

# Beyond Additivity: Sparse Isotonic Shapley Regression toward Nonlinear Explainability

Jialai She

## Abstract

Shapley values, a gold standard for feature attribution in Explainable AI, face two primary challenges. First, the canonical Shapley framework assumes that the worth function is additive, yet real-world payoff constructions—driven by non-Gaussian distributions, heavy tails, feature dependence, or domain-specific loss scales—often violate this assumption, leading to distorted attributions. Secondly, achieving sparse explanations in high-dimensional settings by computing dense Shapley values and then applying ad hoc thresholding is prohibitively costly and risks inconsistency. We introduce Sparse Isotonic Shapley Regression (SISR), a unified nonlinear explanation framework. SISR simultaneously learns a monotonic transformation to restore additivity—obviating the need for a closed-form specification—and enforces an L0 sparsity constraint on the Shapley vector, enhancing computational efficiency in large feature spaces. Its optimization algorithm leverages Pool-Adjacent-Violators for efficient isotonic regression and normalized hard-thresholding for support selection, yielding ease in implementation and global convergence guarantees. Analysis shows that SISR recovers the true transformation in a wide range of scenarios and achieves strong support recovery even in high noise. Moreover, we are the first to demonstrate that irrelevant features and inter-feature dependencies can induce a true payoff transformation that deviates substantially from linearity. Extensive experiments in regression, logistic regression, and tree ensembles demonstrate that SISR stabilizes attributions across payoff schemes, correctly filters irrelevant features; in contrast, standard Shapley values suffer severe rank and sign distortions. By unifying nonlinear transformation estimation with sparsity pursuit, SISR advances the frontier of nonlinear

explainability, providing a theoretically grounded and practical attribution framework.

**Keywords:** Shapley value, machine learning explainability, isotonic regression, sparsity pursuit

## 1 Introduction and Motivation

Let  $F = \{1, 2, \dots, p\}$  denote the set of  $p$  features, and let  $\nu : 2^F \rightarrow \mathbb{R}$  be the characteristic function or payoff function, with  $\nu(A)$  representing the contribution or worth generated by a subset  $A \subseteq F$  of features working together (often referred to as a *coalition* in game theory). A central question in economics and cooperative game theory is how to fairly allocate the value of a coalition to its individual members. The *Shapley value* (Shapley, 1953), a concept from Nobel laureate Lloyd Shapley, offers a theoretically grounded solution by assigning payoffs according to each member’s average marginal contribution across all possible subsets.

In this paper, we denote the Shapley value for feature  $j$  by  $\beta_j$  for  $1 \leq j \leq p$ , quantifying the fair share or importance of feature  $j$ ; for a subset  $A \subseteq F$  we define  $\beta_A$  as the vector  $[\beta_j]_{j \in A} \in \mathbb{R}^{|A|}$ . For brevity, we also write  $\nu_A$  as shorthand for  $\nu(A)$  for any subset  $A \subseteq F$ . Shapley values establish a connection between the payoff function  $\nu(A)$  and the underlying model parameters  $\beta_A$ . To make this dependence explicit, we introduce a function  $V(\beta_1, \dots, \beta_p; A)$ , also denoted by  $V_A(\{\beta_j\}_{j \in A})$ , characterizing the deterministic, *noise-free* contribution associated with subset  $A$ :

$$\nu_A \sim V_A(\{\beta_j\}_{j \in A}), \quad (1)$$

where  $\sim$  denotes approximate equality up to noise, a convention adopted throughout the paper.

In recent years, Shapley values have attracted substantial attention in machine learning, particularly in the field of *Explainable AI* (**XAI**) (Ancona et al., 2019). While assessing variable importance in simple regression is straightforward using traditional tools like  $T$ -tests and  $p$ -values, this task becomes a formidable challenge for the complex, “black-box” models now widely used to analyze sequential data in economics and finance. For sophisticated models—ranging from *tree-based ensembles* like random forests and boosted trees to *deep neural networks* like Long Short-Term Memory (LSTM)

networks—standard inference methods are no longer applicable, making interpretation notoriously difficult.

Shapley values provide a model-agnostic framework for quantifying feature importance. This is applied in two main ways: local explanations, which explain a single prediction (e.g., SHAP by [Lundberg and Lee \(2017\)](#)), and global explanations, which explain the model’s overall behavior (e.g., SAGE by [Covert et al. \(2020\)](#)). Our work focuses on the global setting, where the goal is to find a single, interpretable set of importance values for the entire model. Specifically, for a prediction model  $f(x)$ , where  $x \in \mathbb{R}^p$ , researchers first design a payoff function  $\nu_A$  over subsets  $A \subseteq F$  to quantify the model’s global performance when using only the features in  $A$ . This reframes the explanation task as a “**credit allocation**” problem, enabling the use of Shapley values to quantify feature importance, and perhaps more importantly, to construct interpretable *surrogate models* based on restricted feature sets. The resulting additive structure of individual feature contributions is the *very* property appealing for interpretability, as it provides an intuitive, linear explanation of an intricate model’s behavior, a primary goal in XAI.

However, standard Shapley-based methods also face several limitations that restrict their practical utility in complex modeling scenarios.

**(i) Moving beyond additive frameworks:** Given a prediction model, various methods have been proposed to construct the payoff function  $\nu_A$  for Shapley-value analysis. (a) A fundamental approach involves retraining the model on every subset of features  $A$  and defining  $\nu_A$  based on the reduction in statistical accuracy (such as  $R^2$  in regression) ([Lipovetsky and Conklin, 2001](#)). However, this exhaustive procedure may be computationally prohibitive for modern AI models due to its exponential cost. (b) To circumvent retraining, SAGE ([Covert et al., 2020](#)) provides an efficient alternative: it keeps the trained model fixed and quantifies the expected loss increase when certain features are made unavailable. The approach marginalizes over missing inputs—conditionally in theory and, for scalability, interventionally in practice. (c) In contrast, SHAP and TreeSHAP ([Lundberg and Lee, 2017](#); [Lundberg et al., 2020](#)) define local payoffs for each instance by marginalizing absent features and then derive global importance by aggregating the resulting local Shapley attributions. (d) Other global variants, such as Sobol-Shapley indices ([Owen, 2014](#)) and derivative-based formulations ([Duan and Okten, 2025](#)), approximate risk or variance decomposition

under specific assumptions (feature independence, distributional priors, or model smoothness), often motivated by numerical sensitivity analysis rather than prediction risk. Once  $\nu_A$  is constructed, researchers often mechanically apply the Shapley formula to compute feature attributions.

However, the theoretical justification for Shapley values relies on several foundational “axioms”—efficiency, symmetry, linearity, and nullity (Shapley, 1953)—which are not easily testable and are *rarely* validated in practice. In particular, Shapley’s framework implicitly assumes an **additive structure** (Lundberg and Lee, 2017):

$$\nu_A \sim \sum_{j \in A} \beta_j \quad \text{or} \quad V_A(\{\beta_j\}_{j \in A}) = \sum_{j \in A} \beta_j. \quad (2)$$

But the so-called additive feature attribution is not guaranteed to hold in real-world constructions of coalition values. For example, we can reformulate the abstract Shapley axioms and principles into a multivariate Gaussian assumption (cf. Section 2), but many of the constructions mentioned previously are prone to violating this assumption due to **non-Gaussian** characteristics such as bounded ranges, heavy tails, and skewness. In particular, Fryer et al. (2021) recently proposed a realistic “taxicab” payoff defined by a *winner-takes-all* dynamic that is in stark contrast to (2):

$$V_A(\{\beta_j\}_{j \in A}) = \max_{j \in A} \beta_j, \quad \forall A \subseteq 2^F. \quad (3)$$

Such nonlinear relations are prevalent in applications but fundamentally violate the additive model underpinning standard Shapley value estimation.

**(ii) Embedding sparsity into value attribution:** In many real-world applications with a large number of features, a substantial proportion contribute only negligibly—or are effectively irrelevant—to the overall outcome, making them unnecessary to explain in practice (Strumbelj and Kononenko, 2014; Covert et al., 2021). Exploiting the structural parsimony can enhance both statistical accuracy and interpretability of Shapley values. In implementation, leveraging sparsity helps to reduce iteration complexity, thanks to a substantially smaller effective model size, along with mitigating communication costs and storage requirements in high-dimensional settings.

However, existing approaches adopt a **greedy** strategy to achieve sparsity and have significant drawbacks. Many methods first compute dense Shapley values for the *full* model and then apply post-hoc ranking or thresholding

(Cohen et al., 2007; Jothi et al., 2021; Fryer et al., 2021; Au et al., 2022). For large  $p$ , such multi-step procedures are not only inefficient but may also fail to provide faithful explanations or meaningful selection (see, e.g., Covert et al. (2021); Slack et al. (2020); Ma and Tourani (2020)). An alternative class of older, less efficient approaches resorts to an  $\ell_1$ -penalty (Lundberg and Lee, 2017; Ribeiro et al., 2016), but requires cumbersome parameter tuning and induces unwanted shrinkage on the attribution values, which can distort their magnitude. This often necessitates a multi-step re-fitting procedure to correct for the shrinkage, undermining the goal of a unified estimation. More fundamentally, this entire approach relies on the  $\ell_1$ -norm’s ability to select the correct features, an inherent drawback as it often fails to recover the true support, especially in the presence of correlated features (Zhang, 2010).

To the best of our knowledge, *no* widely adopted framework integrates direct, shrinkage-free sparsity control as an intrinsic property into Shapley-value estimation, let alone in the context of an unknown transformation. These challenges underscore the need for a unified approach that *simultaneously* enforces sparsity and ensures coherent Shapley-based attributions.

This paper aims to develop a novel nonlinear explanation framework for applying the Shapley mechanism in a way that simultaneously aligns individual feature contributions with appropriately transformed worths across all subsets, and promotes sparsity by eliminating irrelevant features to enhance both computational efficiency and statistical accuracy. The contributions of our work are as follows:

- Our research is the first to demonstrate that common factors such as the presence of irrelevant features and inter-feature dependencies can induce a payoff transformation that deviates substantially from linearity, even when using standard payoff constructions (e.g.,  $R^2$ -based worths). This finding underscores the need for nonlinear explainability frameworks.
- We propose Sparse Isotonic Shapley Regression (**SISR**), the first framework to *jointly* address payoff non-additivity and attribution sparsity. By learning a monotonic transformation and enforcing an  $\ell_0$  constraint simultaneously, our integrated approach overcomes the limitations of ad-hoc methods.
- SISR learns the transformation of payoffs without requiring a prede-

finely analytical form. This is achieved through efficiently leveraging the Pool-Adjacent-Violators algorithm, allowing the model to adapt to diverse real-world payoff structures.

- The optimization algorithm developed for SISR features simple, closed-form updates and is accompanied by global convergence guarantees. The incorporation of sparsity improves computational efficiency
- Through extensive experiments across various datasets and payoff schemes, we show that SISR significantly stabilizes feature attributions and correctly identifies relevant features, mitigating the severe rank and sign distortions often observed with standard Shapley value applications.

Our goal with SISR is not to abandon the interpretability of additivity, but rather to restore it. While other valuable frameworks explicitly model higher-order feature interactions (cf. Section 5), this paper proposes a distinct alternative that seeks to *learn* a principled monotonic transformation that maps the payoff function *back* to a domain where a simple, additive main-effect structure holds. This preserves the interpretability of the original Shapley framework while robustly handling payoff constructions driven by non-Gaussian distributions and domain-specific loss scales that violate its core assumptions.

The rest of the paper is organized as follows. Section 2 proposes a novel Sparse Isotonic Shapley Regression model to address challenges related to domain adaptation and high dimensionality. In Section 3, an optimization-based algorithm is developed to address the functional challenge and the nonsmooth sparsification, with established theoretical guarantees. Section 4 provides valuable data-driven insights drawn from experiments in various scenarios. Extensions are given in Section 5. We conclude in Section 6.

## 2 Proposed Method

The Shapley axioms and principles have been interpreted by economists in various ways (Algaba et al., 2019). Here, we recast the Shapley framework as a statistical assumption on the data-generating process. To begin, let us revisit a motivating *weighted least squares* formulation of Shapley value estimation as derived in Lundberg and Lee (2017), echoing earlier developments

in econometrics ([Charnes et al., 1988](#)):

$$\begin{aligned} \min_{\beta \in \mathbb{R}^p, c \in \mathbb{R}} \quad & \sum_{A \subseteq 2^F, A \neq \emptyset, A \neq F} w_{\text{SH}}(A) \left( \nu_A - \sum_{j \in A} \beta_j - c \right)^2 \\ \text{subject to} \quad & c = \nu_{\emptyset}, \quad c + \sum_{j=1}^p \beta_j = \nu_F, \end{aligned} \tag{4}$$

where the Shapley weights are given by

$$w_{\text{SH}}(A) = \frac{p-1}{\binom{p}{|A|} |A| (p-|A|)}. \tag{5}$$

Perhaps surprisingly, it can be shown that the optimal solution  $\hat{\beta}$  to (4) recovers the exact Shapley values ([Lundberg and Lee, 2017](#)), which are traditionally derived based on the concept of *marginal contributions* across all possible feature coalitions.

If we define

$$w_{\text{SH}}(\emptyset) = +\infty, \quad w_{\text{SH}}(F) = +\infty \tag{6}$$

as an extension of (5), then (4) can be written as

$$\min_{\beta, c} \sum_{A \subseteq 2^F} w_{\text{SH}}(A) (\nu_A - \sum_{j \in A} \beta_j - c)^2$$

where  $A$  can take any subset of the power set  $2^F$ . Note that when  $A = \emptyset$ ,  $\sum_{j \in A} \beta_j = 0$  by convention and  $\hat{c} = \nu_{\emptyset}$ . It is thus convenient to define the *baseline-adjusted* coalition values:

$$\nu_A^c = \nu_A - \nu_{\emptyset}, \quad \forall A \subseteq 2^F, \tag{7}$$

which saves one parameter in the subsequent optimization:

$$\min_{\beta} \sum_{A \subseteq 2^F} w_{\text{SH}}(A) \left( \nu_A^c - \sum_{j \in A} \beta_j \right)^2.$$

For notational simplicity, we will write ' $\nu_A^c$ ' as just ' $\nu_A$ ', assuming that all  $\nu$  values have been properly shifted according to (7) unless otherwise specified.

It is helpful to reinterpret the weighted least squares formulation of Shapley values as a probabilistic model:

$$\begin{aligned}\nu_A &\sim \mathcal{N}(\mu_A, \sigma_A^2) \\ \mu_A &= \sum_{j \in A} \beta_j^*, \\ \sigma_A^2 &\propto \binom{p}{|A|} |A| (p - |A|) \left( \propto \frac{1}{w_{\text{SH}}(A)} \right),\end{aligned}\tag{8}$$

and all  $\nu_A$ 's are independent. Here,  $\beta_j^*$  denotes the true Shapley value for the  $j$ th feature.

By reformulating the Shapley axioms and principles into assumption (8), we gain insight into why numerous payoff functions may not meet the model criteria. Indeed, due to issues such as range constraints, skewness, heavy tails, and heterogeneity, it is natural to question the appropriateness of the multivariate Gaussianity across different definitions of coalition values.

In our view, one viable solution is to apply a transformation that promote Gaussianity. Let's consider an alternative Shapley value model in a *transformed domain*:

$$T(\nu_A) \sim \mathcal{N}\left(\sum_{j \in A} T(\beta_j^*), \sigma_A^2\right),\tag{9}$$

where  $T(\cdot)$  is an unknown transformation. Under this model,

$$\mathbb{E}[T(\nu_A)] = \sum_{j \in A} T(\beta_j^*),\tag{10}$$

which defines a “ $T$ -additive” framework for nonlinear settings. To model this structure, we propose a new Shapley framework termed Functional Shapley Regression, which jointly estimates  $\beta$  and  $T(\cdot)$  by solving

$$\min_{\beta, T(\cdot)} \sum_{A \subseteq 2^F} w_{\text{SH}}(A) \left\{ T(\nu_A) - \sum_{j \in A} T(\beta_j) \right\}^2 \text{ subject to } \beta \in \mathcal{C}, T(\cdot) \in \mathcal{T},\tag{11}$$

where the objective minimizes the Shapley-weighted sum of squared differences between the transformed coalition values  $T(\nu_A)$  and the transformed linear sum  $\sum_{j \in A} T(\beta_j)$  over all subsets  $A \subseteq 2^F$ . Here, we use  $\mathcal{C} \subseteq \mathbb{R}^p$



to denote the constraint set for  $\beta$ , and  $\mathcal{T}$  to denote the class of admissible transformation functions. By the notational convention for  $A = \emptyset$ , (11) automatically enforces

$$T(0) = 0,$$

corresponding to  $T(\nu_\emptyset) = 0$  (recall all  $\nu_A$  have been centered).

The remark below illustrates that our framework, in contrast to the common additive main-effect model (see, e.g., [Lundberg and Lee \(2017\)](#)), accommodates a broader range of multivariate payoff structures. By learning a transformation to restore the Shapley framework’s underlying statistical assumptions (Gaussianity), the resulting  $T^{-1}$ -sum- $T$  structure in (12) enables *nonlinear* explainability.

**Remark 1 (Univariate  $T$ -Mappings for Multivariate Structure).** *Introducing a univariate transformation  $T(\cdot)$  enables a remarkably rich class of models capable of capturing complex multivariate relationships between  $\nu_A$  and  $\{\beta_j : j \in A\}$ , well beyond the standard additive form.*

*Specifically, under the  $T$ -transformed model (9), assuming the existence of the inverse transformation  $T^{-1}$  and using the notation  $V_A$  (cf. (1)) we have*

$$V_A(\{\beta_j\}_{j \in A}) = T^{-1}\left(\sum_{j \in A} T(\beta_j)\right) \quad \text{or} \quad \nu_A \sim T^{-1}\left(\sum_{j \in A} T(\beta_j)\right). \quad (12)$$

*If  $T$  is a nondegenerate linear map like the identity map, the model reduces to the conventional additive Shapley game, where the coalition value  $\nu_A$  in (12) is essentially a simple sum of individual contributions. However, a general transformation lends the  $T^{-1}$ -sum- $T$  multivariate structure of (12) the flexibility to model a broad range of application domains. By learning the transformation from the data, our framework acts as a robust generalization of the conventional model, rather than imposing a forced, arbitrary transformation.*

*For instance, consider a monomial transformation  $T(x) = |x|^d$  for some  $d > 0$ , which, under the mild assumption of non-negative payoffs, induces a multivariate “ $d$ -norm” relationship:*

$$V_A(\{\beta_j\}_{j \in A}) = \left(\sum_{j \in A} |\beta_j|^d\right)^{1/d} = \|\beta_A\|_d.$$

*Varying the degree  $d$  recovers a spectrum of geometric structures, e.g.,*

- (i)  $d = 1$ : the  $\ell_1$ -norm polytope,  $V_A(\{\beta_j\}_{j \in A}) = \sum_{j \in A} |\beta_j|$ ;

(ii)  $d = 2$ : the  $\ell_2$ -norm ball,  $V_A(\{\beta_j\}_{j \in A}) = (\sum_{j \in A} \beta_j^2)^{1/2}$ ;

(iii)  $d \rightarrow \infty$ : the  $\ell_\infty$ -norm cube,  $V_A(\{\beta_j\}_{j \in A}) = \max_{j \in A} |\beta_j|$ .

In particular, under nonnegativity constraints ( $\nu_A \geq 0$ ,  $\beta_j \geq 0$ ), the  $\ell_\infty$  case corresponds to the winner-takes-all mechanism as first noted in [Fryer et al. \(2021\)](#), where the coalition value is dominated by the largest individual contribution (practically, monomial transformations with large degrees  $d$  can closely approximate such behavior). This is motivating, as the examples here are highly nonlinear and incompatible with a linear Shapley game. Yet, with a univariate transformation, they can be incorporated into the  $T$ -Shapley framework. Additional examples are the exponential form  $T(x) = \exp(x) - 1$  and the odds form  $T(x) = \Phi(x)/(1 - \Phi(x))$  with  $\Phi$  a distribution function of a continuous random variable.

In sum, an appropriately chosen  $T(\cdot)$  establishes a versatile nonlinear modeling mechanism that enhances the additive expressiveness of Shapley values for XAI. Another advantage of the proposed approach is that it **bypasses** the need for a predefined analytic transformation, instead learning it directly from the data (cf. [Section 3](#)). This novel capability of “learning to be additive” to recover a linear, main-effect Shapley structure is a key contribution of the framework.

In this paper, we focus a specific instance of [\(11\)](#), referred to as the “**S**parse **I**sotonic **S**hapley **R**egression” (**SISR**):

$$\begin{aligned} \min_{\beta, T(\cdot)} \sum_{A \subseteq 2^F} w_{\text{SH}}(A) \left\{ T(\nu_A) - \sum_{j \in A} T(\beta_j) \right\}^2 \\ \text{subject to } \|\beta\|_0 \leq s, T \in \mathcal{M}, \sum_{j=1}^p (T(\beta_j))^2 = 1, \end{aligned} \tag{13}$$

where  $\mathcal{M}$  denotes the class of strictly increasing functions and  $1 \leq s \leq p$  specifies the user-defined upper bound on the true model sparsity. Notably, this objective is defined on the  $T(\nu_A)$  scale, as this is the domain where the Gaussian error assumption holds (cf. [\(9\)](#)); applying a quadratic loss to the original  $\nu_A$  scale would be statistically inconsistent. [\(13\)](#) incorporates three critical modeling considerations.

**Monotonicity.** We impose a monotonicity constraint on  $T(\cdot)$  to preserve the relative ordering of feature importance values:

$$\beta_i \geq \beta_j \quad \Rightarrow \quad T(\beta_i) \geq T(\beta_j).$$

This ensures that the learned transformation respects the relative contribution levels of individual features. The structure closely resembles **isotonic regression** (Robertson et al., 1988), which seeks a weighted least squares fit under monotonicity constraints and has widespread applications in psychometrics, epidemiology, yield curve estimation, risk modeling and credit scoring studies. Compared with enforcing smoothness in  $T$ , our monotonicity approach avoids any need for a basis expansion or other parametric representation. Pursuing  $T$  reinterprets the data in a transformed domain where feature contributions recover an additive Shapley structure.

**Normalization.** A normalization is imposed on the transformed feature contributions,  $\sum_{j=1}^p (T(\beta_j))^2 = 1$ . This prevents degeneracy (e.g., trivial solutions such as  $T \equiv 0$ ) and anchors the scale of the model. An appealing feature of (13) is its invariance to the overall scaling of  $\{\nu_A\}$ , and the normalization constant is fixed at 1 without loss of generality. Moreover, Section 3 will show that imposing such a spherical constraint yields computational benefits, enabling a closed-form solution for the attribution-update and improving implementability.

**Sparsity.** (13) directly incorporates **sparsity** into the Shapley estimation process. Rather than relying on multi-step methods that first estimate a dense Shapley vector and then rank features (Slack et al., 2020), the formulation constrains the support of  $\hat{\beta}$  while pursuing the transformation during the iterative optimization process (cf. Algorithm 1). This *unified* treatment ensures that sparsity, domain adaptation, and Shapley coherence are achieved simultaneously, avoiding inconsistencies by post hoc selection. The popular  $\ell_1$ -penalty  $\lambda \sum |\beta_j|$  requires cumbersome  $\lambda$ -tuning and induces unwanted shrinkage. For example, it is generally unclear *a priori* how many nonzero coefficients result from a particular choice of  $\lambda$ . But our  $\ell_0$  constraint offers direct control over model sparsity and is entirely shrinkage-free, avoiding the attribution-distorting bias of  $\ell_1$ -type methods and the need for multi-step re-fitting procedures. It further remedies a well-documented limitation of  $\ell_1$  selection, which often fails to recover the true support in the

presence of even moderately correlated features (Zhao and Yu, 2006; Zhang, 2010). These properties make (13) particularly appealing for applications such as bioinformatics, where practitioners frequently require a fixed number of interpretable predictors. For cases where  $s$  must be chosen, we find the RIC criterion (Foster and George, 1994) to be effective within our Shapley framework.

Before concluding this section, we introduce a *reparameterization* trick that proves beneficial for both modeling and computation. Define

$$\gamma_j = T(\beta_j). \quad (14)$$

Since  $T$  is strictly increasing and  $T(0) = 0$  (the loss would become infinite if  $T(0) \neq 0$ ), (13) can be rewritten as

$$\min_{\gamma, T(\cdot)} \sum_{A \subseteq 2^F} w_{\text{SH}}(A) (T(\nu_A) - \sum_{j \in A} \gamma_j)^2 \text{ s.t. } \|\gamma\|_0 \leq s, \|\gamma\|_2 = 1, T \in \mathcal{M}. \quad (15)$$

The corresponding model assumption is thus  $T^*(\nu_A) \sim \mathcal{N}(\sum_{j \in A} \gamma_j^*, \sigma_A^2)$  for all  $A \subseteq 2^F$ , where the genuine transformation function  $T^*$  is monotonic with  $T^*(0) = 0$ , and  $\nu_A$  are assumed to be independent across different subsets  $A$ . The *starred* quantities represent the underlying statistical truth of interest to estimate. Assume  $\gamma^* \in \mathbb{R}^p$  satisfies  $\|\gamma^*\|_0 \leq s^*$  and  $\|\gamma^*\|_2 = 1$  with  $1 \leq s^* \leq p$  and  $s$  is specified as an upper bound on  $s^*$ . After estimating  $\hat{\gamma}$  and  $\hat{T}$  from (15), one can recover the  $\beta$ -scores by applying the inverse transformation  $\hat{\beta}_j = \hat{T}^{-1}(\hat{\gamma}_j)$ . This reconstructs the multivariate relationship between  $\nu_A$  and the set of feature contributions in the original scale, yielding  $\nu_A \approx \hat{T}^{-1}(\sum_{j \in A} \hat{T}(\hat{\beta}_j))$ , to offer interpretable Shapley-based attributions.

### 3 Optimization Algorithm

The optimization of SISR involves two main challenges: (i) a functional estimation component, and (ii) a combinatorial sparsity constraint coupled with a nonconvex normalization constraint. We show that the functional challenge can be addressed by a discretization technique, which, rather than introducing an approximation, preserves full equivalence. To handle the two constraints on  $\gamma$ , we develop a surrogate function framework. These efforts lead to an iterative procedure that combines the pool-adjacent-violators with

a normalized hard thresholding. Each step has implementation ease and the sparse structure ensures that the overall algorithm remains efficient in high-dimensional settings.

First, since  $T(\cdot)$  is only evaluated at the observed values  $\nu_A$  in the objective function, we “discretize” (15) by introducing the vector

$$t = [T(\nu_A)]_{A \subseteq 2^F} \in \mathbb{R}^{2^p}.$$

In defining this vector, one should fix a specific order over subsets  $A \subseteq F$ ; we follow the conventional lexicographic binary ordering to arrange the entries of  $t$ . Correspondingly, we define

$$\nu = [\nu_A]_{A \subseteq 2^F} \in \mathbb{R}^{2^p}, \quad \delta = \left[ \sum_{j \in A} \gamma_j \right]_{A \subseteq 2^F} = Z\gamma \in \mathbb{R}^{2^p},$$

where  $Z \in \mathbb{R}^{2^p \times p}$  is the “incidence matrix” indicating which features are active in each subset  $A$ , aligned with the same ordering used for  $t$ . Henceforth, we also write  $\nu_i$  (and likewise  $\delta_i$ ) to denote the entry corresponding to the  $i$ th subset. Additionally, introduce the diagonal weight matrix

$$W = \text{diag}\{w_{\text{SH}}(A)\}_{A \subseteq 2^F}. \quad (16)$$

With this notation in place, we study the following optimization problem:

$$\begin{aligned} \min_{\gamma \in \mathbb{R}^p, t \in \mathbb{R}^{2^p}} \quad & \frac{1}{2}(t - \delta)^\top W(t - \delta) \\ \text{subject to} \quad & \delta = Z\gamma, \quad \|\gamma\|_0 \leq s, \quad \|\gamma\|_2 = 1, \\ & t_i \leq t_j \quad \text{for all } (i, j) \in E(\nu) = \{(i, j) : \nu_i \leq \nu_j\}, \end{aligned} \quad (17)$$

where  $E$  encodes the pairwise ordering constraints induced by  $\nu$ , due to the monotonicity of the transformation  $T$ . This formulation replaces strict monotonicity with a non-decreasing constraint, a mild adjustment that facilitates numerical implementation. In the following of the section, we design a two-block alternating optimization algorithm.

First, with  $\delta$  fixed, the optimization over  $t$  corresponds to the (weighted) *isotonic regression*

$$\min_{t \in \mathbb{R}^{2^p}} \quad \frac{1}{2}(t - \delta)^\top W(t - \delta) \quad \text{subject to} \quad t_i \leq t_j \quad \text{for all } (i, j) \in E, \quad (18)$$

where the goal is to obtain a monotonic fit to  $\delta$  under a weighted squared-error loss defined by  $W$ . The problem can be solved using any standard Quadratic Programming (QP) solver, but it is more efficiently handled by the Pool-Adjacent-Violators Algorithm (**PAVA**)([de Leeuw et al., 2009](#)), which leverages the structure of the monotonicity constraints for improved computational performance.

Next, we focus on the  $\gamma$ -optimization.

**Theorem 1.** *Let  $\mathcal{H}(\cdot; s)$  denote the hard-thresholding operator associated with cardinality  $s$ , defined as follows: for a vector  $y \in \mathbb{R}^p$ ,  $\mathcal{H}(y; s) = z$  where  $z_i = y_i$  if  $|y_i|$  is among the  $s$  largest entries of  $|y_1|, \dots, |y_p|$ , and  $z_i = 0$  otherwise, and the normalized hard-thresholding operator  $\mathcal{H}^\circ(y; s) = \mathcal{H}(y; s) / \|\mathcal{H}(y; s)\|_2$  if  $\|\mathcal{H}(y; s)\|_2 \neq 0$ . Then, for the optimization problem with  $y \neq \vec{0}$ ,  $1 \leq s \leq p$ ,*

$$\min_{\beta} \frac{1}{2} \|y - \beta\|_2^2 \quad \text{subject to} \quad \|\beta\|_0 \leq s, \quad \|\beta\|_2 = 1,$$

*the vector obtained by normalized hard-thresholding,*

$$\hat{\beta} = \mathcal{H}^\circ(y; s) = \frac{\mathcal{H}(y; s)}{\|\mathcal{H}(y; s)\|_2}$$

*is a global optimizer.*

*Proof.* Let  $A \subseteq \{1, \dots, p\}$  and assume  $\beta_{A^c} = 0$  and  $\|\beta\|_2 = 1$ . Because

$$\begin{aligned} \|y - \beta\|_2^2 &= \|y\|_2^2 + \|\beta\|_2^2 - 2\langle y, \beta \rangle \\ &= \|y\|_2^2 + 1 - 2\langle y, \beta \rangle \\ &= \|y\|_2^2 + 1 - 2\langle y_A, \beta_A \rangle \\ &\geq \|y\|_2^2 + 1 - 2\|y_A\|_2 \|\beta_A\|_2 = \|y\|_2^2 + 1 - 2\|y_A\|_2, \end{aligned}$$

where we used the Cauchy-Schwarz inequality and the equality is achieved at  $\beta_A = y_A / \|y_A\|_2$ .

Therefore,  $\min_{\beta: \beta_{A^c}=0, \|\beta\|_2=1} \|y - \beta\|_2^2 = \|y\|_2^2 + 1 - 2\|y_A\|_2$  for any  $A : |A| = s$ . Minimizing over  $A$  gives an index set corresponding to the  $s$  largest entries of  $|y_1|, \dots, |y_p|$ , thereby the normalized hard thresholding operator  $\mathcal{H}^\circ(y; s)$ .  $\square$

We are now ready to develop an iterative algorithm for updating  $\gamma$  with  $t$  held fixed. Define the objective function

$$l(\gamma) = \frac{1}{2}(Z\gamma - t)^\top W(Z\gamma - t).$$

A straightforward calculation yields the gradient:

$$\nabla l(\gamma) = Z^\top W(Z\gamma - t). \quad (19)$$

To facilitate optimization, we construct a new “surrogate function”:

$$g(\gamma, \gamma^-) = l(\gamma^-) + \langle \nabla l(\gamma^-), \gamma - \gamma^- \rangle + \frac{\rho}{2} \|\gamma - \gamma^-\|_2^2. \quad (20)$$

where  $\rho > 0$  should be properly large (cf. Theorem 2). Define an iterative scheme:

$$\gamma^{(k+1)} = \arg \min_{\gamma} g(\gamma, \gamma^{(k)}) \quad \text{subject to } \|\gamma\|_0 \leq s, \|\gamma\|_2 = 1.$$

Using Theorem 1, the update step admits a closed-form expression:

$$\begin{aligned} \gamma^{(k+1)} &= \mathcal{H}^\circ(y; s) = \frac{\mathcal{H}(y; s)}{\|\mathcal{H}(y; s)\|_2}, \quad \text{with} \\ y &= \gamma^{(k)} - \frac{1}{\rho} \nabla l(\gamma^{(k)}) = \gamma^{(k)} - \frac{1}{\rho} Z^\top W(Z\gamma^{(k)} - t). \end{aligned} \quad (21)$$

**Theorem 2.** *Let  $\rho \geq \|Z^\top WZ\|_2$ , where  $\|\cdot\|_2$  denotes the matrix spectral norm. For any initial point  $\gamma^{(0)}$  satisfying  $\|\gamma^{(0)}\|_0 \leq s$  and  $\|\gamma^{(0)}\|_2 = 1$ , the sequence  $\{\gamma^{(k)}\}$  generated by (21) produces non-increasing (and thus convergent) function values:*

$$l(\gamma^{(k+1)}) \leq l(\gamma^{(k)}) \quad \text{for all } k \geq 0.$$

Furthermore, if  $\rho > \|Z^\top WZ\|_2$ ,  $\|\gamma^{(k+1)} - \gamma^{(k)}\|_2 \rightarrow 0$  as  $k \rightarrow \infty$ .

*Proof.* First, simple algebra shows

$$\begin{aligned} g(\gamma, \gamma^-) - l(\gamma) &= \frac{\rho}{2} \|\gamma - \gamma^-\|_2^2 - (l(\gamma) - l(\gamma^-) - \langle \nabla l(\gamma^-), \gamma - \gamma^- \rangle) \\ &= \frac{\rho}{2} \|\gamma - \gamma^-\|_2^2 - \frac{1}{2} (\gamma - \gamma^-)^\top H(\xi) (\gamma - \gamma^-) \\ &= \frac{1}{2} (\gamma - \gamma^-)^\top (\rho I - H(\xi)) (\gamma - \gamma^-), \end{aligned}$$

where we applied the mean-value theorem,  $H(\xi)$  denotes the Hessian matrix of  $l$  at  $\xi$  which is between  $\gamma$  and  $\gamma^-$ . Thus under the choice of  $\rho$ ,  $l(\gamma^{(k+1)}) \leq g(\gamma^{(k+1)}, \gamma^{(k)})$  for any  $k \geq 0$ .

By the optimality of  $\gamma^{(k+1)}$ ,  $g(\gamma^{(k+1)}, \gamma^{(k)}) \leq g(\gamma^{(k)}, \gamma^{(k)}) = l(\gamma^{(k)})$  and the first conclusion follows. Moreover, from the inequality:  $l(\gamma^{(k)}) - l(\gamma^{(k+1)}) \geq \frac{1}{2}(\gamma^{(k+1)} - \gamma^{(k)})^\top (\rho I - H(\xi))(\gamma^{(k+1)} - \gamma^{(k)}) \geq \frac{\rho - \|Z^\top W Z\|_2}{2} \|\gamma^{(k+1)} - \gamma^{(k)}\|_2^2$ , we obtain the second result.  $\square$

A summary of our algorithmic procedure is outlined in Algorithm 1. Some practical implementation notes: (i) The provided values  $\nu_A$  have been baseline adjusted as described in (7) (i.e., a preprocessing  $\nu_A \leftarrow \nu_A - \nu_\emptyset$  for all  $A \subseteq 2^F$  is assumed). We take  $C$  as  $1e+4$  if  $\|\nu\|_\infty \leq 10$ . (ii) For  $A = \emptyset$  or  $A = F$ , although  $w_{\text{SH}}(A)$  take infinite weights in theory, practically one can assign a weight equal to a large multiplier (e.g., 10) times the largest non-infinite weight (cf. (5)), which is often numerically sufficient to enforce  $\hat{T}(\nu_F) \doteq \sum_{j=1}^p \hat{T}(\hat{\beta}_j)$ . (iii) It is unnecessary to explicitly form the diagonal matrix  $W$ ; only the diagonal weights are required. Likewise, the sparsity of matrix  $Z$  can be utilized. Additionally, key quantities such as  $Z^\top W$  and  $Z^\top W t$  can be precomputed prior to the iterative updates to improve computational efficiency. (iv) In step 9), we employ a self-implemented, *stack*-based weighted PAVA for improved efficiency. (v) The paired data  $(\nu_i, \hat{t}_i)$  approximate  $T$  and form the basis for visualizing  $\hat{T}(\cdot)$ . A heuristic for  $\hat{T}^{-1}(\cdot)$  involves interpolating the inverted pairs  $(\hat{t}_i, \nu_i)$ , after averaging the  $\nu_i$  values for any duplicate  $\hat{t}_i$  coordinates to ensure a well-defined mapping. Alternatives include fitting a second weighted isotonic regression  $G(\hat{\delta})$  to  $\nu$ , or enforcing strict monotonicity in Step 9 (e.g., with an  $\epsilon$ -margin constraint) for an invertible  $\hat{T}$ . Overall, Algorithm 1 is straightforward to implement and scales very well in practice.

## 4 Data-Driven Insights

### 4.1 Domain Adaptation

To propose a convenient noisy data generation scheme, let's revisit the statistical model defined in Section 2, where the expectation  $\mathbb{E}[T^*(\nu)] = Z\gamma^*$ , with  $T(\cdot)$  applied componentwise. Assume without loss of generality that  $Z$  is structured using a bit generation process, with each row corresponding to



---

**Algorithm 1** Sparse Isotonic Shapley Regression (**SISR**) Algorithm

---

**Input:**  $\nu = [\nu_A]_{A \subseteq 2^F} \in \mathbb{R}^{2^p}$  (baseline-adjusted, such that  $\nu_\emptyset = 0$ ), sparsity level  $s$ , design matrix  $Z \in \mathbb{R}^{2^p \times p}$ , diagonal weight matrix  $W$  (cf. (16)), and an initial vector  $t^{(0)} \in \mathbb{R}^{2^p}$  (e.g.,  $C\nu$  with a large  $C$  if  $\|\nu\|_\infty$  is small, to improve precision, and  $C = 1$  otherwise).

```

1: Initialize  $t \leftarrow t^{(0)}$ ,  $\gamma \leftarrow 0$ 
2:  $\rho \leftarrow \|Z^\top W Z\|_2$ 
3: repeat
4:   while not converged do
5:      $\xi \leftarrow \mathcal{H}(\gamma - \frac{1}{\rho} Z^\top W (Z\gamma - t); s)$ 
6:      $\gamma \leftarrow \frac{\xi}{\|\xi\|_2}$ 
7:   end while
8:    $\delta \leftarrow Z\gamma$ 
9:   Fit isotonic regression (18) with  $\delta, W, Z$  to update  $t$ 
10: until convergence
11: return  $t, \gamma$ 

```

---

the binary representation of  $i - 1$  (e.g., the second row is  $[1, 0, \dots, 0]$ ). For  $\gamma^* = c_0[2^0, 2^1, \dots, 2^{p-2}, 2^{p-1}]^\top$ , this yields

$$\mathbb{E}[T^*(\nu)] = c_0[0, 1, \dots, 2^p - 2, 2^p - 1]^\top$$

where  $c_0 = \sqrt{\frac{3}{4^p - 1}}$  ensures that  $\gamma^*$  is normalized. To simulate this in experiments, we generate *noisy* versions  $\nu_A$  for all subsets  $A$  using

$$\nu = Q(c_1 \cdot \sigma(U)) \in \mathbb{R}^{2^p}$$

where  $U \in \mathbb{R}^{2^p}$  contains entries uniformly distributed between 0 and  $c_0(2^p - 1)$ , approximately  $\sqrt{3}$  when  $p$  is sufficiently large. Here,  $\sigma$  denotes the permutation that sorts the elements of  $U$  in ascending order. An accurate estimator,  $\hat{T}$  or  $\hat{t}$ , should then closely approximate the inverse transformation

$$T^* = Q^{-1}/c_1.$$

The inclusion of  $c_1$  is to ensure flexibility.

Figure 1 presents the results under 6 different functional forms for the true transformation  $T^*$ : *square root* ( $T^* = (\cdot)^{1/2}$ ), *fifth root* ( $T^* = (\cdot)^{1/5}$ ), *exponential* ( $T^* = \exp(\cdot) - 1$ ), *logarithmic* ( $T^* = \log(\cdot + 1)$ ), *tangent* ( $T^* =$

$\tan(\cdot)/c_1, c_1 = 10$ ), and *normal distribution* ( $T^* = \Phi(\cdot + c_2)/c_1, Q(\cdot) = \Phi^{-1}(c_1 \cdot) - c_2, c_1 = 1/\sqrt{3}, c_2 = Q(c_1 \sigma_{\min})$ ). As pointed out by a reviewer, comparing the estimated  $\hat{T}$  directly to the true  $T^*$ , rather than their inverses on the  $\nu$  scale, is statistically sound, as this aligns with our additive Gaussian assumption (9). Encouragingly, across all cases, the estimated transformation  $\hat{T}(\nu)$  closely aligns the ground-truth  $T^*(\nu)$ , providing strong empirical evidence for the effectiveness of SISR in accurately recovering the underlying transformation structure.

Finally, an additional experiment was conducted using data generated according to  $\nu_A = \max_{j \in A} \beta_j$ , where  $\beta_j = j$ . The resulting estimated transformation is displayed in Figure 2. The recovered transformation exhibits a pronounced increasing trend and nonlinearity. The estimates are nearly perfectly correlated with the transformed ground truths, and the best-fit line passes through the origin up to a scaling factor.

## 4.2 Sparsity Recovery

We generate data according to the sparse- $\gamma$  model in Section 2, with the true transformation set as the cubic root,  $T^*(\cdot) = \sqrt[3]{\cdot}$ . The true coefficient vector is  $\gamma^* = [1/\sqrt{3}, 1/\sqrt{3}, 1/\sqrt{3}, 0, \dots, 0]^\top$ , a relatively weak signal with sparsity level  $s^* = 3$ . The noise variance is defined as  $\sigma_A^2 = \sigma_0^2/w_{\text{SH}}(A)$ , with varying values of  $\sigma_0$ . The sparsity level upper bound  $s$  in running SISR is set to  $1.5 s^*$ . Performance is evaluated using two metrics. The first measures the alignment or affinity between the two unit-norm vectors:  $\langle \hat{\gamma}, \gamma^* \rangle \times 100$  (denoted by **Affn**), serving as an index of estimation accuracy. The second metric is the support recovery rate:  $|\text{supp}(\hat{\gamma}) \cap \text{supp}(\gamma^*)|/s^* \times 100\%$  (denoted by **Supp**), reflecting the proportion of correctly identified nonzero components in the true support. Table 1 reports results for varying values of  $p$  and  $\sigma_0$ . All results are averaged over 100 simulation runs.

As shown in the table, both performance metrics decline with increasing feature dimension  $p$  and noise level  $\sigma_0$ , as expected due to greater model complexity and reduced signal-to-noise ratio (SNR). Although not reported, running the algorithm without sparsity enforcement (i.e.,  $s = p$ ) yields noticeably worse affinity scores in high SNR settings (e.g., the first setting row of Table 1). Compared to the affinity scores, the support recovery rate remains surprisingly strong even under challenging conditions, indicating that SISR consistently identifies the correct features.

We further investigated the impact of the sparsity level  $s$  on computa-

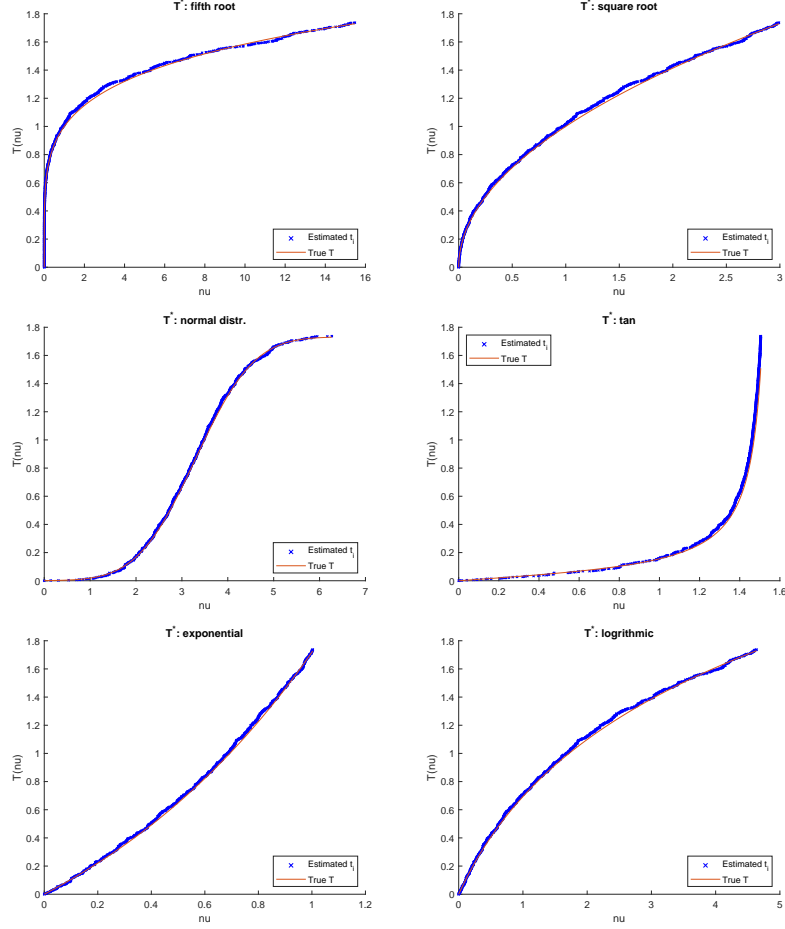


Figure 1: Estimated monotonic transformation  $\hat{T}(\nu)$  (in *blue*) versus the true transformation  $T^*$  (in *red*), for  $p = 10$  under 6 different functional forms for  $T^*$ : the fifth root, square root, normal distribution, tangent, exponential, and logarithmic transformations.

tional time. In this experiment, we varied  $s$  from 5 to 15, while fixing  $s^* = 3$ ,  $p = 15$ ,  $\sigma_0 = 5e-3$ . As illustrated in the figure, lower sparsity levels generally lead to faster computation, highlighting the efficiency gains achievable when enforcing proper sparsity in the model.

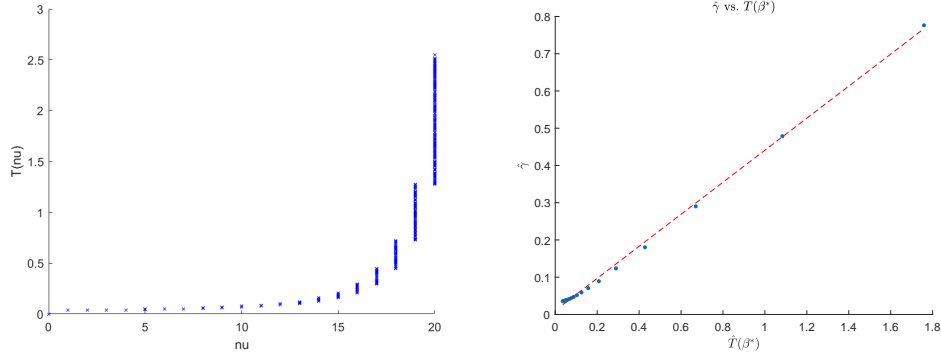


Figure 2: Estimated monotonic transformation  $\hat{T}(\nu)$  (left) and comparison between  $\hat{\gamma}$  vs  $\hat{T}(\beta^*)$  (right, showing an almost perfect correlation of 1.00) for  $p = 20$  under a winner-takes-all setting.

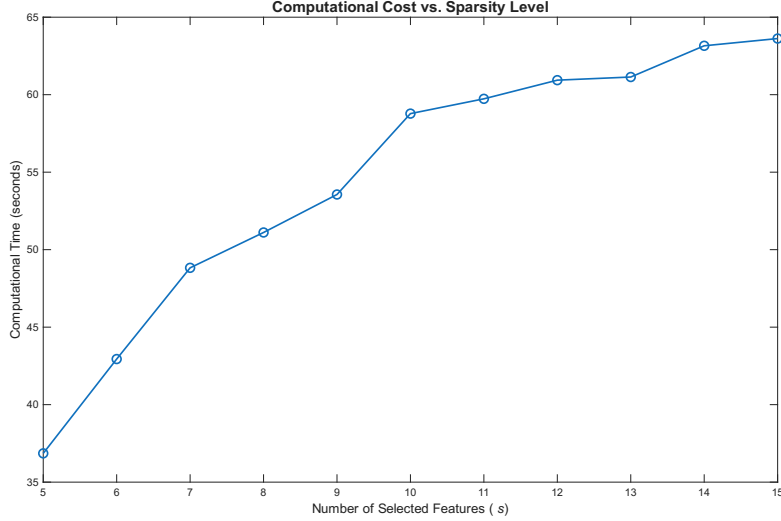


Figure 3: Computational time versus sparsity level.

### 4.3 $R^2$ -Payoffs in Regression/Logistic Regression

In regression settings, coalition values for feature subsets are commonly defined using the *coefficient of determination* ( $R^2$ ) obtained from retraining the model on each subset, reflecting the scaled improvement in model fit (Lipovetsky and Conklin, 2001; Covert et al., 2021). Contrary to conventional expectations, our results unveil a novel insight: such a standard construction

Table 1: Affinity score (**Affn**) and support recovery rate (**Supp**) across different values of  $p$  and noise level  $\sigma_0$ . Larger values reflect better performance.

	$\sigma_0 = 1\text{e-}3$		$\sigma_0 = 5\text{e-}3$		$\sigma_0 = 1\text{e-}2$	
	Affn	Supp	Affn	Supp	Affn	Supp
$p = 10$	99.6	100%	99.6	100%	99.5	100%
$p = 15$	99.8	100%	99.9	100%	97.8	100%
$p = 20$	99.9	100%	87.9	100%	80.3	100%
$p = 25$	87.9	100%	74.0	100%	70.5	100%
	$\sigma_0 = 5\text{e-}2$		$\sigma_0 = 1\text{e-}1$		$\sigma_0 = 2\text{e-}1$	
	Affn	Supp	Affn	Supp	Affn	Supp
$p = 10$	97.9	100%	88.7	98.7%	66.2	80.7%
$p = 15$	79.9	100%	70.9	98.0%	57.6	73.3%
$p = 20$	68.9	100%	63.2	96.0%	54.3	65.3%
$p = 25$	65.5	100%	60.6	90.7%	52.1	62.0%

can fail to yield an inherently additive Shapley framework, especially when features are dependent or include irrelevant ones, which are almost certain to occur in practice.

To illustrate this phenomenon, let's consider the following simulation setup:  $y = X\alpha^* + \epsilon$ , where  $\epsilon_i \sim \mathcal{N}(0, 1)$  and each row of  $X \in \mathbb{R}^{n \times p}$  is drawn from a multivariate normal distribution with mean zero and Toeplitz covariance  $\Sigma_{ij} = \theta^{|i-j|}$  with  $\theta = 0.5$ . The true coefficient vector is set as  $\alpha^* = [3, 3, \dots, 3]^\top$ , and the sample size is  $n = 5p$ . For each subset  $A$ , we fit a regression model using the predictors indexed by  $A$  and define  $\nu_A$  as the resulting  $R^2$  value. Logistic regression is also considered in the classification setting,  $y_i \sim \text{Bernoulli}(\pi_i)$  and  $\text{logit}(\pi) = X\alpha^*$ , where we generate  $y$  according to a Bernoulli model and define  $\nu_A$  using the *deviance*-based pseudo- $R^2$ .

As shown in Figure 4, the estimated transformation  $\hat{T}(\nu)$  deviates significantly from linearity. For instance, in the regression case, even a simple logarithmic transformation fails to linearize the relationship, whereas a log-log transformation produces a nearly linear pattern—suggesting that the underlying transformation is super-exponential.

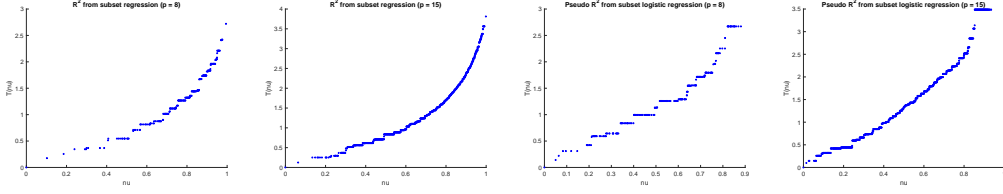


Figure 4: Estimated monotonic transformation  $\hat{T}(\nu)$  using regression-based  $R^2$  and logistic regression-based pseudo- $R^2$  as the coalition worth function, for  $p = 8, 15$ .

To examine whether model sparsity and feature correlation play a role in shaping the transformation  $T$ , we conducted a factorial simulation for linear regression with  $p = 15$  and  $s = p$ . We fixed the coefficient vector to  $\alpha^* = [3, 0, 3, 0, \dots, 0]^\top$  with the sparsity level  $s^* \in \{2, 8, 15\}$ , and varied the correlation parameter  $\theta \in \{0, 0.5, 0.9\}$ . Hence the design ranges from independent ( $\theta = 0$ ) to collinear ( $\theta = 0.9$ ) predictors, and from very sparse to fully dense signals. Figure 5 displays the empirical transformation  $\hat{T}$  recovered in each setting.

Two main patterns emerge. (i) **Correlation drives curvature:** as  $\theta$  increases,  $\hat{T}$  deviates sharply from linearity, even in dense models. (ii) **Sparsity introduces breaks:** at low  $s^*$  the curve becomes piecewise, with segment slopes that differ substantially—another marker of non-additivity. Notably, pronounced nonlinearity appears even in the independent, ultra-sparse case ( $\theta = 0$ ,  $s^* = 2$ ), showing that irrelevant features alone can distort raw worths. Hence a *monotone nonlinear* transformation is indispensable when translating  $R^2$ -based worth measures into the additive Shapley framework; without it, either strong correlation or feature irrelevance breaks the required additivity. To the best of our knowledge, our results are the first to reveal that correlated features and the presence of irrelevant features can substantially undermine the additivity assumption in Shapley frameworks.

## 4.4 Prostate Cancer

We used the prostate data from Tibshirani (1996) and took  $\log(\text{cancer volume})$  (`lcavol`) as the response, with clinical predictors including  $\log(\text{prostate weight})$  (`lweight`), age (`age`),  $\log(\text{benign prostatic hyperplasia})$  (`lbph`), seminal vesicle invasion (`svi`, binary),  $\log(\text{capsular penetration})$  (`lcp`), Gleason

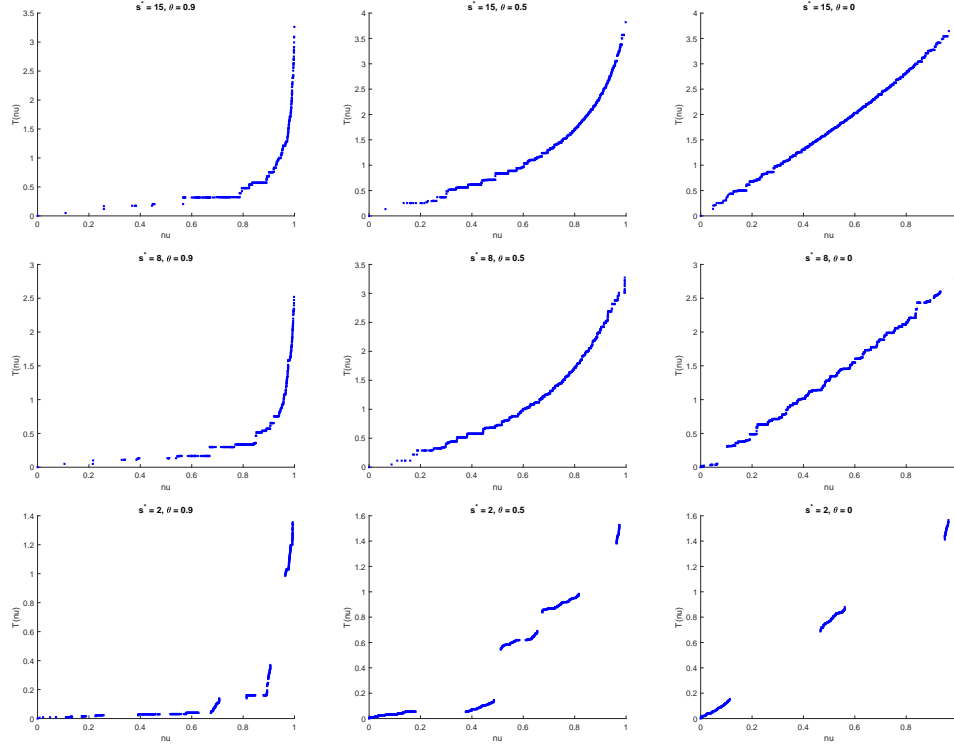


Figure 5: Estimated monotonic transformation  $\hat{T}(\nu)$  across varying sparsity levels ( $s = 15, 8, 2$ , top to down) and feature correlation strengths ( $\theta = 0.9, 0.5, 0$ , left to right). The RIC criterion identifies  $s = 6$  as optimal.

score (`gleason`), percentage Gleason 4/5 cells (`pgg45`), and log(prostate specific antigen) (`lpsa`). We computed SISR calibrated Shapley values  $\hat{\gamma}$  from  $R^2$ -worths at sparsity  $s \in \{8, 6, 4\}$  and contrasted them with the conventional (raw) Shapley values in Fig. 6.

Both schemes agree that `lcp` and `lpsa` dominate the signal. The striking discrepancy concerns `svi`: the naive Shapley ranking elevates it to third place, claiming more than 10% of the total importance, whereas the calibrated  $\hat{\gamma}$  assigns it virtually none.

This SISR finding turns out to be the one that aligns with established evidence. Statistically, independent checks confirm that `svi` contributes little: stepwise AIC and BIC both discard it, it is the final variable selected on the LASSO path (Tibshirani, 1996), and its  $p$ -value is as large as 0.6 in the full model. The conventional Shapley result is also biologically implausible, since





pel et al., 2011), which down-weights the penalty from large, unpredictable errors. Figure 7 displays the resulting Shapley attributions.

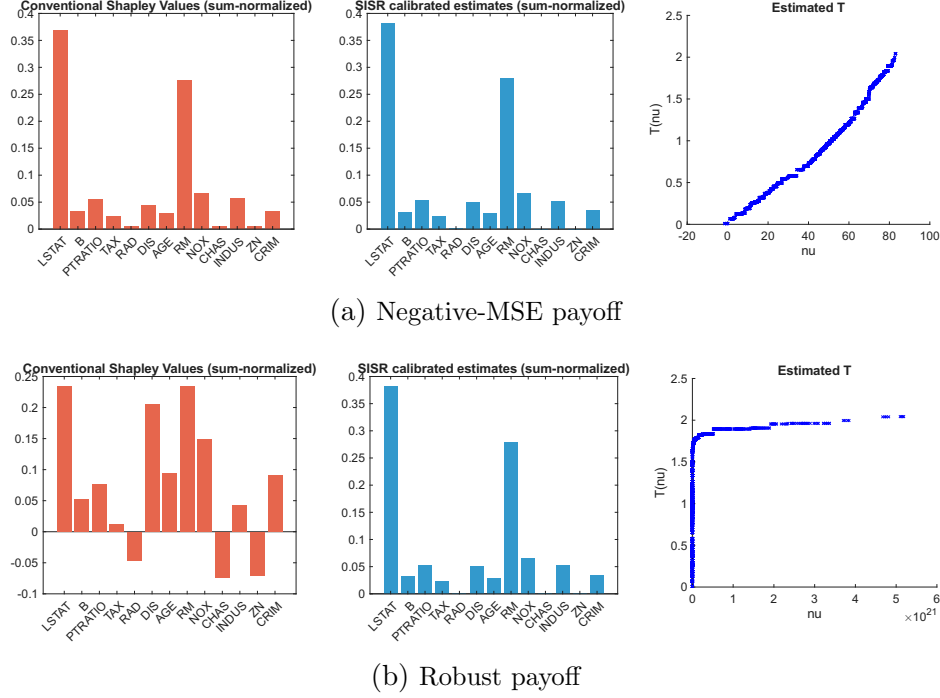


Figure 7: Boston housing: feature attributions computed with conventional Shapley and SISR-calibrated Shapley values for the negative-MSE payoff (top) and the robust payoff (bottom), along with the corresponding estimated monotone transformations.

The figure contrasts how the two attribution schemes respond to different payoff schemes. Under the MSE payoff, SISR has little to adjust—the scale is already compatible with linear additivity. In contrast, for the robust payoff, SISR produces a highly nonlinear transformation, which compensates for the distortions and preserves essentially the same attribution pattern observed under the MSE scale. The conventional Shapley values shift noticeably: the importance of **DIS** increases from minor to leading, and **CHAS** and several other variables even receive negative attributions. These sign and rank changes substantially alter the qualitative interpretation of the game and reveal the standard procedure’s sensitivity to the underlying payoff construction, whereas SISR remains robust.

## 4.6 Bank Credit

We analyze the South German Credit dataset, a benchmark for credit risk classification (Covert et al., 2020). The dataset contains 1,000 observations, where the response is a binary indicator of credit risk, and 20 predictor variables. These features include checking status, duration, credit history, age and so on. Following the experimental setup in Covert et al. (2020), we trained a CATBOOST classification model (Prokhorenkova et al., 2018) on a training set. With  $p = 20$  features, computing the full  $2^{20}$  coalition values is computationally intractable. To approximate this full game, we selected a representative subset of 1,000 coalitions, using the efficient sampling strategy proposed by Covert and Lee (2021). For each of these sampled coalitions, the payoff  $\nu_A$  was then defined as the model’s global performance on a background test set, following the interventional SAGE methodology (Covert et al., 2020). The two payoff functions we considered are the negative cross entropy,  $\nu_A^{\text{ent}}$ , and an exponential utility counterpart,  $\nu_A^{\text{exp}} = -\exp(-c\nu_A^{\text{ent}})$ . The first payoff corresponds to the standard logistic log-likelihood, while the second payoff models a strong risk-averse preference, a foundational principle in modern economics and decision theory (Pratt, 1964; Mas-Colell et al., 1995).

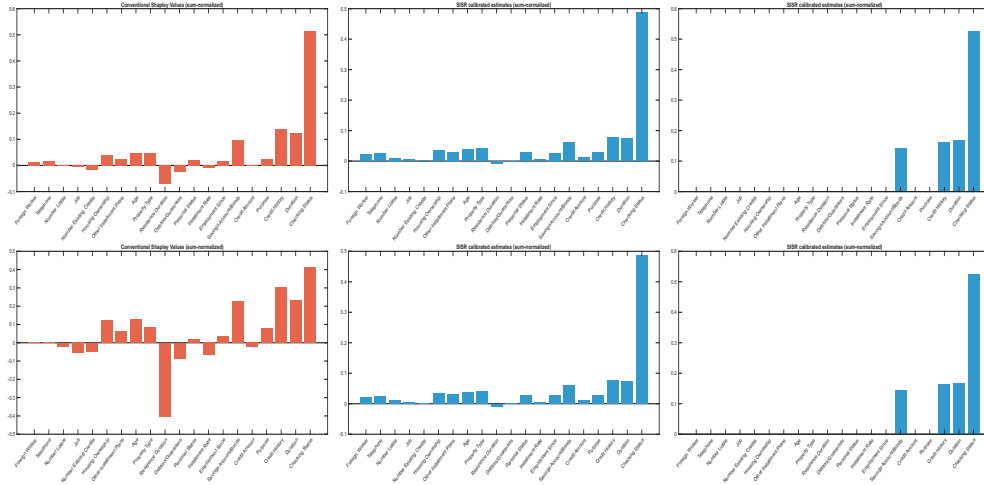


Figure 8: Bank credit data. The top row shows the feature attributions for the negative cross-entropy payoff, displaying (left to right): conventional Shapley values, SISR without sparsity, and SISR with  $s = 4$  (selected by RIC). The bottom row shows the same comparisons for the risk-averse exponential payoff.

Figure 8 displays the resulting feature attributions. The conventional Shapley results (top-left) differ from those in the original SAGE analysis, likely due to sampling-based approximation errors. In particular, **Residence Duration** exhibits a spurious negative attribution. This instability becomes even more pronounced under the risk-averse exponential payoff (bottom-left), where the negative value grows to nearly four times its original magnitude.

In contrast, the SISR-calibrated attributions (with and without sparsity) remain remarkably stable. **Residence Duration** is assigned a near-zero contribution—consistent with the mild effect observed in Covert et al. (2020). By removing nonlinear distortions in the payoff construction, SISR reveals the same sparse and interpretable structure across both settings, effectively filtering out the distortions that undermine conventional Shapley estimates.

## 4.7 Diabetes

The Pima Indians Diabetes dataset (Smith et al., 1988) was collected by the U.S. National Institute of Diabetes and Digestive and Kidney Diseases. It records eight medical measurements for 768 women of Pima heritage near Phoenix, Arizona; the response indicates a physician’s diagnosis of diabetes. Key predictors include plasma-glucose concentration two hours after an oral glucose-tolerance test (**Glucose**), body-mass index (**BMI**), and the diabetes-pedigree function (**DiabetesPedigreeFunction**). We trained an XGBOOST classifier (Chen and Guestrin, 2016) and tuned its hyper-parameters with GridSearchCV in scikit-learn (Pedregosa et al., 2011). Feature-subset worth was measured using two payoff functions: the negative cross-entropy (or logistic log-likelihood) and the likelihood payoff (analogous to the utility variant in Section 4.5). The results are shown in Figure 9.

Using the negative-entropy payoff, the estimated transformation is piecewise linear—an indicator of sparsity in the underlying feature importances (cf. Section 4.3). Although the segment slopes differ only modestly, the resulting shift alters several attributions: the largest discrepancy is for **DiabetesPedigreeFunction**, whose importance is noticeably lower under the standard Shapley calculation than in the SISR-calibrated version. The bottom panel reports results for the likelihood payoff. Here, the conventional Shapley values are severely distorted: numerous features appear almost as influential as **Glucose**, and **Pregnancies** even becomes negative, illustrating how conventional attributions can be distorted by such an utility function. By contrast, the SISR estimates without sparsity nearly replicate those ob-

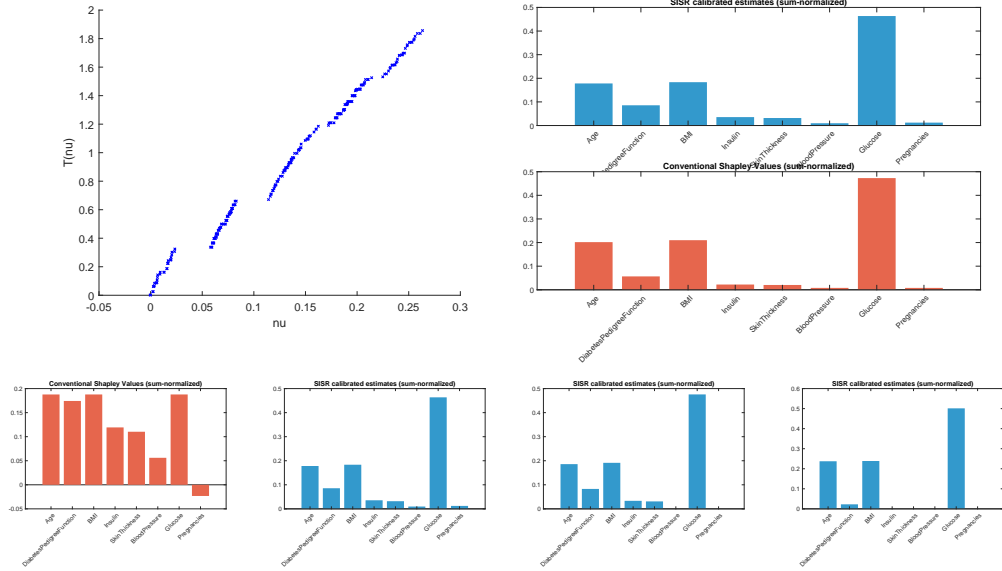


Figure 9: Diabetes data. Top row (negative-entropy payoff): estimated transformation  $\hat{T}(\nu)$  obtained with SISR (left) and feature attributions from SISR-calibrated versus conventional Shapley values (right). Bottom row (likelihood or exponential payoff): feature attributions from the conventional method (leftmost) and from SISR for sparsity levels  $s \in \{8, 6, 4\}$  (left to right).

tained under negative entropy, and introducing sparsity leaves the ranking largely unchanged. Overall, SISR delivers stable attributions across different value (payoff) functions, whereas the standard procedure is highly sensitive to which value function is chosen.

## 5 Discussion and Extensions

Recent work has extended the classical Shapley value to account for higher-order interactions, including the Shapley-GAM framework and a wide family of Shapley Interaction (SI) indices, to provide a finer decomposition of the payoff function  $\nu$  (Grabisch and Roubens, 1999; Marichal, 2000; Sundararajan et al., 2020; Tsai et al., 2023; Bordt and von Luxburg, 2023). Unlike the classical Shapley value, which is uniquely determined by its axioms, interaction indices are not: different extensions of the axiomatic framework or design criteria lead to distinct definitions and thus non-unique attributions,

as exemplified by  $k$ -SII, STII, FSII, among many others (Muschalik et al., 2024). Computationally, all such methods scale unfavorably with the number of features  $p$ , and even approximate variants remain expensive when  $k$  grows beyond two (Sundararajan et al., 2020). From an interpretability perspective, explanatory power typically wanes with order: higher-order effects become increasingly hard to communicate and may overwhelm practitioners with “information overload” (Baniecki et al., 2024). From a statistical standpoint, complex interaction terms may instead reflect noise or dependence-induced artifacts, yielding patterns that appear compelling but lack substantive validity (Hooker, 2007).

In comparison, the SISR framework arises from a different—and often overlooked—misspecification. In many applications, the apparent non-additivity of the payoff does not reflect genuine high-order model logic but rather a nonlinear *distortion* of the value function  $\nu$ , introduced through preprocessing, surrogate construction, or sampling approximations. Existing interaction-based methods implicitly assume that  $\nu_A$  is a clean, faithful, and **Gaussian** measure of model worth, so that any deviation from additivity must represent true interaction. Ignoring the noise, the  $T^{-1}$ - $\Sigma$ - $T$  structure of SISR can also be rephrased (via Taylor expansion) in terms of nonlinear interactions—for example, if  $\beta_j \geq 0$  and  $T(\cdot) = \log(1 + \cdot)$ , the resulting value function  $V_A(\{\beta_j\}_{j \in A}) = \prod_{j \in A} (1 + \beta_j) - 1 = \sum_{j \in A} \beta_j + \sum_{j, j' \in A} \beta_j \beta_{j'} + \dots$ . But payoff constructions frequently violate Gaussianity: skewed, heavy-tailed, or bounded responses, as well as feature dependence or irrelevant covariates, can all distort  $\nu$  and induce *spurious* interaction where no genuine interaction would be found after a stabilization transformation. SISR addresses this fundamental issue by jointly estimating a monotone transformation  $T$  and attribution effects, recovering an additive, Gaussian working model that restores interpretability.

We can unify stabilization and interaction to form a powerful framework for nonlinear XAI. For example, we can generalize (9) by including interaction terms on the transformed scale:

$$T(\nu_A) \sim \mathcal{N}\left(\sum_{j \in A} T(\beta_j^*) + \sum_{\{i, k\} \subseteq A} T(\beta_{ik}^*) + \dots, \sigma_A^2\right),$$

subject to appropriate structural constraints such as monotonicity, sparsity, and normalization. Such a framework would simultaneously correct for non-Gaussian payoff distortions via  $T$  and capture genuine synergistic effects via the interaction coefficients. This represents a promising yet computationally

demanding direction for future research.

As noted by a reviewer, the methodology (e.g., (9)) extends from the squared-loss regression setting to a generalized linear model framework (McCullagh and Nelder, 1989). In particular, the loss in (11) or (13) is now rephrased as

$$\sum_{A \subseteq 2^F} w_{\text{SH}}(A) L(\delta_A; T(\nu_A)), \quad \delta_A = \sum_{j \in A} T(\beta_j), \quad L(\eta, y) = -\eta y + b(\eta),$$

where  $b$  is the log-partition function in the specified natural exponential distribution and  $b^{-1}$  induces the canonical link (Hardin and Hilbe, 2018). The resulting optimization subject to the same monotonicity constraints and sparsity/normalization constraints as in (13) may be termed as the sparse isotonic Shapley GLM. When  $b(\eta) = \eta^2/2$ , this reduces to the squared-loss formulation analyzed in previous sections.

An appealing fact is that the overall computational framework in Section 3 carries over seamlessly. The update of  $t$  (and hence  $\delta$ ) is obtained by an order-restricted fit using the isotonic machinery adapted to the GLM loss; see, e.g., de Leeuw et al. (2009) and Luss and Rosset (2014) for related PAVA and active-set extensions. The update of  $\gamma$  uses exactly the same surrogate function as in (20) with the working loss  $l(\gamma) = \sum_{A \subseteq 2^F} w_{\text{SH}}(A) L(\delta_A; T(\nu_A))$ , for which a direct calculation yields  $\nabla l(\gamma) = Z^\top W(b'(Z\gamma) - t)$ , recovering (19) in the Gaussian case. The nonlinear operators from Theorem 1 are unchanged and Theorem 2 continues to apply with a distribution-specific curvature bound. For example, in logistic regression where  $b(\eta) = \log(1 + \exp(\eta))$ , the lower bound on  $\rho$  becomes  $\|Z^\top W Z\|_2/4$  (since  $b'' \leq 1/4$ ). Derivations follow Section 3 closely and are omitted.

On the other hand, the Shapley interpretability is most natural in the Gaussian working model; moving to non-Gaussian families weakens the link to Shapley’s foundational axioms. Also, in our experience, practical pay-off functions rarely exhibit count-type behavior (as in Poisson models) or binary-valued outcomes (as in Bernoulli models), and we have not identified a compelling application in these settings. Therefore, although the GLM extension is methodologically sound and straightforward to implement, we leave its detailed investigation to future work.

## 6 Conclusion

This paper introduced *Sparse Isotonic Shapley Regression* (SISR) to address two pressing limitations of Shapley values in the context of XAI: the blind application of an additive main-effect valuation to real-world payoff constructions driven by non-Gaussian distributions or domain-specific loss scales, and the lack of native sparsity control in high-dimensional attribution tasks.

Rather than abandoning the interpretability provided by a simple, additive structure of individual feature contributions, SISR is designed to restore it. Our framework achieves nonlinear explainability by jointly learning a data-driven monotonic transformation via weighted isotonic regression and enforcing an  $\ell_0$  sparsity constraint through normalized hard-thresholding. Remarkably, the monotonic ordering constraint allows SISR to bypass a closed-form specification of the transformation, and sparsity both improves interpretability and accelerates computation in large feature spaces. Each alternating step admits a closed-form update, enjoys global convergence guarantees, and is straightforward to implement.

Our  $T$ -additive Shapley model is built upon this novel capability of “learning to be additive,” which is able to simultaneously recover the true payoff transformation and sparse support. We reveal for the first time that the mere existence of irrelevant features or inter-feature dependence can induce a payoff transformation that departs substantially from linearity, highlighting the necessity of nonlinear explainability. Extensive experiments validate that SISR stabilizes attributions across radically different payoff constructions, correctly filters out spurious features, and aligns with established diagnostics, whereas standard Shapley values suffer severe rank and sign distortions.

By unifying domain adaptation and sparsity pursuit within the Shapley framework, SISR advances the frontier of nonlinear explainability, providing a theoretically grounded, robust, and scalable attribution methodology.

## References

- Algaba, E., Fragnelli, V., and Sánchez-Soriano, J. (2019). *Handbook of the Shapley Value*. Chapman & Hall/CRC Series in Operations Research. CRC Press.
- Ancona, M., Ceolini, E., Öztireli, C., and Gross, M. (2019). Explaining deep neural networks with a polynomial time algorithm for shapley value ap-

- proximation. *Proceedings of the 36th International Conference on Machine Learning*.
- Au, Q., Herbringer, J., Stachl, C., Bischl, B., and Casalicchio, G. (2022). Grouped feature importance and combined features effect plot. *Data Min. Knowl. Discov.*, 36(4):1401–1450.
- Baniecki, H., Parzych, D., and Biecek, P. (2024). The grammar of interactive explanatory model analysis. *Data Mining and Knowledge Discovery*, 38:2596–2632.
- Bordt, S. and von Luxburg, U. (2023). From shapley values to generalized additive models and back. In *Proceedings of the 26th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 206 of *Proceedings of Machine Learning Research*, pages 709–745.
- Charnes, A., Golany, B., Keane, M., and Rousseau, J. (1988). *Extremal Principle Solutions of Games in Characteristic Function Form: Core, Chebyshev and Shapley Value Generalizations*, pages 123–133. Springer Netherlands, Dordrecht.
- Chen, T. and Guestrin, C. (2016). Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 785–794.
- Cohen, S., Dror, G., and Ruppin, E. (2007). Feature selection via coalitional game theory. *Neural Computation*, 19(7):1939–1961.
- Covert, I. and Lee, S.-I. (2021). Improving kernelshap: Practical shapley value estimation via linear regression. In *Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 130 of *Proceedings of Machine Learning Research*, pages 3457–3465.
- Covert, I., Lundberg, S. M., and Lee, S.-I. (2020). Understanding global feature contributions with additive importance measures. *Advances in neural information processing systems*, 33:17212–17223.
- Covert, I. C., Lundberg, S., and Lee, S.-I. (2021). Explaining by removing: a unified framework for model explanation. *The Journal of Machine Learning Research*, 22(209):9477–9566.



- de Leeuw, J., Hornik, K., and Mair, P. (2009). Isotone optimization in R: Pool-adjacent-violators algorithm (PAVA) and active set methods. *Journal of Statistical Software*, 32(5):1–24.
- Debras, B., Guillonnet, B., Bougaran, J., Chambon, E., and Vallancien, G. (1998). Prognostic significance of seminal vesicle invasion on the radical prostatectomy specimen. rationale for seminal vesicle biopsies. *European Urology*, 33(3):271–277.
- Duan, H. and Okten, G. (2025). Derivative-based shapley value for global sensitivity analysis and machine learning explainability. *International Journal for Uncertainty Quantification*, 15(1):1–16.
- Foster, D. P. and George, E. I. (1994). The risk inflation criterion for multiple regression. *The Annals of Statistics*, 22(4):1947–1975.
- Fryer, D., Strümke, I., and Nguyen, H. (2021). Shapley values for feature selection: The good, the bad, and the axioms. *Ieee Access*, 9:144352–144360.
- Grabisch, M. and Roubens, M. (1999). An axiomatic approach to the concept of interaction among players in cooperative games. *International Journal of Game Theory*, 28(4):547–565.
- Hampel, F. R., Ronchetti, E. M., Rousseeuw, P. J., and Stahel, W. A. (2011). *Robust statistics: the approach based on influence functions*, volume 114. John Wiley & Sons.
- Hardin, J. W. and Hilbe, J. M. (2018). *Generalized Linear Models and Extensions*. Stata Press, College Station, TX, 4 edition.
- Harrison, D. and Rubinfeld, D. L. (1978). Hedonic housing prices and the demand for clean air. *Journal of Environmental Economics and Management*, 5(1):81–102.
- Hooker, G. (2007). Generalized functional anova diagnostics for high dimensional functions of dependent variables. *Journal of Computational and Graphical Statistics*, 16(3):709–732.
- Huber, P. J. and Ronchetti, E. M. (2009). *Robust statistics*. John Wiley & Sons, Hoboken, NJ, Second edition.

- Jothi, N., Husain, W., and Rashid, N. A. (2021). Predicting generalized anxiety disorder among women using shapley value. *Journal of Infection and Public Health*, 14(1):103–108.
- Kristiansen, A., Wiklund, F., Wiklund, P., and Egevad, L. (2013). Prognostic significance of patterns of seminal vesicle invasion in prostate cancer. *Histopathology*, 62(7):1049–1056.
- Lipovetsky, S. and Conklin, M. (2001). Analysis of regression in game theory approach. *Applied Stochastic Models in Business and Industry*, 17(4):319–330.
- Lundberg, S. M., Erion, G. G., Chen, H., DeGrave, A. J., Prutkin, J. M., Nair, B., Katz, R., Himmelfarb, J., Bansal, N., and Lee, S.-I. (2020). From local explanations to global understanding with explainable ai for trees. *Nature Machine Intelligence*, 2(1):56–67.
- Lundberg, S. M. and Lee, S.-I. (2017). A unified approach to interpreting model predictions. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 30, pages 4768–4777. Curran Associates, Inc.
- Luss, R. and Rosset, S. (2014). Generalized isotonic regression. *Journal of Computational and Graphical Statistics*, 23(1):192–210.
- Ma, S. and Tourani, R. (2020). Predictive and causal implications of using shapley value for model interpretation. In *Proceedings of the 2020 KDD Workshop on Causal Discovery*, volume 127 of *Proceedings of Machine Learning Research*, pages 23–38. PMLR.
- Maniar, A. (2023). Xgboost model optimization – boston housing. <https://www.kaggle.com/code/advikmaniar/xgboost-model-optimization-94-boston-housing/notebook>.
- Marichal, J.-L. (2000). An axiomatic approach of the discrete choquet integral as a tool to aggregate interaction effects. *IEEE Transactions on Fuzzy Systems*, 8(6):800–807.
- Mas-Colell, A., Whinston, M. D., and Green, J. R. (1995). *Microeconomic Theory*. Oxford University Press.

- McCullagh, P. and Nelder, J. A. (1989). *Generalized Linear Models*. Chapman and Hall/CRC Monographs on Statistics and Applied Probability Series. Chapman & Hall/CRC, London.
- Muschalik, M., Fumagalli, F., Hammer, B., and Hüllermeier, E. (2024). shapiq: Shapley interactions for machine learning. *Advances in Neural Information Processing Systems*, 37:130324–130357.
- Owen, A. B. (2014). Sobol’ indices and Shapley value. *SIAM/ASA Journal on Uncertainty Quantification*, 2(1):245–251.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Pratt, J. W. (1964). Risk aversion in the small and in the large. *Econometrica*, 32(1-2):122–136.
- Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A. V., and Gulin, A. (2018). Catboost: unbiased boosting with categorical features. In *Advances in Neural Information Processing Systems*, volume 31, pages 6638–6648.
- Ribeiro, M. T., Singh, S., and Guestrin, C. (2016). “Why should I trust you?” Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144.
- Robertson, T., Wright, F. T., and Dykstra, R. L. (1988). *Order restricted statistical inference*. Wiley Series in Probability and Mathematical Statistics: Probability and Mathematical Statistics. John Wiley & Sons Ltd., Chichester.
- Shapley, L. S. (1953). A value for n-person games. In Kuhn, H. W. and Tucker, A. W., editors, *Contributions to the Theory of Games II*, pages 307–317. Princeton University Press, Princeton.
- Slack, D., Hilgard, S., Jia, E., Singh, S., and Lakkaraju, H. (2020). Fooling lime and shap: Adversarial attacks on post hoc explanation methods. *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*.

- Smith, J. W., Everhart, J. E., Dickson, W. C., Knowler, W. C., and Johannes, C. E. (1988). Using the ADA diagnostic criteria to classify individuals and test glycemic models. *Proceedings of the IEEE Symposium on Computer-Based Medical Systems*, pages 261–264. Dataset popularly known as the Pima Indians Diabetes dataset.
- Strumbelj, E. and Kononenko, I. (2014). Explaining prediction models and individual predictions with feature contributions. *Knowledge and Information Systems*, 41:647–665.
- Sundararajan, M., Dhamdhere, K., and Agarwal, A. (2020). The shapley Taylor interaction index. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *JRSSB*, 58:267–288.
- Tsai, C. J., Yeh, C.-W., and Ravikumar, P. (2023). Faith-shap: An axiomatic framework for faithful interaction attribution. *Journal of Machine Learning Research*, 24(281):1–64.
- Zhang, C.-H. (2010). Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics*, 38(2):894–942.
- Zhao, P. and Yu, B. (2006). On model selection consistency of lasso. *Journal of Machine Learning Research*, 7(81):2541–2563.