Dino Sejdinovic, and Yee Whye Teh. Causal inference via kernel deviance measures. *Advances in neural information processing systems*, 31, 2018.
- Siu Lun Chau, Jean-Francois Ton, Javier González, Yee Teh, and Dino Sejdinovic. Bayesimp: Uncertainty quantification for causal data fusion. *Advances in Neural Information Processing Systems*, 34:3466–3477, 2021.

540 Dino Sejdinovic. An overview of causal inference using kernel embeddings. *arXiv preprint*
541 *arXiv:2410.22754*, 2024.

542 Stan Lipovetsky and Michael Conklin. Analysis of regression in game theory approach. *Applied*
543 *stochastic models in business and industry*, 17(4):319–330, 2001.

544 I Elizabeth Kumar, Suresh Venkatasubramanian, Carlos Scheidegger, and Sorelle Friedler. Problems
545 with shapley-value-based explanations as feature importance measures. In *International conference*
546 *on machine learning*, pages 5491–5500. PMLR, 2020.

547 Krikamol Muandet, Kenji Fukumizu, Bharath Sriperumbudur, Bernhard Schölkopf, et al. Kernel
548 mean embedding of distributions: A review and beyond. *Foundations and Trends® in Machine*
549 *Learning*, 10(1-2):1–141, 2017.

550 Zoltán Szabó and Bharath K Sriperumbudur. Characteristic and universal tensor product kernels.
551 *Journal of Machine Learning Research*, 18(233):1–29, 2018.

552 Pkex-shapley: Reproducibility code for neurips 2025 submission. [https://anonymous.4open.](https://anonymous.4open.science/r/PKeX-Shapley-B1C0/)
553 [science/r/PKeX-Shapley-B1C0/](https://anonymous.4open.science/r/PKeX-Shapley-B1C0/), 2025.

554 Dominik Janzing, Lenon Minorics, and Patrick Blöbaum. Feature relevance quantification in explain-
555 able ai: A causal problem. In *International Conference on artificial intelligence and statistics*,
556 pages 2907–2916. PMLR, 2020.

557 Christoph Molnar. *Interpreting Machine Learning Models with SAP: A Guide with Python Examples*
558 *and Theory on Shapley Values*. Christoph Molnar c/o MUCBOOK, Heidi Seibold, 2023.

559 Jacob Gardner, Geoff Pleiss, Ruihan Wu, Kilian Weinberger, and Andrew Wilson. Product kernel
560 interpolation for scalable gaussian processes. In *International Conference on Artificial Intelligence*
561 *and Statistics*, pages 1407–1416. PMLR, 2018.

562 Eric S Egge. *An introduction to symmetric functions and their combinatorics*, volume 91. American
563 *Mathematical Soc.*, 2019.

564 David K Duvenaud, Hannes Nickisch, and Carl Rasmussen. Additive gaussian processes. *Advances*
565 *in neural information processing systems*, 24, 2011.

566 Maximilian Fleissner, Leena Chennuru Vankadara, and Debarghya Ghoshdastidar. Explaining kernel
567 clustering via decision trees. In *The Twelfth International Conference on Learning Representations*,
568 2024.

569 Arthur Gretton, Kenji Fukumizu, Choon Teo, Le Song, Bernhard Schölkopf, and Alex Smola. A
570 kernel statistical test of independence. *Advances in neural information processing systems*, 20,
571 2007.

572 Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A
573 kernel two-sample test. *The Journal of Machine Learning Research*, 13(1):723–773, 2012.

574 Wenju Zhang, Xiang Zhang, Long Lan, and Zhigang Luo. Maximum mean and covariance discrep-
575 ancy for unsupervised domain adaptation. *Neural Processing Letters*, 51:347–366, 2020.

576 Ian Covert and Su-In Lee. Improving kernelshap: Practical shapley value estimation using linear
577 regression. In *International Conference on Artificial Intelligence and Statistics*, pages 3457–3465.
578 PMLR, 2021.

579 Majid Mohammadi, Ilaria Tiddi, and Annette Ten Teije. Unlocking the game: Estimating games in
580 möbius representation for explanation and high-order interaction detection. In *Proceedings of the*
581 *AAAI Conference on Artificial Intelligence*, volume 39, pages 19512–19519, 2025a.

582 Aria Masoomi, Davin Hill, Zhonghui Xu, Craig P Hersh, Edwin K Silverman, Peter J Castaldi,
583 Stratis Ioannidis, and Jennifer Dy. Explanations of black-box models based on directional feature
584 interactions. In *International Conference on Learning Representations*, 2021.

585 Erik Štrumbelj and Igor Kononenko. Explaining prediction models and individual predictions with
586 feature contributions. *Knowledge and information systems*, 41:647–665, 2014.

594 Héctor Climente-González, Chloé-Agathe Azencott, Samuel Kaski, and Makoto Yamada. Block hsic
595 lasso: model-free biomarker detection for ultra-high dimensional data. *Bioinformatics*, 35(14):
596 i427–i435, 2019.

597 Mukund Sundararajan, Kedar Dhamdhere, and Ashish Agarwal. The shapley taylor interaction index.
598 In *International conference on machine learning*, pages 9259–9268. PMLR, 2020.

600 Maximilian Muschalik, Hubert Baniecki, Fabian Fumagalli, Patrick Kolpaczki, Barbara Hammer,
601 and Eyke Hüllermeier. shapiq: Shapley interactions for machine learning. *Advances in Neural*
602 *Information Processing Systems*, 37:130324–130357, 2024.

603 Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna:
604 A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM*
605 *SIGKDD international conference on knowledge discovery & data mining*, pages 2623–2631,
606 2019.

607 Majid Mohammadi, Krikamol Muandet, Ilaria Tiddi, Annette Ten Teije, and Siu Lun Chau.
608 Exact shapley attributions in quadratic-time for fanova gaussian processes. *arXiv preprint*
609 *arXiv:2508.14499*, 2025b.
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

APPENDIX

The appendix provides additional information and proofs related to the material presented in the main paper. It includes detailed explanations, proofs, algorithms, and experiments relevant to explaining product-kernel models. The structure of the appendix is as follows:

- **Product kernels in kernel methods:** §A discusses the construction and examples of product kernels, including their definitions and properties.
- **Newton’s identities :** §C presents the main result of Newton’s identities for elementary symmetric polynomials.
- **Proofs:** §D provides detailed proofs of key theorems and lemmas used in the main paper.
- **Recursive and numerically stable algorithms for computing Shapley values for product-kernel learning models:** §E describes a recursive algorithm for computing Shapley values, as well as an adjusted algorithm for numerical stability and efficiency.
- **Additivity of explanations for learning models, MMD and HSIC:** §F discusses the additivity of explanations for the product-kernel methods studied in this paper.
- **Shapley value attribution for HSIC with two multivariate variables** §G discusses the Shapley value attribution for HSIC when both random variables are multivariate.
- **Experiments:** §H reports experimental setup and extra experiments on explaining product-kernel models.

A PRODUCT KERNELS IN KERNEL METHODS

Product kernels provide a powerful mechanism for constructing high-dimensional similarity measures by combining kernels defined on individual dimensions or feature subsets. This section discusses key examples of product kernels.

Radial Basis Function (RBF) Kernels as Product Kernels The RBF kernel is a canonical example of product kernels. The RBF kernel is defined based on a distance metric between two instances, with the two well-known metrics being Euclidean (L2 norm) and Manhattan distances (L1 norm). We refer to the former as RBF, and the latter as Laplacian RBF to distinguish these two kernel functions. In addition, when we have only one kernel bandwidth parameter σ , the RBF kernel is referred to as *isotropic*. The RBF kernel with both distance metrics inherently decomposes into products of univariate kernels across dimensions:

- **RBF Kernel:**

$$K_{\text{RBF}}(\mathbf{x}, \mathbf{z}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{z}\|^2}{2\sigma^2}\right) = \prod_{i=1}^d \exp\left(-\frac{(x_i - z_i)^2}{2\sigma^2}\right).$$

- **Laplacian RBF Kernel:**

$$K_{\text{Laplacian RBF}}(\mathbf{x}, \mathbf{z}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{z}\|_1}{\sigma}\right) = \prod_{i=1}^d \exp\left(-\frac{|x_i - z_i|}{\sigma}\right).$$

When alternative distance metrics are incorporated into the RBF kernel, such as the Mahalanobis distance, which involves a covariance matrix, the resulting kernel might lose its product decomposition.

Automatic Relevance Determination (ARD) in Gaussian Processes ARD extends RBF kernels by assigning independent length-scale parameters σ_i to each dimension:

$$K_{\text{ARD}}(\mathbf{x}, \mathbf{z}) = \exp\left(-\sum_{i=1}^d \frac{(x_i - z_i)^2}{2\sigma_i^2}\right) = \prod_{i=1}^d \exp\left(-\frac{(x_i - z_i)^2}{2\sigma_i^2}\right).$$

The ARD is extensively used in Gaussian processes for feature selection via learned σ_i , and to enhance interpretability and adaptability. ARD is also referred to as *anisotropic*.

Cauchy Kernel The Cauchy kernel, inspired by the Cauchy distribution, is another example of a product kernel:

$$K_{\text{Cauchy}}(\mathbf{x}, \mathbf{z}) = \prod_{i=1}^d \frac{1}{1 + \frac{(x_i - z_i)^2}{\sigma^2}}.$$

Product of Base Kernels A popular way for constructing product kernels is by first defining a base kernel for each individual feature and then computing the overall kernel function by multiplying the base kernels over individual features. The product of PSD kernels remains PSD by the Schur product theorem:

$$K(\mathbf{x}, \mathbf{z}) = \prod_{i=1}^d K_i(x_i, z_i).$$

This type of kernel introduces flexibility in combining base kernels while maintaining validity as a (product) kernel function.

B VALUE FUNCTIONS FOR PRODUCT KERNELS: BASELINE AND INTERVENTIONAL

We now discuss the interventional and baseline value functions for product-kernel methods and their relevance and affinity with the proposed value function. For a product-kernel model, the decision function can be expressed as

$$f(\mathbf{x}) = \boldsymbol{\alpha}^\top (k_{\mathcal{S}}(\mathbf{X}_{\mathcal{S}}, \mathbf{x}_{\mathcal{S}}) \odot k_{\mathcal{S}^c}(\mathbf{X}_{\mathcal{S}^c}, \mathbf{x}_{\mathcal{S}^c})),$$

where \odot denotes the Hadamard product, $k_{\mathcal{S}}$ is the product kernel restricted to subset \mathcal{S} , and \mathcal{S}^c is its complement.

B.1 BASELINE VALUE FUNCTION

The baseline value function removes the effect of features outside the coalition \mathcal{S} by replacing them with fixed baseline values. Formally, let $\mathbf{x}_{\mathcal{S}^c}^b$ denote the baseline values for the excluded features. Then the baseline value function is defined as

$$\begin{aligned} v_{\mathbf{x}}^{\text{base}}(\mathcal{S}) &= f([\mathbf{x}_{\mathcal{S}}, \mathbf{x}_{\mathcal{S}^c}^b]) \\ &= \boldsymbol{\alpha}^\top \mathbf{K}([\mathbf{x}_{\mathcal{S}}, \mathbf{x}_{\mathcal{S}^c}^b], \mathbf{X}) = \boldsymbol{\alpha}^\top (k_{\mathcal{S}}(\mathbf{X}_{\mathcal{S}}, \mathbf{x}_{\mathcal{S}}) \odot k_{\mathcal{S}^c}(\mathbf{X}_{\mathcal{S}^c}, \mathbf{x}_{\mathcal{S}^c}^b)). \end{aligned}$$

In model-agnostic settings, the choice of baseline values is often application-dependent (e.g., mean input, zero vector, or reference sample). For product-kernel methods, however, we have an explicit multiplicative structure. Instead of fixing raw feature values, we can directly mask the kernel factors corresponding to excluded features. Specifically, setting a kernel factor to one corresponds to the least informative contribution, since it makes the kernel computation independent of that feature. Thus, the baseline operation in product kernels amounts to replacing $k_j(x_j, x_j^{(i)})$ with 1 for all $j \in \mathcal{S}^c$. This yields

$$v_{\mathbf{x}}(\mathcal{S}) = \boldsymbol{\alpha}^\top k_{\mathcal{S}}(\mathbf{X}_{\mathcal{S}}, \mathbf{x}_{\mathcal{S}}),$$

which is exactly the *functional baseline value function*.

In words, *rather than choosing baseline feature values in input space, product-kernel methods define the baseline at the level of the kernel: excluded features are masked by setting their kernel contributions to 1*. This ensures that only features in \mathcal{S} influence the similarity measure.

B.2 INTERVENTIONAL VALUE FUNCTION

An alternative approach is the *interventional value function*, which removes the effect of features outside \mathcal{S} by replacing them with their *average contribution* under the data distribution. Formally, it

756 is defined as
757

$$758 \quad v_{\mathbf{x}}^{\text{int}}(\mathcal{S}) = \mathbb{E}_{\mathbf{x}_{S^c}} [f(\mathbf{x}) \mid \mathbf{x}_S] \\ 759 \quad \quad \quad = \mathbb{E}_{\mathbf{x}_{S^c}} [\boldsymbol{\alpha}^\top k([\mathbf{x}_S, \mathbf{x}_{S^c}], \mathbf{X})] \\ 760 \quad \quad \quad = \boldsymbol{\alpha}^\top \mathbb{E}_{\mathbf{x}_{S^c}} [k([\mathbf{x}_S, \mathbf{x}_{S^c}], \mathbf{X})].$$

761 Exploiting the product structure of the kernel, this simplifies to

$$762 \quad v_{\mathbf{x}}^{\text{int}}(\mathcal{S}) = \boldsymbol{\alpha}^\top \left(k_S(\mathbf{x}_S, \mathbf{X}_S) \odot \mathbb{E}_{\mathbf{x}_{S^c}} [k_{S^c}(\mathbf{x}_{S^c}, \mathbf{X}_{S^c})] \right). \quad (5)$$

763 To summarize, *the interventional value function removes a feature's contribution by averaging it out,*
764 *while our baseline-inspired value function directly eliminates the feature's influence by setting its*
765 *kernel component to one.*

766 This highlights the conceptual and practical difference: the interventional function is defined through
767 expectations with respect to the data distribution, whereas our proposed value function exploits the
768 explicit multiplicative decomposition of product kernels and provides an efficient, distribution-free
769 alternative.

770 C NEWTON'S IDENTITIES

771 To explain Newton's identities, we begin with a specific set of variables $\mathcal{Z}_4 = \{z_1, z_2, z_3, z_4\}$ before
772 generalizing it to sets of arbitrary size $\mathcal{Z}_d = \{z_1, z_2, \dots, z_d\}$. The *elementary symmetric polynomials*
773 (ESPs) of degree q is defined as

$$774 \quad e_q(\mathcal{Z}_4) = \sum_{1 \leq i_1 < i_2 < \dots < i_q \leq 4} z_{i_1} z_{i_2} \dots z_{i_q},$$

775 with the conventions $e_0(\mathcal{Z}_4) = 1$ and $e_q(\mathcal{Z}_4) = 0$ for $q > 4$. For example:

$$776 \quad e_1(\mathcal{Z}_4) = z_1 + z_2 + z_3 + z_4, \\ 777 \quad e_2(\mathcal{Z}_4) = z_1 z_2 + z_1 z_3 + z_1 z_4 + z_2 z_3 + z_2 z_4 + z_3 z_4, \\ 778 \quad e_4(\mathcal{Z}_4) = z_1 z_2 z_3 z_4.$$

779 The *power sum* of degree r is given by

$$780 \quad p_r(\mathcal{Z}_4) = z_1^r + z_2^r + z_3^r + z_4^r.$$

781 In particular, $p_1(\mathcal{Z}_4) = e_1(\mathcal{Z}_4)$ and, for example, $p_2(\mathcal{Z}_4) = z_1^2 + z_2^2 + z_3^2 + z_4^2$. Then, Newton's
782 identities relate the ESPs to the power sum recursively. For $q \geq 1$,

$$783 \quad e_q(\mathcal{Z}_4) = \frac{1}{q} \sum_{r=1}^q (-1)^{r-1} e_{q-r}(\mathcal{Z}_4) p_r(\mathcal{Z}_4).$$

784 For the set $\mathcal{Z}_4 = \{z_1, z_2, z_3, z_4\}$, the identities yield:

$$785 \quad e_1(\mathcal{Z}_4) = \frac{1}{1} [e_0(\mathcal{Z}_4) p_1(\mathcal{Z}_4)] = p_1(\mathcal{Z}_4) = z_1 + z_2 + z_3 + z_4, \\ 786 \quad e_2(\mathcal{Z}_4) = \frac{1}{2} [e_1(\mathcal{Z}_4) p_1(\mathcal{Z}_4) - e_0(\mathcal{Z}_4) p_2(\mathcal{Z}_4)] = \frac{(z_1 + z_2 + z_3 + z_4)^2 - (z_1^2 + z_2^2 + z_3^2 + z_4^2)}{2}, \\ 787 \quad e_3(\mathcal{Z}_4) = \frac{1}{3} [e_2(\mathcal{Z}_4) p_1(\mathcal{Z}_4) - e_1(\mathcal{Z}_4) p_2(\mathcal{Z}_4) + e_0(\mathcal{Z}_4) p_3(\mathcal{Z}_4)], \\ 788 \quad e_4(\mathcal{Z}_4) = \frac{1}{4} [e_3(\mathcal{Z}_4) p_1(\mathcal{Z}_4) - e_2(\mathcal{Z}_4) p_2(\mathcal{Z}_4) + e_1(\mathcal{Z}_4) p_3(\mathcal{Z}_4) - e_0(\mathcal{Z}_4) p_4(\mathcal{Z}_4)].$$

789 The identities presented above can be extended to sets of arbitrary size $\mathcal{Z}_d = \{z_1, z_2, \dots, z_d\}$ as

$$790 \quad e_q(\mathcal{Z}_d) = \frac{1}{q} \sum_{r=1}^q (-1)^{r-1} e_{q-r}(\mathcal{Z}_d) p_r(\mathcal{Z}_d), \quad \text{for } q \geq 1,$$

791 with $e_0(\mathcal{Z}_d) = 1$ and $e_q(\mathcal{Z}_d) = 0$ if $q > d$ or $q < 0$.

D PROOFS

This section provides the detailed proofs of the theoretical results presented in the main paper. First of all, we present a theorem that plays a key role in the other proofs.

D.1 PROOF OF THEOREM 2

We can write the Shapley value as follows:

$$\phi_j^{\mathbf{x}} := \sum_{S \subseteq \mathcal{D} \setminus \{j\}} \mu(|S|) (v_{\mathbf{x}}(S \cup \{j\}) - v_{\mathbf{x}}(S)) = \sum_{q=0}^{d-1} \mu(q) \sum_{\substack{S \subseteq \mathcal{D} \setminus \{j\} \\ |S|=q}} (v_{\mathbf{x}}(S \cup \{j\}) - v_{\mathbf{x}}(S)).$$

We now substitute $v_{\mathbf{x}}(S)$ for the product kernel into the Shapley value formula, one gets:

$$\begin{aligned} \phi_j^{\mathbf{x}} &= \boldsymbol{\alpha}^\top \left(\sum_{q=0}^{d-1} \mu(q) \sum_{\substack{S \subseteq \mathcal{D} \setminus \{j\} \\ |S|=q}} k_{S \cup \{j\}}(\mathbf{X}_{S \cup \{j\}}, \mathbf{x}_{S \cup \{j\}}) - k_S(\mathbf{X}_S, \mathbf{x}_S) \right) \\ &= \boldsymbol{\alpha}^\top \left(\sum_{q=0}^{d-1} \mu(q) \sum_{\substack{S \subseteq \mathcal{D} \setminus \{j\} \\ |S|=q}} k_j(\mathbf{X}_j, x_j) \odot k_S(\mathbf{X}_S, \mathbf{x}_S) - k_S(\mathbf{X}_S, \mathbf{x}_S) \right) \\ &= \boldsymbol{\alpha}^\top \left(\left(k_j(\mathbf{X}_j, x_j) - \mathbf{1} \right) \odot \left(\sum_{q=0}^{d-1} \mu(q) \sum_{\substack{S \subseteq \mathcal{D} \setminus \{j\} \\ |S|=q}} k_S(\mathbf{X}_S, \mathbf{x}_S) \right) \right), \end{aligned} \quad (6)$$

where $\mathbf{1}$ is a vector of one. Let $\mathbf{z}_i = k_i(\mathbf{X}_i, x_i)$ and $\mathcal{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_d\}$, one can write:

$$\sum_{\substack{S \subseteq \mathcal{D} \setminus \{j\} \\ |S|=q}} k_S(\mathbf{X}, \mathbf{x}) = \sum_{\substack{S \subseteq \mathcal{D} \setminus \{j\} \\ |S|=q}} \odot_{i \in S} k_i(\mathbf{X}, \mathbf{x}) = \sum_{\substack{1 \leq i_1 < i_2 < \dots < i_q \leq d-1 \\ j \notin \{i_1, \dots, i_q\}}} \mathbf{z}_{i_1} \odot \mathbf{z}_{i_2} \odot \dots \odot \mathbf{z}_{i_q} = e_q(\mathcal{Z}_{-j}),$$

where $e_q(\mathcal{Z}_{-j})$ is the elementary symmetric polynomials. This equation means that the inner sum can be recursively computed by Newton's identities formulation. It then follows:

$$\phi_j^{\mathbf{x}} = \boldsymbol{\alpha}^\top \left(\left(k_j(\mathbf{X}_j, x_j) - \mathbf{1} \right) \odot \sum_{q=0}^{d-1} \mu(q) e_q(\mathcal{Z}_{-j}) \right)$$

where

$$e_q(\mathcal{Z}_{-j}) = \frac{1}{q} \sum_{r=1}^q (-1)^{r-1} e_{q-r}(\mathcal{Z}_{-j}) \odot p_r(\mathcal{Z}_{-j}),$$

and $p_r(\mathcal{Z}) = \sum_{\mathbf{z}_i \in \mathcal{Z}} \mathbf{z}_i^r$ is the power sum, with the power working element-wise. This completes the proof. \square

D.2 PROOF OF LEMMA 3

By the efficiency property of Shapley values (Shapley, 1953), the sum of Shapley values equals the difference between the value of the grand coalition and the empty coalition:

$$\sum_{j=1}^d \phi_j^{\mathbf{x}} = v_{\mathbf{x}}(\mathcal{D}) - v_{\mathbf{x}}(\emptyset).$$

For the functional baseline value function, we have:

$$\begin{aligned} v_{\mathbf{x}}(\mathcal{D}) &= \boldsymbol{\alpha}^\top k_{\mathcal{D}}(\mathbf{X}_{\mathcal{D}}, \mathbf{x}_{\mathcal{D}}) = f(\mathbf{x}), \\ v_{\mathbf{x}}(\emptyset) &= \boldsymbol{\alpha}^\top k_{\emptyset}(\mathbf{X}_{\emptyset}, \mathbf{x}_{\emptyset}) = \boldsymbol{\alpha}^\top \mathbf{1} = \sum_{i=1}^n \alpha_i = f_{\emptyset}(\mathbf{x}). \end{aligned}$$

The result follows immediately by substitution.

864 D.3 PROOF OF PROPOSITION 5

865 By substituting $v_{\text{MMD}}(\mathcal{S})$ into the Shapley value formula, we obtain:

$$\begin{aligned}
866 \phi_q^{\text{MMD}} &= \sum_{r=0}^{d-1} \mu(r) \sum_{\substack{\mathcal{S} \subseteq \mathcal{D} \setminus \{q\} \\ |\mathcal{S}|=r}} \left(v_{\text{MMD}}(\mathcal{S} \cup \{q\}) - v_{\text{MMD}}(\mathcal{S}) \right) \\
867 &= \sum_{r=0}^{d-1} \mu(r) \sum_{\substack{\mathcal{S} \subseteq \mathcal{D} \setminus \{q\} \\ |\mathcal{S}|=r}} \left(\frac{1}{n(n-1)} \sum_{i \neq j} (k_q(x_q^{(i)}, x_q^{(j)}) - 1) k_{\mathcal{S}}(\mathbf{x}_{\mathcal{S}}^{(i)}, \mathbf{x}_{\mathcal{S}}^{(j)}) \right. \\
868 &\quad \left. + \frac{1}{m(m-1)} \sum_{i \neq j} (k_q(z_q^{(i)}, z_q^{(j)}) - 1) k_{\mathcal{S}}(\mathbf{z}_{\mathcal{S}}^{(i)}, \mathbf{z}_{\mathcal{S}}^{(j)}) \right. \\
869 &\quad \left. - \frac{2}{nm} \sum_{i,j} (k_q(x_q^{(i)}, z_q^{(j)}) - 1) k_{\mathcal{S}}(\mathbf{x}_{\mathcal{S}}^{(i)}, \mathbf{z}_{\mathcal{S}}^{(j)}) \right).
\end{aligned}$$

870 Let $\mathcal{Z}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \{k_1(x_1^{(i)}, x_1^{(j)}), \dots, k_d(x_d^{(i)}, x_d^{(j)})\}$ and define the elementary symmetric poly-
871 nomial as

$$872 e_r(\mathcal{Z}_{-q}^{(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})}) = \sum_{\substack{\mathcal{S} \subseteq \mathcal{D} \setminus \{q\} \\ |\mathcal{S}|=r}} k_{\mathcal{S}}(\mathbf{x}_{\mathcal{S}}^{(i)}, \mathbf{x}_{\mathcal{S}}^{(j)}).$$

873 Then, it follows from Newton's identities recurrence:

$$874 e_r(\mathcal{Z}_{-q}^{(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})}) = \frac{1}{r} \sum_{s=1}^r (-1)^{s-1} e_{r-s}(\mathcal{Z}_{-q}^{(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})}) p_s(\mathcal{Z}_{-q}^{(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})}),$$

875 where $p_s(\mathcal{Z}) = \sum_{z \in \mathcal{Z}} z^s$. Finally, we obtain

$$\begin{aligned}
876 \phi_q^{\text{MMD}} &= \frac{1}{n(n-1)} \sum_{i \neq j} \left((k_q(x_q^{(i)}, x_q^{(j)}) - 1) \sum_{r=0}^{d-1} \mu(r) e_r(\mathcal{Z}_{-q}^{(x_q^{(i)}, x_q^{(j)})}) \right) \\
877 &\quad + \frac{1}{m(m-1)} \sum_{i \neq j} \left((k_q(z_q^{(i)}, z_q^{(j)}) - 1) \sum_{r=0}^{d-1} \mu(r) e_r(\mathcal{Z}_{-q}^{(z_q^{(i)}, z_q^{(j)})}) \right) \\
878 &\quad - \frac{2}{nm} \sum_{i,j} \left((k_q(x_q^{(i)}, z_q^{(j)}) - 1) \sum_{r=0}^{d-1} \mu(r) e_r(\mathcal{Z}_{-q}^{(x_q^{(i)}, z_q^{(j)})}) \right),
\end{aligned}$$

879 and that completes the proof. \square

902 D.4 PROOF OF PROPOSITION 7

903 By substituting $v_{\text{HSIC}}(\mathcal{S}) = \frac{1}{(n-1)^2} \text{tr}(\mathbf{K}_{\mathcal{S}} \mathbf{H} \mathbf{L} \mathbf{H})$ into the Shapley value formula, we obtain:

$$\begin{aligned}
904 \phi_j^{\text{HSIC}} &= \frac{1}{(n-1)^2} \text{tr} \left(\mathbf{H} \mathbf{L} \mathbf{H} \sum_{q=0}^{d-1} \mu(q) \sum_{\substack{\mathcal{S} \subseteq \mathcal{D} \setminus \{j\} \\ |\mathcal{S}|=q}} \left(\mathbf{K}_{\mathcal{S} \cup \{j\}} - \mathbf{K}_{\mathcal{S}} \right) \right) \\
905 &= \frac{1}{(n-1)^2} \text{tr} \left(\mathbf{H} \mathbf{L} \mathbf{H} \sum_{q=0}^{d-1} \mu(q) \sum_{\substack{\mathcal{S} \subseteq \mathcal{D} \setminus \{j\} \\ |\mathcal{S}|=q}} \left(\mathbf{K}_j \odot \mathbf{K}_{\mathcal{S}} - \mathbf{K}_{\mathcal{S}} \right) \right) \\
906 &= \frac{1}{(n-1)^2} \text{tr} \left(\mathbf{H} \mathbf{L} \mathbf{H} (\mathbf{K}_j - \mathbf{1} \mathbf{1}^{\top}) \odot \sum_{q=0}^{d-1} \mu(q) \sum_{\substack{\mathcal{S} \subseteq \mathcal{D} \setminus \{j\} \\ |\mathcal{S}|=q}} \mathbf{K}_{\mathcal{S}} \right).
\end{aligned}$$

Letting $\mathcal{K} = \{\mathbf{K}_1, \dots, \mathbf{K}_d\}$, we express:

$$\sum_{\substack{\mathcal{S} \subseteq \mathcal{D} \setminus \{j\} \\ |\mathcal{S}|=q}} \mathbf{K}_{\mathcal{S}} = E_q(\mathcal{K}_{-j}),$$

which is computed recursively via Newton's identities formulation:

$$E_q(\mathcal{K}_{-j}) = \frac{1}{d-1} \sum_{r=1}^{d-1} (-1)^{r-1} E_{q-r}(\mathcal{K}_{-j}) \odot P_r(\mathcal{K}_{-j}),$$

where $P_r(\mathcal{K}) = \sum_{\mathbf{K}_i \in \mathcal{K}} \mathbf{K}_i^r$ is the element-wise power sum polynomial. Substituting this back, we obtain:

$$\phi_j^{\text{HSIC}} = \frac{1}{(n-1)^2} \text{tr} \left(\mathbf{H} \mathbf{L} \mathbf{H} \left((\mathbf{K}_j - \mathbf{1} \mathbf{1}^\top) \odot \sum_{q=0}^{d-1} \mu(q) E_q(\mathcal{K}_{-j}) \right) \right),$$

which completes the proof. \square

E RECURSIVE AND NUMERICALLY STABLE ALGORITHMS FOR COMPUTING SHAPLEY VALUES FOR PRODUCT-KERNEL LEARNING MODELS

Algorithm 1 Recursive Computation of Shapley Values for Product-Kernel Learning Models

Require: Trained model with product kernels (SVM/SVR/GP), instance $\mathbf{x} \in \mathbb{R}^d$

Ensure: Shapley values $\phi_1^{\mathbf{x}}, \dots, \phi_d^{\mathbf{x}}$ for each feature

- 1: Retrieve training data X , coefficients α , and kernel function k
 - 2: Compute feature-wise kernel vectors $\mathbf{z}_j = k(\mathbf{x}, X)$, $\forall j \in \{1, \dots, d\}$
 - 3: Precompute coefficients $\mu(q) = \frac{q!(d-q-1)!}{d!}$, $q = 0, \dots, d-1$
 - 4: **for** $j \in \{1, \dots, d\}$ **do**
 - 5: Let $\mathcal{Z}_{-j} \leftarrow \{\mathbf{z}_1, \dots, \mathbf{z}_d\} \setminus \{\mathbf{z}_j\}$
 - 6: Compute power sums $p_r(\mathcal{Z}_{-j}) = \sum_{\mathbf{z} \in \mathcal{Z}_{-j}} \mathbf{z}^r$ for $r = 1, \dots, d-1$
 - 7: Initialize $e_0 \leftarrow \mathbf{1}$
 - 8: **for** $q = 1$ **to** $d-1$ **do**
 - 9: $e_q(\mathcal{Z}_{-j}) \leftarrow \frac{1}{q} \sum_{r=1}^q (-1)^{k-1} e_{q-r}(\mathcal{Z}_{-j}) \odot p_r(\mathcal{Z}_{-j})$
 - 10: **end for**
 - 11: $\psi_j \leftarrow \sum_{q=0}^{d-1} \mu(q) \cdot e_q(\mathcal{Z}_{-j})$
 - 12: $\phi_j^{\mathbf{x}} \leftarrow \alpha^\top ((\mathbf{z}_j - \mathbf{1}) \odot \psi_j)$
 - 13: **end for**
 - 14:
 - 15: **return** $(\phi_1^{\mathbf{x}}, \dots, \phi_d^{\mathbf{x}})$
-

We present the algorithm for computing the Shapley values for product-kernel learning models. Let $\mathcal{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_d\}$ be a collection of kernel vectors. The elementary symmetric polynomial (ESP) of degree q is defined as:

$$e_q(\mathcal{Z}) = \sum_{1 \leq i_1 < \dots < i_q \leq d} \mathbf{z}_{i_1} \odot \dots \odot \mathbf{z}_{i_q}$$

Traditional computation uses Newton's identities:

$$e_q(\mathcal{Z}) = \frac{1}{q} \sum_{k=1}^q (-1)^{k-1} e_{q-k}(\mathcal{Z}) p_k(\mathcal{Z}), \quad p_k(\mathcal{Z}) = \sum_{\mathbf{z}_i \in \mathcal{Z}} \mathbf{z}_i^k$$

We used this recursion to compute Shapley values in product-kernel learning models. Algorithm 1 summarizes the overall algorithm based on this recursive formulation.

Algorithm 2 Numerically Stable Shapley Values Computation

Require: Trained model with product kernels (SVM/SVR/GP), instance $\mathbf{x} \in \mathbb{R}^d$

Ensure: Shapley values $\phi_1^{\mathbf{x}}, \dots, \phi_d^{\mathbf{x}}$ for each feature

```

1: Retrieve training data  $X$ , coefficients  $\alpha$ , and kernel function  $k$ 
2: Compute feature-wise kernel vectors  $\mathbf{z}_j = k(\mathbf{x}, X), \forall j \in \{1, \dots, d\}$ 
3: Precompute coefficients  $\mu(q) = \frac{q!(d-q-1)!}{d!}, q = 0, \dots, d-1$ 
4: for  $j = 1$  to  $d$  do
5:    $\mathcal{Z}_{-j} \leftarrow \{\mathbf{z}_1, \dots, \mathbf{z}_d\} \setminus \{\mathbf{z}_j\}$ 
6:   Scale  $\mathcal{Z}_{-j}$ :  $s \leftarrow \max_{i,j} \mathbf{z}_{ij} \forall \mathbf{z}_i \in \mathcal{Z}_{-j}$  (or 1 if all 0)
7:    $\tilde{\mathcal{Z}}_{-j} \leftarrow \{\mathbf{z}_i/s \mid \mathbf{z}_i \in \mathcal{Z}_{-j}\}$ 
8:   Coeff  $\leftarrow [1] \{P_0(\mathbf{x}) = 1\}$ 
9:   for  $\tilde{\mathbf{z}}_i \in \tilde{\mathcal{Z}}_{-j}$  do
10:    new_coeff  $\leftarrow []$ 
11:    Append  $-\tilde{\mathbf{z}}_i \cdot \text{Coeff}[0]$  to new_coeff {Term for  $\mathbf{x}^0$ }
12:    for  $k = 1$  to  $\text{len}(\text{Coeff}) - 1$  do
13:      term  $\leftarrow \text{Coeff}[k-1] - (\tilde{\mathbf{z}}_i \odot \text{Coeff}[k])$ 
14:      Append term to new_coeff
15:    end for
16:    Append Coeff[-1] to new_coeff
17:    Coeff  $\leftarrow$  new_coeff
18:  end for
19:  Extract  $e_q(\mathcal{Z}_{-j})$ :  $e_q(\mathcal{Z}_{-j}) \leftarrow (-1)^q \cdot \text{Coeff}[d-q-1] \cdot s^q$  for  $q = 0, \dots, d-1$ 
20:   $\psi_j \leftarrow \sum_{q=0}^{d-1} \mu(q) \odot e_q(\mathcal{Z}_{-j})$ 
21:   $\phi_j^{\mathbf{x}} \leftarrow \alpha^\top ((\mathbf{z}_j - \mathbf{1}) \odot \psi_j)$ 
22: end for
23: return  $(\phi_1^{\mathbf{x}}, \dots, \phi_d^{\mathbf{x}})$ 

```

Though being efficient, the recursive formulation suffers from numerical instability due to alternating sign cancellation, power operations, and error amplification through division. To develop a stable approach, we use the fact that symmetric elementary polynomials emerge naturally as coefficients of the characteristic polynomial (Egge, 2019):

$$P(\mathbf{x}) = \bigodot_{i=1}^d (\mathbf{x} - \mathbf{z}_i) = \sum_{q=0}^d (-1)^{d-q} e_{d-q}(\mathcal{Z}) \mathbf{x}^q.$$

Initialize $P_0(\mathbf{x}) = 1$. For each \mathbf{z}_i , the update rule is

$$P_i(\mathbf{x}) = P_{i-1}(\mathbf{x}) \odot (\mathbf{x} - \mathbf{z}_i),$$

and the coefficients evolve as

$$\text{Coeff}_m^{(i)} = \text{Coeff}_{m-1}^{(i-1)} - \left(\mathbf{z}_i \odot \text{Coeff}_m^{(i-1)} \right),$$

where $\text{Coeff}_m^{(i)}$ denotes the coefficient of \mathbf{x}^m after processing i elements. The elementary symmetric polynomials are then:

$$e_q(\mathcal{Z}) = (-1)^{d-q} \cdot \text{Coeff}_q^{(d)}.$$

Using this relationship, we can avoid the power and division operations in the standard recursive formulation and develop a more stable algorithm to compute ESPs. Aside from this, we scale kernel vectors $\mathbf{z}_i \leftarrow \mathbf{z}_i/s$ where $s = \max_{i,j} \mathbf{z}_{ij}$ to prevent overflow, with final correction $e_q \leftarrow e_q \cdot s^q$. Algorithm 2 shows the numerically stable algorithm for computing Shapley values for product-kernel learning models. This algorithm has the same complexity as Algorithm 1, but it avoids the power and division, which makes it more numerically stable (especially for high values of d).

Computing Shapley values for product-kernel models requires evaluating the marginal contribution of each feature across all possible coalitions, which naïvely involves computing d different sets of ESPs for the leave-one-out collections \mathcal{Z}_{-j} . It makes the overall time-complexity $O(d^3)$ for all d features. To overcome these limitations, we exploit a key structural insight: the ESPs of any subset \mathcal{Z}_{-j} appear as coefficients of the characteristic polynomial obtained by removing a single linear factor. Thus, instead of recomputing ESPs independently, we construct the global polynomial $P(t) = \prod_{i=1}^d (t - \tilde{z}_i)$ once, and then obtain each leave-one-out polynomial $Q^{(j)}(t) = P(t)/(t - \tilde{z}_j)$ using synthetic division. This procedure yields all ESPs $e_q(\mathcal{Z}_{-j})$ for every feature j at the cost of a single backward pass over the polynomial coefficients. The resulting algorithm is both numerically stable—avoiding powers, divisions, and alternating-sign recurrences—and achieves a quadratic $O(d^2)$ runtime, matching the efficiency of TreeSHAP while being fully compatible with product-kernel methods. In summary, the synthetic-division algorithm provides an exact, scalable, and stable method for Shapley value computation in multiplicative kernel models, enabling efficient instance-wise explainability without approximations.

The algorithm first constructs the global characteristic polynomial $P(t) = \prod_{j=1}^d (t - \tilde{z}_j)$, whose coefficient array is obtained by successively multiplying by linear factors via synthetic division. At step j , the current polynomial has degree $j - 1$, so the update touches $O(j)$ coefficient vectors. Summing over $j = 1, \dots, d$ yields $O(d^2)$ operations in the feature dimension. Next, for each feature j , we perform a single synthetic division of $P(t)$ by $(t - \tilde{z}_j)$, requiring $O(d)$ operations and producing the leave-one-out polynomial $Q^{(j)}(t)$. Repeating this for all d features again costs $O(d^2)$. Finally, extracting ESPs, forming the terms ψ_j , and computing each Shapley value $\phi_j^{\mathbf{x}}$ requires $O(d)$ per feature, again totalling $O(d^2)$. Overall, the full procedure computes all Shapley values in $O(d^2)$ time, while maintaining numerical stability and avoiding the cubic $O(d^3)$ cost of recomputing ESPs independently for each feature. Algorithm 3 summarizes the algorithm

F ADDITIVITY OF EXPLANATIONS FOR LEARNING MODELS, MMD AND HSIC

F.1 EXPLANATION ADDITIVITY FOR LEARNING MODELS

In Lemma 3, we established the additivity property of the explanation in product-kernel learning models. In particular, we demonstrated that

$$f(\mathbf{x}) - \boldsymbol{\alpha}^\top \mathbf{1} = \sum_j \phi_j^{\mathbf{x}},$$

where $\boldsymbol{\alpha}^\top \mathbf{1}$ represents the value of the null game according to the value function in equation 3. This result is useful since $\boldsymbol{\alpha}^\top \mathbf{1} = 0$ for several kernel methods such as support vector machines (we have the constraint $\sum \hat{\alpha}_j y_j = 0$ in the dual problem and $f(\mathbf{x}) = \sum_j \hat{\alpha}_j y_j k(\mathbf{x}, \mathbf{x}_j)$ with $\hat{\boldsymbol{\alpha}}$ be the solution to the dual problem) and support vector regression. However, in general, it distributes the value of $f(\mathbf{x})$ only after subtracting the null game’s value rather than allocating the full output $f(\mathbf{x})$ directly among the features. This is because $k_\emptyset = 1$ by definition, and this will lead to $v(\emptyset) = \boldsymbol{\alpha}^\top \mathbf{1}$. We now show that by redefining the value function in equation 3 so that $v(\emptyset) = 0$, the corresponding Shapley values will sum to $f(\mathbf{x})$, with each Shapley value augmented by $\frac{\boldsymbol{\alpha}^\top \mathbf{1}}{n}$.

Proposition 8. *Define a normalized value function by setting the kernel component for the empty set to zero, i.e., $k_\emptyset = 0$. Let $\hat{v}_{\mathbf{x}}$ be the corresponding value function and $\hat{\phi}_j$ the resulting Shapley values, which satisfy*

$$f(\mathbf{x}) = \sum_{j=1}^n \hat{\phi}_j.$$

Then, assuming an equal allocation of the baseline value, the following relationship holds for every feature j :

$$\hat{\phi}_j = \phi_j + \frac{\boldsymbol{\alpha}^\top \mathbf{1}}{n}.$$

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115

Algorithm 3 Shapley Values via Synthetic Division for Product-Kernel Methods

Require: Trained product-kernel model (SVM/SVR/GP), instance $\mathbf{x} \in \mathbb{R}^d$

Ensure: Shapley values $\phi_1^{\mathbf{x}}, \dots, \phi_d^{\mathbf{x}}$

- 1: Retrieve training data X , coefficients $\boldsymbol{\alpha}$, and kernel function k
 - 2: Compute feature-wise kernel vectors $\mathbf{z}_j = k_j(x_j, X_{:,j})$, for all $j = 1, \dots, d$
 - 3: Precompute Shapley weights $\mu(q) = \frac{q!(d-q-1)!}{d!}$, for $q = 0, \dots, d-1$
 - 4: **(Global scaling and characteristic polynomial)**
 - 5: $s \leftarrow \max_{i,j} z_{ij}$ {set $s = 1$ if all entries are zero}
 - 6: $\tilde{\mathbf{z}}_j \leftarrow \mathbf{z}_j / s$ for all $j = 1, \dots, d$
 - 7: Initialize coefficient list `Coeff` with `Coeff = [1]`
 - 8: **for** $j = 1$ **to** d **do**
 - 9: Let `new_coeff` be an empty list
 - 10: `new_coeff[0]` $\leftarrow -\tilde{\mathbf{z}}_j \odot \text{Coeff}[0]$
 - 11: **for** $m = 1$ **to** $\text{len}(\text{Coeff}) - 1$ **do**
 - 12: `new_coeff[m]` $\leftarrow \text{Coeff}[m-1] - \tilde{\mathbf{z}}_j \odot \text{Coeff}[m]$
 - 13: **end for**
 - 14: `new_coeff` $[\text{len}(\text{Coeff})] \leftarrow \text{Coeff}[\text{len}(\text{Coeff}) - 1]$
 - 15: `Coeff` \leftarrow `new_coeff`
 - 16: **end for**
 - 17: **(Synthetic division for leave-one-out ESPs)**
 - 18: **for** $j = 1$ **to** d **do**
 - 19: Let $r_j = \tilde{\mathbf{z}}_j$
 - 20: Initialize $Q^{(j)}[0..d-1]$
 - 21: $Q^{(j)}[d-1] \leftarrow \text{Coeff}[d]$
 - 22: **for** $m = d-1$ **to** 1 **do**
 - 23: $Q^{(j)}[m-1] \leftarrow \text{Coeff}[m] + r_j \odot Q^{(j)}[m]$
 - 24: **end for**
 - 25: $\{Q^{(j)}(t) = \prod_{i \neq j} (t - \tilde{\mathbf{z}}_i)\}$
 - 26: **Extract ESPs:**
 - 27: **for** $q = 0$ **to** $d-1$ **do**
 - 28: $e_q(\mathcal{Z}_{-j}) \leftarrow (-1)^q Q^{(j)}[d-1-q] \cdot s^q$
 - 29: **end for**
 - 30: $\psi_j \leftarrow \sum_{q=0}^{d-1} \mu(q) \odot e_q(\mathcal{Z}_{-j})$
 - 31: $\phi_j^{\mathbf{x}} \leftarrow \boldsymbol{\alpha}^\top ((\mathbf{z}_j - \mathbf{1}) \odot \psi_j)$
 - 32: **end for**
 - 33: **return** $(\phi_1^{\mathbf{x}}, \dots, \phi_d^{\mathbf{x}})$
-

Proof Using the efficiency axioms for $\hat{v}_{\mathbf{x}}$, the results will follow.

F.2 EXPLANATION ADDITIVITY FOR MMD

Lemma 9. For the MMD with product kernel, the sum of Shapley values satisfies:

$$\sum_{j=1}^d \phi_j^{\text{MMD}} = \widehat{\text{MMD}}^2(\mathbb{P}, \mathbb{Q}).$$

Proof By the efficiency property of Shapley values (Shapley, 1953):

$$\sum_{j=1}^d \phi_j^{\text{MMD}} = v_{\text{MMD}}(\mathcal{D}) - v_{\text{MMD}}(\emptyset)$$

For the product-kernel decomposition:

$$v_{\text{MMD}}(\mathcal{D}) = \frac{1}{n(n-1)} \sum_{i \neq j} k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) + \frac{1}{m(m-1)} \sum_{i \neq j} k(\mathbf{y}^{(i)}, \mathbf{y}^{(j)}) - \frac{2}{nm} \sum_{i,j} k(\mathbf{x}^{(i)}, \mathbf{y}^{(j)}) = \widehat{\text{MMD}}^2(\mathbb{P}, \mathbb{Q})$$

1134 For the empty coalition (all features removed):

$$1135 v_{\text{MMD}}(\emptyset) = \frac{1}{n(n-1)} \sum_{i \neq j} 1 + \frac{1}{m(m-1)} \sum_{i \neq j} 1 - \frac{2}{nm} \sum_{i,j} 1 = \frac{n(n-1)}{n(n-1)} + \frac{m(m-1)}{m(m-1)} - \frac{2nm}{nm} = 0$$

1138 Thus, $\sum_{j=1}^d \phi_j^{\text{MMD}} = \widehat{\text{MMD}}^2(\mathbb{P}, \mathbb{Q}) - 0 = \widehat{\text{MMD}}^2(\mathbb{P}, \mathbb{Q})$. \square

1142 F.3 EXPLANATION ADDITIVITY FOR HSIC

1143 **Lemma 10.** *For the HSIC dependence measure with a product kernel, the sum of Shapley values satisfies:*

$$1146 \sum_{j=1}^d \phi_j^{\text{HSIC}} = \widehat{\text{HSIC}}(X, y).$$

1149 **Proof** By the efficiency property of Shapley values (Shapley, 1953), the sum of Shapley values equals the difference between the value of the grand coalition and the empty coalition:

$$1152 \sum_{j=1}^d \phi_j^{\text{HSIC}} = v_{\text{HSIC}}(\mathcal{D}) - v_{\text{HSIC}}(\emptyset).$$

1155 From Definition 6, we have:

$$1157 v_{\text{HSIC}}(\mathcal{D}) = \frac{1}{(n-1)^2} \text{tr} \left(\mathbf{HLH} \bigodot_{j \in \mathcal{D}} \mathbf{K}_j \right) = \widehat{\text{HSIC}}(X, y),$$

$$1161 v_{\text{HSIC}}(\emptyset) = \frac{1}{(n-1)^2} \text{tr}(\mathbf{HLH}\mathbf{1}\mathbf{1}^\top).$$

1163 One can simply realize that $\mathbf{HLH}\mathbf{1}\mathbf{1}^\top = 0$ as the sum of rows and columns in H is zero, and substituting these into the efficiency property completes the proof. \square

1166 G HSIC ATTRIBUTION WITH TWO MULTIVARIATE VARIABLES

1169 We studied the Shapley value computation for HSIC when it measures dependence between a multivariate variable \mathbf{x} and a univariate target y . We now extend this framework to two multivariate variables $X \in \mathbb{R}^d$ and $Z \in \mathbb{R}^d$ with product kernels function k and l . Given a sample $\{(\mathbf{x}^{(i)}, \mathbf{z}^{(i)})\}_{i=1}^n \sim \mathbb{P}(X, Z)$, $\text{HSIC}(X, Z)$ can be estimated as:

$$1173 \widehat{\text{HSIC}}(X, Z) = \frac{1}{(n-1)^2} \text{tr}(\mathbf{HKHL}),$$

1176 where $\mathbf{K} \in \mathbb{R}^{n \times n}$ is the kernel matrix computed using the kernel k , i.e., $\mathbf{K}_{ij} = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$, $\mathbf{L} \in \mathbb{R}^{n \times n}$ is the kernel matrix computed using the kernel l , i.e., $\mathbf{L}_{ij} = l(\mathbf{z}^{(i)}, \mathbf{z}^{(j)})$, and $\mathbf{H} = I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$ is the centering matrix ensuring zero mean in the feature space.

1179 To address the attribution to variables in both X and Z , we establish two cooperative games to attribute dependence contributions: For product kernels k and l , we define two value functions for the two games:

- 1183 (i) We first assume that \mathbf{L} is fixed, and try to attribute the total HSIC to the variables in X . Since k is a product kernel, we can define a value function similar to Definition 6. Let \mathcal{D}_X be the set of variables of X , we therefore define the value function for attributing the total HSIC to X variables as:

$$1184 v_{\text{HSIC}_X}(\mathcal{S}) = \frac{1}{(n-1)^2} \text{tr}(\mathbf{HK}_\mathcal{S}\mathbf{HL}), \quad \forall \mathcal{S} \subseteq \mathcal{D}_X. \quad (7)$$

(ii) By the same token, we take \mathbf{K} fixed and try to attribute the total HSIC to the variables in Z . Since l is a product kernel, we define the value function as:

$$v_{\text{HSIC}_Z}(\mathcal{S}) = \frac{1}{(n-1)^2} \text{tr}(\mathbf{H}\mathbf{K}\mathbf{H}\mathbf{L}_{\mathcal{S}}), \quad \forall \mathcal{S} \subseteq \mathcal{D}_Z, \quad (8)$$

where \mathcal{D}_Z is the set of variables of Z .

Building on the two value functions, we can compute Shapley values for variables in X and Z separately. We denote the Shapley value of variable j for X and Z as $\phi_j^{\text{HSIC}_X}$ and $\phi_j^{\text{HSIC}_Z}$, respectively. These values are interpreted as:

- $\phi_j^{\text{HSIC}_X}$ quantifies the contribution of the j^{th} variable in X to the total dependence between X and Z ;
- $\phi_j^{\text{HSIC}_Z}$ quantifies the contribution of the j^{th} variable in Z to the total dependence between X and Z .

H EXPERIMENTS

H.1 EXPERIMENTAL SETUP

SVM Optimization Using Optuna (Akiba et al., 2019) When using SVM in our experiments, we optimized the Support Vector Machine (SVM) with a Radial Basis Function (RBF) kernel using Optuna (Akiba et al., 2019), a robust hyperparameter optimization framework. The target type (either 'regression' or 'classification') was determined to guide the selection of the appropriate SVM model (SVR for regression and SVC for classification). The hyperparameters C and γ , critical for the RBF kernel's performance, were optimized within an extensive range using a log-uniform distribution. Specifically, we defined the hyperparameters: C between 10^{-5} and 10^5 and γ between 10^{-5} and 10^3 , and utilized 5-fold cross-validation to ensure reliable evaluation. The optimization process aimed to minimize the mean squared error for regression tasks and maximize accuracy for classification tasks. After conducting the specified number of trials ($n=100$), the best hyperparameters were used to train a final SVM model on the entire dataset, yielding both the optimal model configuration and the best cross-validation score achieved during the optimization process.

Training Gaussian Process (GP) Using K-Fold Cross-Validation When using GP, we trained a model using k-fold cross-validation to ensure robust evaluation and generalization performance. We defined the GP kernel as $C(1.0, (1e-4, 1e1)) * \text{RBF}(1.0, (1e-4, 10))$, suitable for both regression and classification tasks. For classification problems, `GaussianProcessClassifier` was utilized, while `GaussianProcessRegressor` was used for regression tasks. We employed `K-Fold` with $K = 5$ for cross-validation to evaluate the model's performance across different folds. All the hyperparameters, including the kernel width of the RBF kernel, are determined in the training process using optimization. During the cross-validation process, the model was trained on each fold, and predictions were made on the validation fold. Performance metrics were chosen based on the problem type: accuracy for classification models and mean absolute percentage error (MAPE) for regression models. The scores from each fold were aggregated to compute the average and standard deviation of the scores.

H.2 EXECUTION TIME

To assess the computational efficiency of our recursive algorithm, we conducted a simulation study using a randomly generated kernel function of the form $\alpha^T k(\mathbf{x}, \mathbf{X})$ with 1000 samples. We computed Shapley values under both the brute-force enumeration and our recursive method by varying the number of features, as plotted in Figure 1. With a runtime budget of 300 seconds, the brute-force computation was feasible only for up to 30 features, already requiring more than 200 seconds at this scale, while exceeding the budget beyond that. In contrast, our recursive algorithm consistently completed the same computations in a fraction of the time, from milliseconds at lower dimensions to under 100 seconds, even with 500 features. These results clearly highlight the substantial efficiency gains and scalability of our method for high-dimensional settings.

1242 H.3 MMD EXPERIMENTS 1243

1244 In addition to the synthetic experiments for MMD, we first provide another experiment for the
1245 cases when there is no distribution discrepancy. To that end, X and Z are sampled from the same
1246 multivariate normal distribution across all 20 variables. The MMD is near zero, indicating the
1247 distributions are equivalent. Shapley values are computed and replicated 1000 times, with histograms
1248 plotted for each variable in Figure 6. The near-identical distributions of Shapley values across all
1249 variables reflect the uniform contribution of these variables to the MMD close to zero, consistent with
1250 the absence of any significant difference between the distributions.

1251 We extend our analysis to the UCI Diabetes dataset, consisting of 442 samples and 10 baseline
1252 variables, including age, sex, body mass index (BMI), average blood pressure, and six blood serum
1253 measurements (shown by $s1$ to $s6$ features). The dataset is split into male and female subsets using
1254 the second variable (sex), which is excluded from the analysis, leaving nine variables for comparison.

1255 Using MMD, we calculate the dissimilarity between male and female groups and then compute
1256 Shapley values to attribute variable contributions to the MMD. Figure 7 displays the Shapley values
1257 for the nine variables in the Diabetes dataset. The results show that $s3$ and $s4$ contribute most
1258 significantly to the MMD, followed by bp , $s6$, age , $s5$, and $s2$. In contrast, bmi and $s1$ reduce the
1259 MMD, indicating their alignment across the two groups.

1260 To validate these results, we analyze the marginal distributions of variables for males and females,
1261 as shown in Figure 8. The analysis confirms that variables with distinct marginal distributions
1262 between males and females (e.g., $s3$ and $s4$) have high positive Shapley values, reflecting their role
1263 in increasing the MMD. Conversely, variables with similar distributions (e.g., $s1$) exhibit negative
1264 Shapley values, highlighting their role in reducing the MMD.

1266 I HOW DOES PKE X-SHAPLEY DIFFER FROM MOHAMMADI ET AL. (2025B)? 1267

1268 In Mohammadi et al. (2025b), the authors propose an exact Shapley computation method for stochastic
1269 attribution in FANOVA GP models. Similar to PKE X-Shapley, their approach leverages Newton’s
1270 identities and symmetric polynomial representations to obtain closed-form attributions. However, the
1271 two methods differ in both scope and capability.

1272 PKE X-Shapley is designed specifically for kernel methods with product kernels. It presents a new
1273 value function, i.e., a functional baseline value function, to compute Shapley values. While this
1274 formulation enables exact computation of the mean Shapley values, it does not extend naturally to
1275 higher-order moments. In particular, when applied to Gaussian processes, computing the variance
1276 or covariance structure of Shapley values is nontrivial, since orthogonality does not generally hold
1277 for product kernels. As a result, our method cannot be directly extended to explain GP models in a
1278 stochastic way.

1279 By contrast, Mohammadi et al. (2025b) focus on FANOVA GPs, whose kernels admit a functional
1280 ANOVA decomposition. Compared to product kernels, this class of kernels is more restrictive and
1281 often harder to train in practice. Nonetheless, the additional structure provides a key advantage:
1282 it enables exact polynomial-time computation not only of the mean but also of the variance and
1283 covariance of stochastic Shapley values—something that product kernels cannot achieve. This
1284 distinction highlights the central difference between their method and ours.

1286 J HOW DOES PKE X-SHAPLEY DIFFER FROM RKHS-SHAP (CHAU ET AL., 1287 2022)? 1288

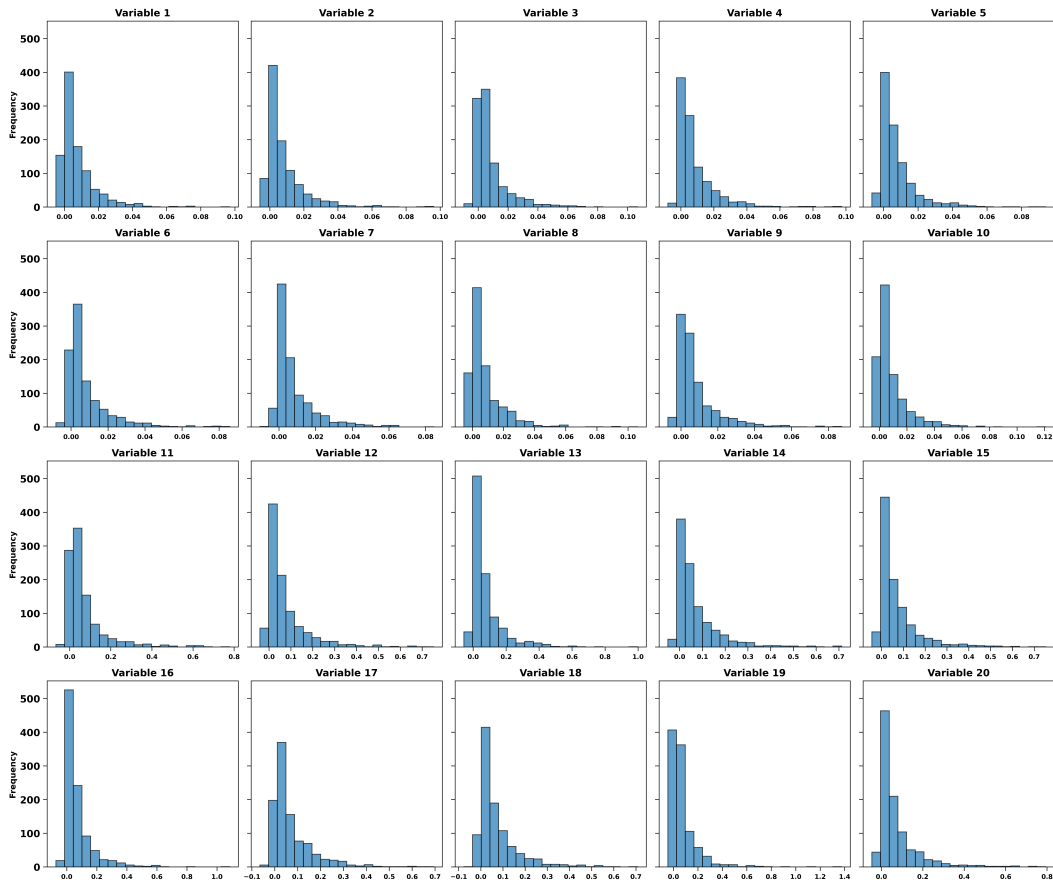
1289 RKHS-SHAP is the first kernel method-specific SHAP-based algorithm. While the author still
1290 employs the conditional expectation value function $\tilde{\nu}_{\mathbf{x}}(S) = \mathbb{E}[f(X) \mid X_S = \mathbf{x}_S]$, they used the fact
1291 that for function f in the RKHS \mathcal{H}_k , it is possible to estimate $\tilde{\nu}_{\mathbf{x}}(S)$ non-parametrically utilizing a tool
1292 known as conditional kernel mean embedding. Specifically, this leads to the following expression:

$$1293 \tilde{\nu}_{\mathbf{x}}(S) = \boldsymbol{\alpha}^\top \mathbf{K}(\mathbf{K}_S + n\lambda I)^{-1} k_S(\mathbf{X}_S, \mathbf{x}_S)$$

1294 which of course, is different from our functional baseline value function

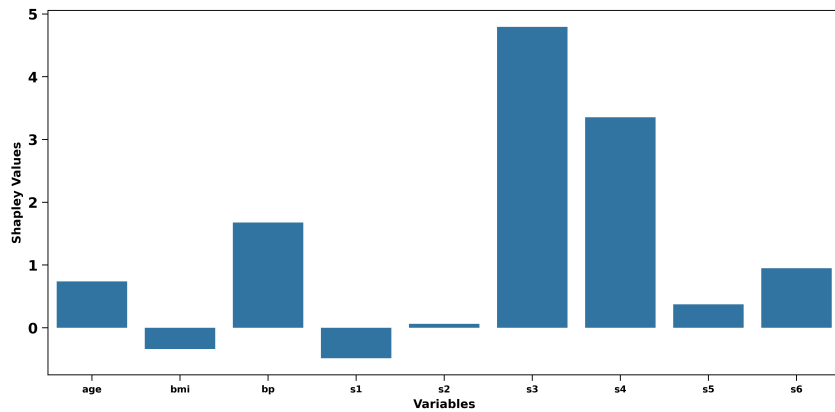
$$1295 \nu_{\mathbf{x}}(S) = \boldsymbol{\alpha}^\top k_S(\mathbf{X}_S, \mathbf{x}_S).$$

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325



1326 Figure 6: Shapley values for the synthetic data sets with equal distributions. All variables contribute
 1327 equally to the near-zero MMD.
 1328

1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349



1346 Figure 7: Shapley values explaining MMD between male and female subsets in the UCI Diabetes
 1347 data set.
 1348

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

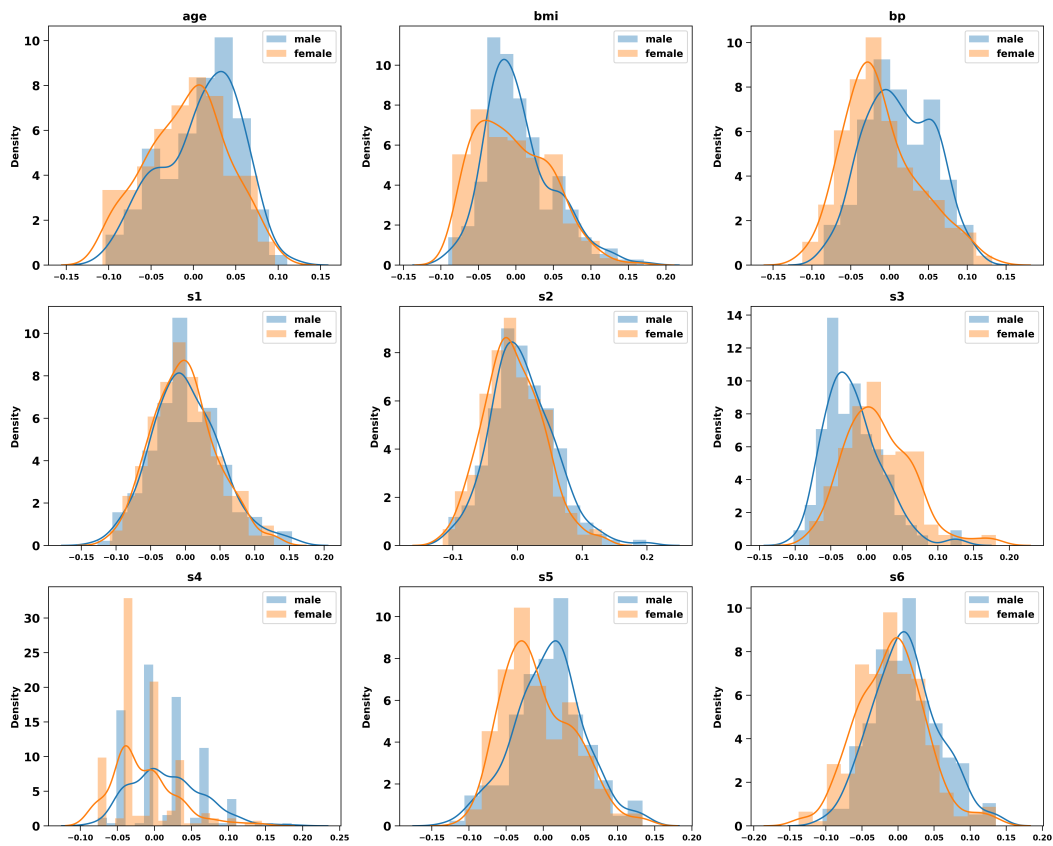


Figure 8: Marginal distributions of variables for male and female subsets in the UCI Diabetes dataset.

1404 Although looking different in their expression, there is an interesting mathematical connection that
1405 is not apparent at first glance. First of all, we need to understand that conditional expectations
1406 can be interpreted as orthogonal projections to the space of nice-behaving functions defined on the
1407 conditioning variable. Specifically, let $\mathcal{L}^2(\mathcal{X})$ be the space of square-integrable functions on X , then
1408 for a random variable Y , the conditional expectation $\mathbb{E}[Y | X]$ is the unique element in $\mathcal{L}^2(\mathcal{X})$ that is
1409 closest to Y in mean-square error:

$$1410 \mathbb{E}[Y | X] = \arg \min_{Z \in \mathcal{L}^2(\mathcal{X})} \mathbb{E}[(Y - Z)^2].$$

1412
1413 As a result, this interpretation of conditional expectations, and thus of the conditional-expectation-
1414 based value function $\tilde{v}_{\mathbf{x}}$, allows us to directly connect it to our functional baseline value function.
1415 As established in Proposition ?? and Appendix ??, the latter can also be viewed as an orthogonal
1416 projection, albeit onto a different function space.

1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457