

SPLINES-BASED FEATURE IMPORTANCE IN KOLMOGOROV-ARNOLD NETWORKS: A FRAMEWORK FOR SUPERVISED TABULAR DATA DIMENSIONALITY REDUCTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Feature selection is a key step in many tabular prediction problems, where multiple candidate variables may be redundant, noisy, or weakly informative. We investigate feature selection based on Kolmogorov-Arnold Networks (KANs), which parameterize feature transformations with splines and expose per-feature importance scores in a natural way. From this idea we derive four KAN-based selection criteria (coefficient norms, gradient-based saliency, and knockout scores) and compare them with standard methods such as LASSO, Random Forest feature importance, Mutual Information, and SVM-RFE on a suite of real and synthetic classification and regression datasets. Using average F1 and R^2 scores across three feature-retention levels (20%, 40%, 60%), we find that KAN-based selectors are generally competitive with, and sometimes superior to, classical baselines. In classification, KAN criteria often match or exceed existing methods on multi-class tasks by removing redundant features and capturing nonlinear interactions. In regression, KAN-based scores provide robust performance on noisy and heterogeneous datasets, closely tracking strong ensemble predictors; we also observe characteristic failure modes, such as overly aggressive pruning with an ℓ_1 criterion. Stability and redundancy analyses further show that KAN-based selectors yield reproducible feature subsets across folds while avoiding unnecessary correlation inflation, ensuring reliable and non-redundant variable selection. Overall, our findings demonstrate that KAN-based feature selection provides a powerful and interpretable alternative to traditional methods, capable of uncovering nonlinear and multivariate feature relevance beyond sparsity or impurity-based measures.

1 INTRODUCTION

Feature selection plays a central role in many machine learning pipelines, where models must cope with redundant, noisy, or weakly informative input variables. By identifying and retaining only the most relevant features, it can improve both interpretability and computational efficiency. Classical approaches to feature selection are often grouped into three categories: *filter* methods, which rank features using data-driven scores independent of the predictor; *wrapper* methods, which evaluate subsets of features through repeated model training; and *embedded* methods, which perform selection as part of the model fitting procedure itself.

Filter methods include statistical criteria such as Pearson correlation and mutual information (MI) (Peng et al., 2005), which are fast and model-agnostic but ignore feature interactions and often select redundant features. Embedded methods integrate selection into the learning algorithm itself, such as Least Absolute Shrinkage and Selection Operator (LASSO) (Tibshirani, 1996a), spline-penalized regression (Marx & Eilers, 1996), tree-based models like Random Forests (Breiman, 2001), and gradient-boosted methods such as XGBoost (Chen & Guestrin, 2016), all of which provide model-aware importance scores but may sacrifice interpretability in complex non-linear settings. Wrapper methods, such as Recursive Feature Elimination (RFE) for Support Vector Machines (SVMs) (Guyon et al., 2002), directly search for subsets of features by iteratively training and evaluating

054 models; these often yield strong predictive performance but are computationally expensive and
055 prone to overfitting. Recent heuristic ensemble strategies such as the Minimum Union (Min) method
056 (Radovic et al., 2017) have also been proposed to combine outputs of multiple selectors, improving
057 stability but adding complexity. Finally, unsupervised projection techniques like Principal Com-
058 ponent Analysis (PCA) (Jolliffe, 2002) reduce dimensionality by constructing orthogonal linear
059 combinations of the original variables, which improves computational efficiency but discards di-
060 rect interpretability and neglects target information. Overall, while traditional methods offer a rich
061 toolbox for feature selection, they involve trade-offs among computational cost, predictive perfor-
062 mance, stability, and interpretability. Filter methods are fast but myopic, wrappers are accurate but
063 expensive, and embedded methods are model-aware but often opaque. These limitations motivate
064 the exploration of new architectures, such as Kolmogorov-Arnold Networks (KANs), which offer
065 both model-awareness and direct interpretability by parameterizing each feature transformation as a
066 trainable spline.

067 Kolmogorov-Arnold Networks (KANs) are a recently proposed neural architecture (Liu et al.,
068 2024b) inspired by the Kolmogorov-Arnold representation theorem, which states that any contin-
069 uous multivariate function can be expressed as a finite superposition of univariate functions and
070 addition. KANs operationalize this idea by replacing traditional scalar weights with trainable spline
071 functions: each input feature is passed through one-dimensional, edge-specific splines, aggregated
072 at hidden nodes, and then transformed by learnable outer splines. This “weights-as-functions”
073 paradigm yields compact yet expressive models and, crucially for feature selection, a structured
074 parameterization in which all parameters associated with a given feature form an explicit block. In-
075 tuitively, if the learned splines attached to feature j are nearly flat, then the network’s output varies
076 little with respect to that coordinate, whereas large or highly curved splines signal a strong, nonlinear
077 dependence. KANs therefore provide a natural basis for defining feature-importance scores directly
078 from a trained model.

078 Despite this appealing structure, existing work on KANs for feature selection remains sparse and
079 does not fully exploit these properties. Zheng et al. (2025) use KAN as a surrogate fitness evaluator
080 within a Whale Optimization Algorithm for high-dimensional medical data, inheriting the computa-
081 tional cost and sensitivity of meta-heuristic search and treating KAN largely as a black-box scorer.
082 Other preliminary studies rely on manual spline inspection to prune features (Wang et al., 2025) or
083 use KAN as a front-end encoder for IMU time series (Liu et al., 2024a), without turning its spline
084 parameters and gradients into systematic importance measures.

085 In this work, we close this gap by introducing a family of KAN-based feature-importance criteria
086 that explicitly leverage the weights-as-functions view. We derive coefficient-based norms of spline
087 parameters (*KAN-L1*, *KAN-L2*) to capture the global magnitude of each feature’s learned univariate
088 transformation, a sensitivity integral based on input gradients (*KAN-SI*) to quantify local influence,
089 and a knock-out score (*KAN-KO*) that measures the increase in loss when a feature’s spline block is
090 ablated. We then systematically evaluate these KAN-based selectors as supervised feature-selection
091 methods on tabular classification and regression benchmarks, comparing them with representative
092 filter, wrapper, and embedded baselines. Our analysis covers predictive performance (macro- F_1
093 and R^2), robustness to redundancy (average pairwise correlation among selected features), stability
094 (Jaccard similarity of selected sets across the five cross-validation folds), and interpretability (plots
095 of KAN layer responses and class logits for top-ranked features)¹. Empirically, *KAN-L2*, *KAN-*
096 *SI*, and *KAN-KO* are competitive with or superior to classical baselines on structured and synthetic
097 datasets, while remaining robust on noisy real-world tasks; *KAN-L1* can be highly effective in some
098 classification settings but tends to over-prune in regression. Overall, our results indicate that KAN-
099 based selectors provide a practical and interpretable alternative to traditional methods.

100
101 The paper is organized as follows. Section 2 surveys the related work on feature selection meth-
102 ods. Section 3 provides background and introduces Kolmogorov-Arnold Networks (KANs). Sec-
103 tion 4 presents the methodology: coefficient-based feature importance (*KAN-L1*, *KAN-L2*), *KAN-*
104 *Knockout* (*KAN-KO*), and *KAN-Sensitivity Integral* (*KAN-SI*), and it specifies the assessment
105 pipeline from feature selection to prediction (predictors, evaluation metrics, and leakage-safe cross-
106

107
¹See appendix 6

validation) together with the baseline methods. Section 5 describes the dataset, experimental results and analysis along with a discussion of different results.

2 RELATED WORK

In supervised learning with tabular data, feature selection is a critical step for improving predictive performance, enhancing interpretability, and reducing computational cost. High-dimensional datasets, common in domains such as finance, bioinformatics, and environmental monitoring often contain irrelevant or redundant variables that can degrade model accuracy and increase overfitting risk. Over time, a variety of supervised feature selection techniques have been developed.

One widely used embedded method is the Least Absolute Shrinkage and Selection Operator (LASSO) (Tibshirani, 1996b), which adds an ℓ_1 penalty to shrink small coefficients to zero, yielding sparse and interpretable linear models at modest computational cost. However, in the presence of strong collinearity it typically keeps only one predictor from a correlated group, potentially discarding relevant variables (Zou & Hastie, 2005). Generalized Additive Models (GAMs) (Wood & Augustin, 2002) extend linear models with penalized splines to capture nonlinear effects while retaining some interpretability, but fitting and selecting splines becomes computationally demanding as dimensionality grows. Tree-based ensembles such as Random Forests (Breiman, 2001) provide impurity-based importance scores and naturally capture interactions, yet their importance can be biased toward high-variance or high-cardinality features and need not transfer well to other model classes (Strobl et al., 2007). Wrapper methods like Recursive Feature Elimination (RFE) (Guyon et al., 2002) iteratively retrain a model while pruning low-importance features, often achieving strong accuracy but at significant computational cost. Finally, information-theoretic filter methods such as Mutual Information (MI) ranking (Peng et al., 2005) can detect nonlinear dependencies without repeated model training, but because they score features individually, they do not account for redundancy and tend to select correlated variables with overlapping information.

Extreme Gradient Boosting (XGBoost) (Chen & Guestrin, 2016) is a widely used tree-ensemble method noted for its efficiency and built-in feature importance metrics. During training, XGBoost performs embedded feature selection by evaluating candidate splits and produces importance scores (e.g., gain or split counts) that can be used to rank and filter variables; this has been exploited in numerous applications, such as reducing a 42-dimensional intrusion-detection dataset to 19 informative features with improved accuracy (Kasongo & Sun, 2020). SHAP (SHapley Additive exPlanations) (Lundberg & Lee, 2017) extends this idea by providing theoretically grounded, instance-level attributions whose magnitudes can be aggregated for feature ranking, often yielding more interpretable importance profiles and supporting domain-facing explanations (e.g., in medical risk models). Empirical studies, however, highlight trade-offs: Wang et al. (2024) report that simple XGBoost importance-based selection can outperform SHAP-based selection in both AUPRC and computational efficiency on credit fraud data, reflecting the additional overhead of SHAP computation.

Minimum Union (Min) Method (Seijo-Pardo et al., 2019) aggregates multiple feature ranking results to improve stability in high-dimensional tabular data. The Min method computes each feature’s best (lowest) rank across an ensemble of selectors, yielding a combined ranking that favors features identified as important by any selector (Seijo-Pardo et al., 2016). However, the evaluation of Min has largely been confined to such high-dimensional biological data, and it has not been extensively benchmarked on diverse datasets or against a wide range of modern feature selection methods. Thus, while Min can improve selection stability and performance in some scenarios, its general efficacy across other domains and in combination with different learning algorithms remains an open question in the literature.

Beyond the original KAN formulation (Liu et al., 2024b; Akazan et al., 2025), recent variants such as GKAN lifts KANs to graph-structured domains, learning node/graph functions via Kolmogorov-Arnold compositions on topology-aware inputs (Kiamari et al., 2024). In biomedicine, interpretable Graph KANs fuse multi-omics with graph priors to achieve multi-cancer classification and biomarker discovery, emphasizing model transparency through graph-aware architectures (Alharbi et al., 2025). Complementarily, KAN-guided Whale Optimization uses KAN outputs to steer a metaheuristic for feature selection on medical datasets, reporting empirical gains but relying on heuristic search rather than first-principles criteria (Zheng et al., 2025). Our contribution is orthog-

onal: we introduce a generic KAN feature-importance framework (KAN-SI,KO, L1 and L2) that derives importance from explicit chain-rule sensitivities and spline energy, treats categoricals in a reference-invariant manner, and separates relevance from redundancy through a derivative-Gram audit enabling diversity-aware selection.

3 BACKGROUND

This part provides a detailed analysis of the vanilla KAN, which we used for our study.

3.1 KOLMOGOROV-ARNOLD NETWORKS

KANs (Liu et al., 2024b; Akazan et al., 2025) are inspired by Kolmogorov’s superposition theorem, which states that any multivariate continuous function $f(x_1, \dots, x_n)$ can be represented as a finite composition of continuous univariate functions. KAN models $f : \mathbb{R}^n \rightarrow \mathbb{R}$ as:

$$f(\mathbf{x}) = \sum_{j=1}^{2n+1} \phi_j \left(\sum_{i=1}^n \psi_{ij}(x_i) \right), \quad (1)$$

where $\psi_{i,j} : [0, 1] \rightarrow \mathbb{R}$ are the univariate functions (B-splines in this case) $\phi_j : \mathbb{R} \rightarrow \mathbb{R}$ are continuous functions.

Each input feature x_i is passed through a spline basis transformation:

$$\psi_{ij}(x_i) = \sum_{k=1}^K w_{ijk} B_k(x_i), \quad (2)$$

where B_k is a fixed spline basis (e.g., B-splines or sinusoidal), and w_{ijk} are learnable coefficients. These coefficients are learned via backpropagation.

KAN layer (matrix form). Let $x \in \mathbb{R}^d$ be an input vector and m the number of outputs. A KAN layer produces

$$y = W_{\text{base}} \phi(x) + \sum_{j=1}^d W_{\text{spline}}^{(j)} b_j(x_j) \in \mathbb{R}^m, \quad (3)$$

where: $W_{\text{base}} \in \mathbb{R}^{m \times d}$ is the base (linear) weight matrix; $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is applied componentwise (e.g. SiLU), so $(\phi(x))_j = \phi(x_j)$; $b_j(x_j) \in \mathbb{R}^K$ is the K -dimensional B-spline basis vector for feature j ; $W_{\text{spline}}^{(j)} \in \mathbb{R}^{m \times K}$ are the spline weights attached to feature j .

Stacking all spline bases into a single vector

$$b(x) = [b_1(x_1)^\top b_2(x_2)^\top \dots b_d(x_d)^\top]^\top \in \mathbb{R}^{d \times K}, \quad (4)$$

and concatenating weights

$$W_{\text{spline}} = [W_{\text{spline}}^{(1)} W_{\text{spline}}^{(2)} \dots W_{\text{spline}}^{(d)}] \in \mathbb{R}^{m \times d \times K}, \quad (5)$$

the layer is simply

$$y = W_{\text{base}} \phi(x) + W_{\text{spline}} b(x). \quad (6)$$

Batch form. For $X \in \mathbb{R}^{n \times d}$ (rows are samples), define $\Phi(X) \in \mathbb{R}^{n \times d}$ by $(\Phi(X))_{ij} = \phi(X_{ij})$, and $B(X) \in \mathbb{R}^{n \times d \times K}$ by concatenating $b_j(X_{:,j})$ columnwise. Then the output matrix $Y \in \mathbb{R}^{n \times m}$ is

$$Y = \Phi(X) W_{\text{base}}^\top + B(X) W_{\text{spline}}^\top. \quad (7)$$

4 METHODOLOGY

In a KAN, each input coordinate is passed through a small set of one-dimensional spline functions and then combined linearly, so all parameters associated with a given feature form an explicit, low-dimensional block. This structure naturally supports several ways of quantifying how much each input dimension contributes to the learned mapping. We therefore define four KAN-based feature-importance criteria derived from a trained model. This section discusses our four KAN-based selectors.²

4.1 COEFFICIENT-BASED FEATURE IMPORTANCE (KAN-L1 AND KAN-L2)

In a KAN layer, each input feature x_i is expanded into a set of K B-spline basis functions,

$$z_i = \sum_{k=1}^K w_{ik} B_k(x_i), \quad (8)$$

where w_{ik} are the learned spline coefficients for feature x_i . These coefficients capture how strongly the model relies on different local regions of x_i 's domain.

To quantify feature importance for x_i , we aggregate each coefficient vector $\mathbf{w}_i = (w_{i1}, \dots, w_{iK})$ using its ℓ_1 and ℓ_2 norms:

$$I_{L_1}(x_i) = \|\mathbf{w}_i\|_1 = \sum_{k=1}^K |w_{ik}|, \quad I_{L_2}(x_i) = \|\mathbf{w}_i\|_2 = \left(\sum_{k=1}^K w_{ik}^2 \right)^{1/2}. \quad (9)$$

To enable comparability across features, we normalize these scores to create the final importance scores:

$$\tilde{I}_{L_1}(x_i) = \frac{I_{L_1}(x_i)}{\sum_{j=1}^d I_{L_1}(x_j)}, \quad \tilde{I}_{L_2}(x_i) = \frac{I_{L_2}(x_i)}{\sum_{j=1}^d I_{L_2}(x_j)}. \quad (10)$$

These normalized quantities provide a direct and interpretable measure of how much each input contributes, on average, through its spline expansion.

4.2 KAN-KNOCKOUT (KAN-KO)

Let the trained KAN layer (see Subsection 3.1) be parameterized by $W := (W_{\text{base}}, W_{\text{spline}})$. For a given feature index $j \in \{1, \dots, d\}$, define the knockout operator \mathcal{K}_j acting on the first KAN layer by

$$\mathcal{K}_j(W) := (\widetilde{W}_{\text{base}}, \widetilde{W}_{\text{spline}}), \quad \text{where } \widetilde{W}_{\text{base}}(:, j) = 0, \widetilde{W}_{\text{spline}}^{(j)} = 0. \quad (11)$$

Equivalently, at the level of the layer output

$$y^{(-j)}(x) = y(x) - W_{\text{base}}(:, j) \phi(x_j) - W_{\text{spline}}^{(j)} b_j(x_j). \quad (12)$$

To quantify the contribution of each input feature, we evaluate how much the task loss increases when that feature is removed from the model. Let $\ell(\hat{y}, y)$ denote the task loss and P the data distribution over (x, y) . We write f_W for the full KAN model, where W collects the base and spline weights (see Subsection 3.1). The expected or population risk of the model is

$$L(W) = \mathbb{E}_{(x,y) \sim P} [\ell(f_W(x), y)]. \quad (13)$$

For a given feature index j , we define a knock-out operator $\mathcal{K}_j(W)$ that sets to zero the j th column of W_{base} and the j th spline block of W_{spline} . The corresponding risk when feature j is removed is

$$L_j(W) = \mathbb{E}_{(x,y) \sim P} [\ell(f_{\mathcal{K}_j(W)}(x), y)]. \quad (14)$$

The *KAN-KO importance score* for feature j is the nonnegative increase in risk,

$$\Delta_j = \max\{0, L_j(W) - L(W)\}. \quad (15)$$

To make scores comparable across features, we normalize them as

$$I_j = \frac{\Delta_j}{\sum_{k=1}^d \Delta_k + \delta}, \quad (16)$$

where $\delta > 0$ is a small constant to avoid division by zero.

²See the appendix A for more details

4.3 KAN- SENSITIVITY INTEGRAL (KAN-SI) FEATURE IMPORTANCE

For an input feature x_i , the instantaneous (local) sensitivity of f is the magnitude of its partial derivative,

$$S_i(\mathbf{x}) = \left| \frac{\partial f(\mathbf{x})}{\partial x_i} \right|. \quad (17)$$

The KAN-SI importance is the (data) expectation of S_i :

$$I_i = \mathbb{E}_{\mathbf{x}}[S_i(\mathbf{x})], \quad \hat{I}_i = \frac{1}{N} \sum_{n=1}^N \left| \frac{\partial f(\mathbf{x}^n)}{\partial x_i} \right| \quad (18)$$

where \hat{I}_i is the empirical estimator over a held-out set $\{\mathbf{x}^n\}_{n=1}^N$ (validation or out-of-fold) to avoid optimistic bias.

From equation 1, let the inner pre-activations be

$$t_j(\mathbf{x}) = \sum_{i=1}^n \psi_{ij}(x_i).$$

Then by the chain rule,

$$\frac{\partial f(\mathbf{x})}{\partial x_i} = \sum_{j=1}^{2n+1} \phi'_j(t_j(\mathbf{x})) \frac{\partial \psi_{ij}(x_i)}{\partial x_i}. \quad (19)$$

Because ψ_{ij} is a spline expansion,

$$\frac{\partial \psi_{ij}(x_i)}{\partial x_i} = \sum_{k=1}^K w_{ijk} B'_k(x_i), \quad (20)$$

and for B-splines of degree p , only $p + 1$ basis derivatives $B'_k(x_i)$ are non-zero at a given x_i (local support), which makes equation 19 efficient to evaluate.

Since the sensitivities depend on the measurement of x_i , we report a normalized score.

$$\tilde{I}_i = s_i \hat{I}_i, \quad s_i \in \{\text{Std}(x_i), \text{IQR}(x_i)\} \quad (21)$$

where Std(standart deviation) or IQR(interquartile range) is computed on the same split used for equation 18. For one-hot encoded categoricals, compute \tilde{I} per dummy and *sum* back to the original category. KAN-SI quantifies the importance of feature x_i as the expected on-data magnitude of the model's directional derivative given by equation 18 and report the scale-invariant score using equation 21 to neutralize unit effects.

4.4 ASSESSMENT PROCEDURE (FROM FEATURE SELECTION TO PREDICTION)

Let $D = \{(x_i, y_i)\}_{i=1}^n$, where $x_i \in \mathbb{R}^d$ and

$$y_i \in \begin{cases} \{1, \dots, C\}, & \text{classification,} \\ \mathbb{R}, & \text{regression.} \end{cases} \quad (22)$$

Let $\mathcal{S} = \{s_1, \dots, s_L\}$ denote feature selectors (e.g., KAN-L1/L2, KAN-SI/KO, MI, LASSO, etc.). Each $s \in \mathcal{S}$ fitted on training data returns a nonnegative importance vector $a_s \in \mathbb{R}_+^d$, normalized so $\sum_{j=1}^d a_{s,j} = 1$. For a retention ratio $k \in \{20\%, 40\%, 60\%\}$, set

$$n_k = \max \left\{ 1, \left\lceil \frac{k \cdot d}{100} \right\rceil \right\}, \text{ number equivalent of } k\% \text{ of the features} \quad (23)$$

Let $J_{s,k} \subset \{1, \dots, d\}$ be the indices of the n_k largest entries of a_s (ties arbitrary), and define the *projection* $\Pi_J(x) \in \mathbb{R}^{|J|}$ that restrict x to the matrix defined by the coordinates J .

4.4.1 PREDICTORS AND EVALUATION METRICS.

For each task $t \in \{\text{classification, regression}\}$, we consider a family of predictors $\mathcal{P}_t = \{P_{t,1}, \dots, P_{t,M}\}$, including Logistic Regression / Ridge, Random Forests, Gradient Boosted Trees (GBT), and XGBoost. Predictive performance is quantified by

$$S_t(y, \hat{y}) = \begin{cases} \text{Macro-}F_1(y, \hat{y}), & t = \text{classification}, \\ R^2(y, \hat{y}), & t = \text{regression}, \end{cases} \quad (24)$$

where the per-class F_1 score is defined one-vs-rest and averaged across C classes:

$$F1_c = \frac{2TP_c}{2TP_c + FP_c + FN_c}, \quad \text{Macro-}F_1 = \frac{1}{C} \sum_{c=1}^C F1_c. \quad (25)$$

4.4.2 LEAKAGE-SAFE CROSS-VALIDATION.

We evaluate each selector-predictor pair using leakage-safe F -fold cross-validation. For each fold f , we fit preprocessing and the feature selector s only on the training split T_f , obtain the selected feature set $J_{s,k}^{(f)}$, train the predictor on the projected training data, and then compute the score $S_{s,k,t,m,f}$ (macro- F_1 or R^2) on the projected validation split V_f . The overall performance $\text{Score}(s, k, t, m)$ is the average of these fold-level scores. This procedure ensures that both feature selection and model training see only training data in each fold, preventing information leakage from the validation sets and yielding an unbiased estimate of generalization performance.³

Comparison with Baseline Methods We compare our spline-based feature selectors against Mutual Information (MI) ranking (Peng et al., 2005), Random Forest (Breiman, 2001), LASSO-based feature selection (Tibshirani, 1996b), and Recursive Feature Elimination (RFE) for Support Vector Machine (SVM) (SVM-RFE)(Guyon et al., 2002).

5 EXPERIMENTS

Reproducibility details and data sets information can be found at Appendix B.1 and Appendix B.2.

The comparative analysis across datasets reveals that averaging F1 and R^2 scores across retention levels k reveals consistent yet dataset-specific interactions between selectors, predictors, and the underlying data structure.

Classification datasets On the *Breast Cancer* dataset (Figure 1), Gradient Boosted Trees perform best with *KAN-L1*, LASSO, and RF importance, often matching the full-feature baseline. In contrast, *KAN-KO* and Mutual Information yield poor subsets, while Logistic Regression and tree ensembles (RF, XGBoost) benefit most from sparsity- or impurity-based selectors. This indicates that selectors capturing either linear sparsity or tree-based splits are most effective, whereas knock-out and sensitivity KAN variants are less suited. In the *Digits* dataset (Figure 1), the pattern shifts: Gradient Boosted Trees and XGBoost achieve higher F1 scores with LASSO, SVM-RFE, and several KAN variants (*KAN-KO*, *KAN-SI*, *KAN-L1*), sometimes surpassing the full-feature baseline. Here, dimensionality reduction removes redundancy and sharpens predictive power. By contrast, *KAN-L2* and Mutual Information underperform, reflecting their limits in capturing complex feature interactions. For the synthetic *make_classification* dataset (Figure 1), all KAN-based selectors (*KAN-L1/L2*, *KAN-SI*, *KAN-KO*) consistently outperform classical baselines, confirming their advantage in structured, high-signal settings where ground-truth feature relevance extends beyond sparsity or univariate dependence. Finally, on the *Wine* dataset (Figure 1), the advantage shifts back toward classical selectors: RF importance, LASSO, Mutual Information, and SVM-RFE provide the strongest subsets, while most KAN-based selectors lag behind. Only *KAN-SI* remains competitive, suggesting that when features are few and relationships are well-structured, simpler selectors aligned with linear sparsity or impurity-based measures are more effective.

³See the mathematical formulation of this process in appendix B.3

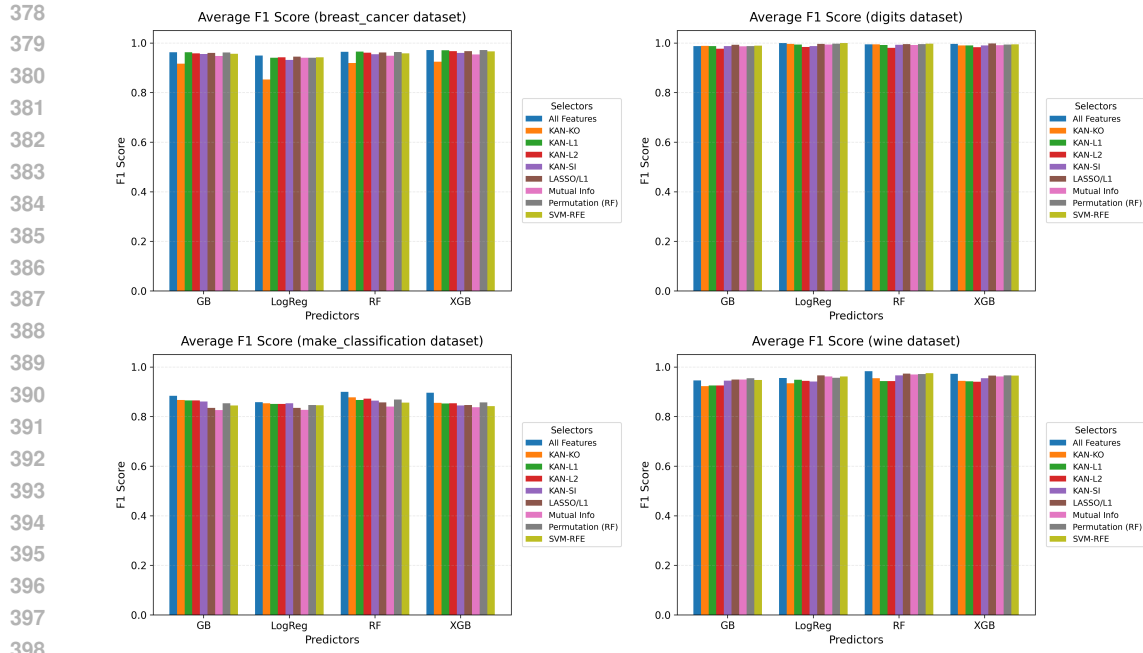


Figure 1: Average $F1$ Score (averaged across 20/40/60% retained features) per Selectors for Each Classifiers

Regression datasets On the *California dataset (nonlinear, heterogeneous)*, Figure. 2, *KAN-L2* and *KAN-SI* preserve tree-ensemble performance: with Random Forest and XGBoost they stay close to the prediction using all Features (minor drops $\sim 0.01-0.03$), while Gradient Boosted Tree is slightly more sensitive. In contrast, *KAN-L1* is overly aggressive: it depresses Ridge markedly

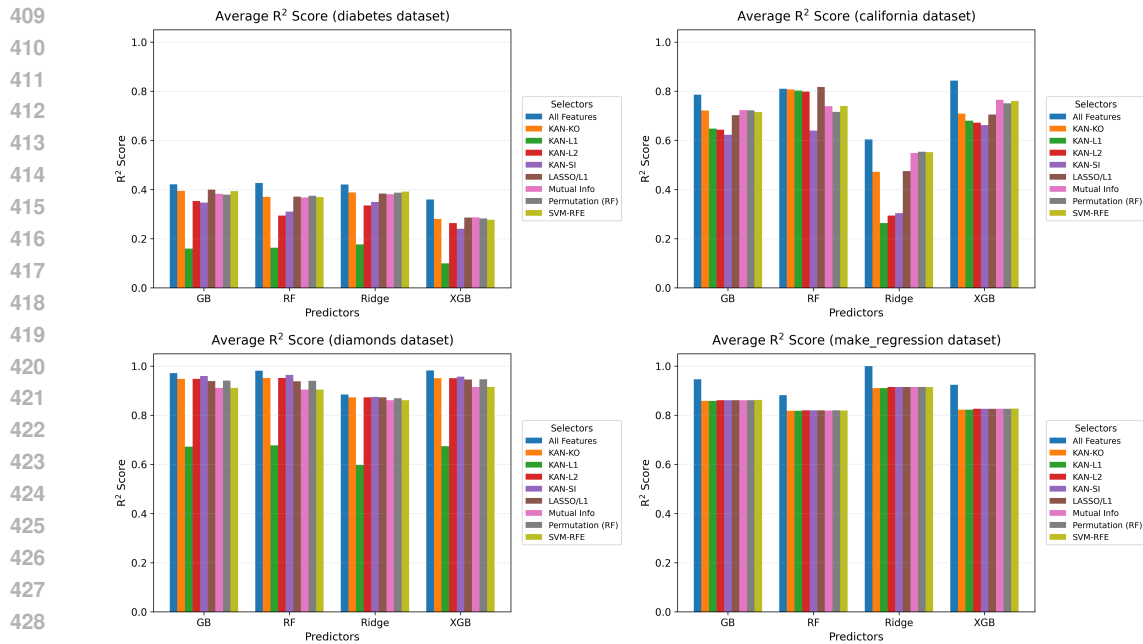


Figure 2: Relative Average R^2 Score (averaged across 20/40/60% retained features) per Selectors for Each Regressors

(large R^2 gap vs. All Features) and also trims XGBoost and Gradient Boosted Tree more than other KAN variants, indicating loss of weak but complementary signals. *KAN-KO* feature selector is generally competitive, trailing All Features by a small margin and often matching the classical Mutual Info, Random Forest and SVM-RFE subsets. For *Diabetes dataset (small-n, noisy)*, Figure. 2), *KAN-L1* again under-selects (lowest bars for Gradient Boosted Tree, Random Forest, XGBoost), whereas *KAN-L2, KAN-SI* yield the most stable KAN performance across predictors and typically track LASSO and Mutual Info. *KAN-KO* behaves as a conservative KAN subset, usually second-best among KANs. In the *Diamond dataset (strong signal, structured, mixed types)*, Figure. 2, All models are high- R^2 ; here *KAN-KO, KAN-L2* and *KAN-SI* are nearly indistinguishable from All Features for Random Forest, XGBoost and Gradient Boosted Tree and even outperform classical features selectors (Mutual Info, Random Forest and SVM-RFE subsets, LASSO and Random Forest). *KAN-L1* is the only KAN variant that noticeably drops on some predictors, consistent with over-pruning informative but correlated attributes. For the Synthetic *make regression dataset (well-conditioned)*, Figure. 2, All KAN variants cluster tightly with baselines across predictors (tiny spreads): when the signal is clean and features are already informative, KAN selection neither helps nor hurts much; *KAN-L2/KAN-SI* remain the safest choices.

Taken together, these findings suggest that KAN-based selectors are particularly attractive when feature interactions and nonlinearity play a central role, or when one wishes to couple selection and interpretability (see Appendix C). Their main failure mode is over-pruning correlated yet informative features under an ℓ_1 -style criterion (*KAN-L1*), in regression. More broadly, the results reinforce that dimensionality reduction is not uniformly harmful: when a selector is aligned with the structure of the data and the bias of the predictor, removing redundant or spurious variables can improve both accuracy and interpretability.

6 CONCLUSION

This work presented, to our knowledge, the first systematic study of Kolmogorov-Arnold Network (KAN) based feature selection on tabular classification and regression benchmarks. We defined four KAN-derived criteria (*KAN-L1*, *KAN-L2*, *KAN-SI*, *KAN-KO*) and compared them against widely used baselines, including LASSO, Random Forest feature importance, mutual information, and SVM-RFE, across multiple datasets and feature-retention levels. Empirically, *KAN-L2*, *KAN-SI*, and *KAN-KO* provide competitive and often superior performance in structured or strong-signal settings, while remaining robust on real-world tasks; *KAN-L1* can be effective in classification but tends to over-prune in noisy or correlated regression problems. Classical selectors such as LASSO and Random Forest remain strong choices, particularly when they align with the assumptions of the downstream model. Furthermore, our stability and redundancy analyses show that KAN-based criteria produce reproducible feature subsets across folds and do not inflate correlation among selected features indicating that they offer both reliable and non-redundant selection even without enforcing sparsity constraints. Beyond aggregate scores, our spline-based case studies show that KANs can yield smooth, one-dimensional response functions that link feature values to model behaviour in a transparent way, offering an interpretable view of nonlinear and multivariate relevance that is difficult to obtain from purely sparsity- or impurity-based methods. Overall, our findings indicate that KAN-based feature selection is a practical, interpretable alternative to traditional approaches, and they provide guidance on when specific KAN criteria are most beneficial (e.g., *KAN-L2/KAN-SI* for noisy regression, *KAN-SI/KAN-KO* for interaction-heavy classification). We see this as a step toward feature-selection methods that jointly offer strong predictive performance, robustness to irrelevant variables, and clear mechanistic insight into how individual predictors influence model outputs. At the same time, training KANs is noticeably slower than fitting sparse linear models or tree ensembles, which limits their use in large-scale screening or AutoML settings. An important direction for future work is therefore to develop faster, yet still interpretable, KAN architectures and training schemes (e.g., lightweight KAN layers, pruning or distillation of spline components, or hybrid models that reuse KAN-based scores to guide simpler selectors).

REFERENCES

Ange-Clement Akazan, Verlon Roel Mbingui, Gnankan Landry Regis N’guessan, and Issa Karambal. Localized weather prediction using kolmogorov-arnold network-based models and deep rnns,

- 486 2025. URL <https://arxiv.org/abs/2505.22686>.
487
- 488 Fadi Alharbi, Nishant Budhiraja, Aleksandar Vakanski, Boyu Zhang, Murtada K Elbashir, Harshith
489 Guduru, and Mohanad Mohammed. Interpretable graph kolmogorov–arnold networks for multi-
490 cancer classification and biomarker identification using multi-omics data. *Scientific Reports*, 15
491 (1):27607, 2025.
- 492 Leo Breiman. Random forests. *Machine learning*, 45:5–32, 2001.
493
- 494 Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the*
495 *22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD*
496 *'16)*, pp. 785–794. ACM, 2016. doi: 10.1145/2939672.2939785.
- 497 Bradley Efron, Trevor Hastie, Iain Johnstone, and Robert Tibshirani. Least angle regression. *Annals*
498 *of Statistics*, 32(2):407–499, 2004. doi: 10.1214/009053604000000067.
499
- 500 M. Forina et al. Uci machine learning repository: Wine recognition dataset. [https://archive.](https://archive.ics.uci.edu/ml/datasets/wine)
501 [ics.uci.edu/ml/datasets/wine](https://archive.ics.uci.edu/ml/datasets/wine), 1991.
502
- 503 Isabelle Guyon, Jason Weston, Stephen Barnhill, and Vladimir Vapnik. Gene selection for cancer
504 classification using support vector machines. *Machine learning*, 46(1):389–422, 2002.
- 505 Ian Jolliffe. *Principal component analysis*. Springer, 2002.
506
- 507 Sydney M Kasongo and Yanxia Sun. Performance analysis of intrusion detection systems using a
508 feature selection method on the unsw-nb15 dataset. *Journal of Big Data*, 7(1):105, 2020.
- 509 Mehrdad Kiamari, Mohammad Kiamari, and Bhaskar Krishnamachari. Gkan: Graph kolmogorov-
510 arnold networks. *arXiv preprint arXiv:2406.06470*, 2024.
511
- 512 Mengxi Liu, Daniel Geißler, Dominique Nshimiyimana, Sizhen Bian, Bo Zhou, and Paul Lukowicz.
513 Initial investigation of kolmogorov-arnold networks (kans) as feature extractors for imu based
514 human activity recognition, 2024a. URL <https://arxiv.org/abs/2406.11914>.
- 515 Ziming Liu, Yixuan Wang, Sachin Vaidya, Fabian Ruehle, James Halverson, Marin Soljačić,
516 Thomas Y Hou, and Max Tegmark. Kan: Kolmogorov-arnold networks. *arXiv preprint*
517 *arXiv:2404.19756*, 2024b.
518
- 519 Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Proceed-*
520 *ings of the 31st International Conference on Neural Information Processing Systems, NIPS'17*,
521 pp. 4768–4777, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- 522 Brian D. Marx and Paul H. C. Eilers. Multivariate adaptive regression splines. *Annals of Statistics*,
523 24(1):89–120, 1996.
524
- 525 Nicolai Meinshausen and Peter Bühlmann. Stability selection. *Journal of the Royal Statistical*
526 *Society Series B: Statistical Methodology*, 72(4):417–473, 2010.
- 527 Sarah Nogueira, Konstantinos Sechidis, and Gavin Brown. On the stability of feature selection
528 algorithms. *Journal of Machine Learning Research*, 18(174):1–54, 2018.
529
- 530 Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier
531 Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas,
532 Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duch-
533 esnay. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12:
534 2825–2830, 2011.
- 535 Hanchuan Peng, Fuhui Long, and Chris Ding. Feature selection based on mutual information criteria
536 of max-dependency, max-relevance, and min-redundancy. *IEEE Transactions on pattern analysis*
537 *and machine intelligence*, 27(8):1226–1238, 2005.
538
- 539 M. Radovic, M. Ghalwash, N. Filipovic, and Z. Obradovic. Minimum redundancy feature selection
from microarray gene expression data. *Pattern Recognition*, 61:12–24, 2017.

- 540 Borja Seijo-Pardo, Verónica Bolón-Canedo, and Amparo Alonso-Betanzos. Using data complexity
541 measures for thresholding in feature selection rankers. In *Conference of the Spanish association*
542 *for artificial intelligence*, pp. 121–131. Springer, 2016.
- 543 Borja Seijo-Pardo, Verónica Bolón-Canedo, and Amparo Alonso-Betanzos. On developing an auto-
544 matic threshold applied to feature selection ensembles. *Information Fusion*, 45:227–245, 2019.
- 546 W. N. Street, W. H. Wolberg, and O. L. Mangasarian. Nuclear feature extraction for breast tu-
547 mor diagnosis. In *Biomedical Image Processing and Biomedical Visualization*, volume 1905 of
548 *Proceedings of SPIE*, pp. 861–870. SPIE, 1993. doi: 10.1117/12.148698.
- 549 Carolin Strobl, Anne-Laure Boulesteix, Achim Zeileis, and Torsten Hothorn. Bias in random forest
550 variable importance measures: Illustrations, sources and a solution. *BMC Bioinformatics*, 8(1):
551 25, 2007. doi: 10.1186/1471-2105-8-25.
- 552 Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical*
553 *Society: Series B (Methodological)*, 58(1):267–288, 1996a.
- 554 Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical*
555 *Society Series B: Statistical Methodology*, 58(1):267–288, 1996b.
- 556 Huanjing Wang, Qianxin Liang, John T Hancock, and Taghi M Khoshgoftaar. Feature selection
557 strategies: a comparative analysis of shap-value and importance-based methods. *Journal of Big*
558 *Data*, 11(1):44, 2024.
- 559 Shuaibo Wang, Wenhao Luo, Sixing Yin, Jie Zhang, Zhuohang Liang, Yihua Zhu, and Shufang
560 Li. Interpretable state estimation in power systems based on the kolmogorov–arnold networks.
561 *Electronics*, 14(2), 2025. ISSN 2079-9292. doi: 10.3390/electronics14020320. URL <https://www.mdpi.com/2079-9292/14/2/320>.
- 562 Hadley Wickham. Diamonds dataset. <https://www.kaggle.com/datasets/shivam2503/diamonds>, 2008. Accessed: 2025-09-24.
- 563 Simon N Wood and Nicole H Augustin. Gams with integrated model selection using penalized
564 regression splines and applications to environmental modelling. *Ecological modelling*, 157(2-3):
565 157–177, 2002.
- 566 Boli Zheng, Yi Chen, Chaofan Wang, Ali Asghar Heidari, Lei Liu, Huiling Chen, and Guoxi Liang.
567 Kolmogorov-arnold networks guided whale optimization algorithm for feature selection in medi-
568 cal datasets. *Journal of Big Data*, 12(1):69, 2025.
- 569 Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the*
570 *Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005. doi: 10.1111/
571 j.1467-9868.2005.00503.x.
- 572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
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

Part I**Table of Contents**

A Mathematical Motivation of KAN-Based Feature Scores	13
Mathematical Motivation of KAN-Based Feature Scores	13
A.1 Functional view of the first KAN layer	13
A.2 Coefficient norms: <i>KAN-L1</i> and <i>KAN-L2</i>	13
A.3 Sensitivity integral: <i>KAN-SI</i>	13
A.4 Knock-out loss: <i>KAN-KO</i>	14
A.5 Summary	14
B Reproducibility	14
Reproducibility	14
B.1 Reproducibility Details	14
Reproducibility Details	14
B.2 Data Sets	14
Data Sets	14
B.3 Leakage-safe cross-validation	15
Leakage-safe cross-validation	15
C KAN-Based Selection Interpretability	15
KAN-Based Selection Interpretability	15
C.1 Mean Concave Points Interpretability Study	15
Mean Concave Points Interpretability Study	15
C.2 Worst Concave Points Interpretability Study	16
Worst Concave Points Interpretability Study	16
C.3 Radius Error Points Interpretability Study	16
Radius Error Points Interpretability Study	16
D Stability and Redundancy Analysis	17
Stability and Redundancy Analysis	17
D.1 Stability	17
Stability	17
D.2 Redundancy	18
Redundancy	18
E Prediction Performances per feature retention levels	19
Accuracy per feature retention levels	19
E.1 Average predictors performances at different retention levels	19
Average Performances at Different Retention Levels	19
E.2 Insight at 60% retention Level	19
Classification Datasets	19
Regression Datasets	21
F Runtime Profiling of Feature Selectors	22
Runtime Profiling of Feature Selectors	22

A MATHEMATICAL MOTIVATION OF KAN-BASED FEATURE SCORES

This section gives a brief functional motivation for the four KAN-based feature-importance criteria used in the main paper.

A.1 FUNCTIONAL VIEW OF THE FIRST KAN LAYER

For clarity, consider a single-output KAN with one hidden layer. The first layer applies, for each unit $r = 1, \dots, R$ and feature $j = 1, \dots, d$, a univariate spline

$$g_{j,r}(x_j) = \sum_{k=1}^K \theta_{j,r,k} B_k(x_j), \quad (26)$$

where $\{B_k\}_{k=1}^K$ are fixed B-spline basis functions and $\theta_{j,r,k}$ are learned coefficients. Collecting all coefficients attached to feature j in the first layer gives a parameter block

$$\theta_j = (\theta_{j,r,k})_{r=1,\dots,R; k=1,\dots,K}. \quad (27)$$

The network output can be written schematically as

$$f(x) = \sum_{r=1}^R w_r \sigma\left(\sum_{j=1}^d g_{j,r}(x_j)\right), \quad (28)$$

so that feature j influences f only through its spline family $\{g_{j,r}\}_r$ controlled by θ_j . The scores below quantify the importance of feature j by measuring, in different ways, how much the learned function f depends on this block.

A.2 COEFFICIENT NORMS: *KAN-L1* AND *KAN-L2*

For fixed j and r , define the spline $g_{j,r}(x_j) = \sum_k \theta_{j,r,k} B_k(x_j)$ and consider its $L^2(P_X)$ norm with respect to the data distribution P_X :

$$\|g_{j,r}\|_{L^2(P_X)}^2 = \mathbb{E}_X[g_{j,r}(X_j)^2] = \theta_{j,r}^\top G_j \theta_{j,r}, \quad (29)$$

where $G_j = \mathbb{E}_X[B(X_j)B(X_j)^\top]$ is the Gram matrix of the B-spline basis for feature j . For standardized inputs and regular knots, G_j is well-conditioned, so there exist constants $0 < \lambda_{\min} \leq \lambda_{\max} < \infty$ such that

$$\lambda_{\min} \|\theta_{j,r}\|_2^2 \leq \|g_{j,r}\|_{L^2(P_X)}^2 \leq \lambda_{\max} \|\theta_{j,r}\|_2^2. \quad (30)$$

Thus, up to fixed constants, the ℓ_2 norm of the spline coefficients is proportional to the L^2 “energy” of the learned univariate transformation of feature j . Summing over hidden units r yields that $\|\theta_j\|_2$ is a proxy for the total contribution of feature j to the first-layer representation. This motivates the *KAN-L2* score as

$$\text{score}_{L2}(j) = \|\theta_j\|_2. \quad (31)$$

The *KAN-L1* score $\text{score}_{L1}(j) = \|\theta_j\|_1$ emphasizes sparsity over basis functions: features whose effect is concentrated on a small number of knots receive higher values, encouraging compact, interpretable univariate responses.

A.3 SENSITIVITY INTEGRAL: *KAN-SI*

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be the trained KAN predictor and assume f is differentiable almost everywhere. For a small perturbation δ along coordinate j we have the first-order approximation

$$f(x + \delta e_j) - f(x) \approx \delta \frac{\partial f}{\partial x_j}(x), \quad (32)$$

where e_j is the j -th canonical basis vector. Taking expectations over the data distribution and random perturbations with $\mathbb{E}[|\delta|] = c$ gives

$$\mathbb{E}[|f(X + \delta e_j) - f(X)|] \approx c \mathbb{E}_X\left[\left|\frac{\partial f}{\partial x_j}(X)\right|\right]. \quad (33)$$

702 The quantity

$$703 S_j^{\text{SI}} = \mathbb{E}_X \left[\left| \frac{\partial f}{\partial x_j}(X) \right| \right] \quad (34)$$

704 therefore measures the expected local sensitivity of the prediction to perturbations in feature j . Our
 705 *KAN-SI* score is the empirical estimate of S_j^{SI} computed on a validation set, using the gradient of
 706 the trained KAN. Features with large S_j^{SI} are those for which small changes in x_j typically induce
 707 large changes in the output, making *KAN-SI* a derivative-based global importance measure.
 708
 709

710 A.4 KNOCK-OUT LOSS: *KAN-KO*

711 Let $\ell(f(X), Y)$ denote the training loss (cross-entropy for classification or squared error for regres-
 712 sion) and $R(f) = \mathbb{E}[\ell(f(X), Y)]$ its risk. Given a trained KAN f , we define $f^{(-j)}$ as the same
 713 network but with all spline parameters associated with feature j in the first layer set to zero, so that
 714 feature j no longer contributes to the first-layer representation. The *knock-out* importance of feature
 715 j is then

$$716 \Delta_j = R(f^{(-j)}) - R(f) \geq 0, \quad (35)$$

717 the increase in risk incurred when the contribution of feature j is ablated. In practice, we estimate
 718 Δ_j by the empirical risk difference on a held-out set, using the same trained parameters and only
 719 modifying the first-layer spline block of feature j . When f approximately minimizes R , larger
 720 Δ_j indicate that the model relies heavily on feature j for accurate predictions. This motivates the
 721 *KAN-KO* score as a leave-one-feature-out, loss-based measure of relevance.
 722
 723

724 A.5 SUMMARY

725 In summary, the four KAN-based criteria exploit the spline-based parameterization of the first layer
 726 to approximate complementary notions of feature relevance: *KAN-L2/KAN-L1* measure the “en-
 727 ergy” and sparsity of the learned univariate transformations, *KAN-SI* captures gradient-based global
 728 sensitivity, and *KAN-KO* quantifies the increase in risk when a feature’s contribution is removed. All
 729 are computed directly from a trained KAN, without external wrapper search or additional surrogate
 730 models.
 731
 732

733 B REPRODUCIBILITY

734 This section provides necessary reproducibility details

735 B.1 REPRODUCIBILITY DETAILS

736 We implement all models in PyTorch using a Kolmogorov-Arnold Network (KAN) whose archi-
 737 tecture is defined by the experiment-specific list `layers_hidden=[n_input, 2n + 1, n_output]`. Across
 738 all experiments, the KAN modules rely on cubic spline bases over a grid of five knots on the in-
 739 terval $[-1, 1]$, with a small grid offset of 0.02 to avoid boundary artefacts. The base branch uses
 740 the *SiLU* nonlinearity, while the spline and base components are initialized with unit scaling and a
 741 modest amount of injected noise (scale 0.1) to encourage exploration during early training. Models
 742 are trained for 100 epochs with a mini-batch size of 64 using the Adam optimizer with PyTorch’s
 743 default hyperparameters. All runs are executed on Kaggle’s GPU environment with an NVIDIA
 744 Tesla P100, using fixed random seeds for Python, NumPy, and PyTorch and identical train-test splits
 745 across methods to ensure reproducibility.
 746
 747

748 B.2 DATA SETS

749 To evaluate the effectiveness of the Kolmogorov-Arnold Networks dimensionality reduction methods,
 750 experiments were conducted on a diverse collection of benchmark datasets covering both classifica-
 751 tion and regression tasks. These datasets includes well-established repositories from scikit-learn as
 752 well as synthetically generated datasets and are summarized in Table 1.
 753
 754
 755

Table 1: Benchmark and synthetic datasets

Task	Dataset	Observations	Numbers of features	References
Classification	Wine	178	13	(Forina et al., 1991).
	Breast Cancer	569	30	(Street et al., 1993)
	Wisconsin Digits	11,797	64	(Pedregosa et al., 2011)
	Make classification	500	10	(Pedregosa et al., 2011)
Regression	California Housing	20,640	8	(Efron et al., 2004)
	Diabetes	442	10	(Efron et al., 2004).
	Diamonds	53,940	10	(Wickham, 2008)
	Make regression	500	10	(Pedregosa et al., 2011)

B.3 LEAKAGE-SAFE CROSS-VALIDATION

We use F -fold cross-validation with splits (T_f, V_f) , $f = 1, \dots, F$. For each selector s , retention level k , predictor $P_{t,m}$, and fold f :

1. fit preprocessing and selector s on $\{(x_i, y_i)\}_{i \in T_f}$ to obtain $J_{s,k}^{(f)}$;
2. train $P_{t,m}$ on the projected training set $\{(\Pi_{J_{s,k}^{(f)}}(x_i), y_i)\}_{i \in T_f}$;
3. predict on the projected validation set $\{\Pi_{J_{s,k}^{(f)}}(x_i)\}_{i \in V_f}$ to obtain $\{\hat{y}_i^{(s,k,t,m,f)}\}_{i \in V_f}$.

The fold-level score is

$$S_{s,k,t,m,f} = S_t\left(\{y_i\}_{i \in V_f}, \{\hat{y}_i^{(s,k,t,m,f)}\}_{i \in V_f}\right), \quad (36)$$

where S_t is macro- F_1 for classification and R^2 for regression. The cross-validated score is then

$$\text{Score}(s, k, t, m) = \frac{1}{F} \sum_{f=1}^F S_{s,k,t,m,f}. \quad (37)$$

C KAN-BASED SELECTION INTERPRETABILITY

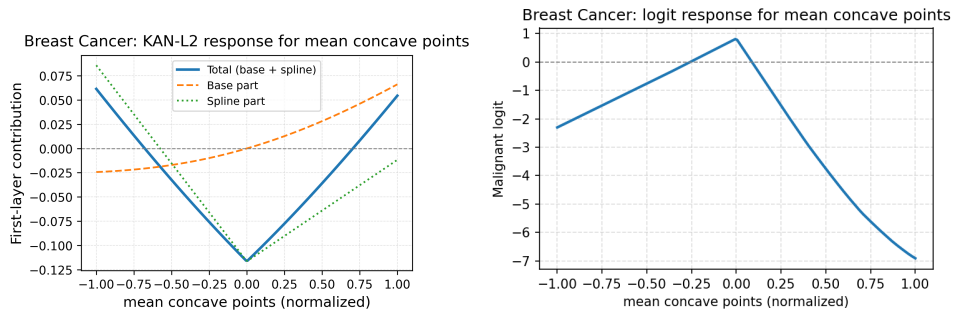
In this study, we visualize the internal KAN response functions $f_j(x_j)$ for the three top features on the Breast Cancer dataset Street et al. (1993), *worst concave points*, *mean concave points*, and *radius error*, as identified by **KAN-L2**, together with the corresponding malignant-class logit. This dataset is based on digitized images of fine-needle aspirates of breast masses, and the prediction task is to distinguish malignant from benign tumours from cell-nucleus morphology. Here, *concave points* quantify how many segments of a nucleus boundary are inwardly curved (concave), so *mean concave points* is the average number of such concave segments across all nuclei in an image, and *worst concave points* is the largest value observed among them. The *radius error* feature measures the variability of the nucleus radius (standard error of the mean radius) across measurements, capturing irregular or heterogeneous tumour shapes. By plotting the learned responses for these features and the resulting malignant-class logit, we can see how changes in concavity and boundary irregularity are transformed inside the network and how they ultimately influence the predicted malignancy score.

C.1 MEAN CONCAVE POINTS INTERPRETABILITY STUDY

For *mean concave points*, the first-layer KAN response in figure 3a is roughly V-shaped. The total contribution (blue curve) is most negative around the centre of the normalized range and rises toward positive values at both low and high extremes. The base term (orange) varies smoothly and almost linearly, while the spline term (green) introduces most of the curvature, indicating that the non-linear spline component is responsible for emphasizing intermediate levels of average concavity in the hidden representation.

When this representation is propagated through the second KAN layer, the malignant-class logit in figure 4b becomes strongly skewed. Starting from low *mean concave points* (very smooth boundaries), the logit is

810 moderately negative (the model leans benign), increases to a clear maximum around the middle of the range
 811 (where the model assigns the highest malignancy risk), and then drops sharply for very high values of *mean*
 812 *concave points*. Thus, in this dataset the KAN regards lesions with intermediate average boundary concavity
 813 as most suspicious, while very smooth or extremely concave boundaries are treated as lower risk. Together,
 814 the two plots illustrate how a V-shaped internal code for this feature is mapped by the final layer into a peaked
 815 malignancy score that is highest at intermediate concavity levels.

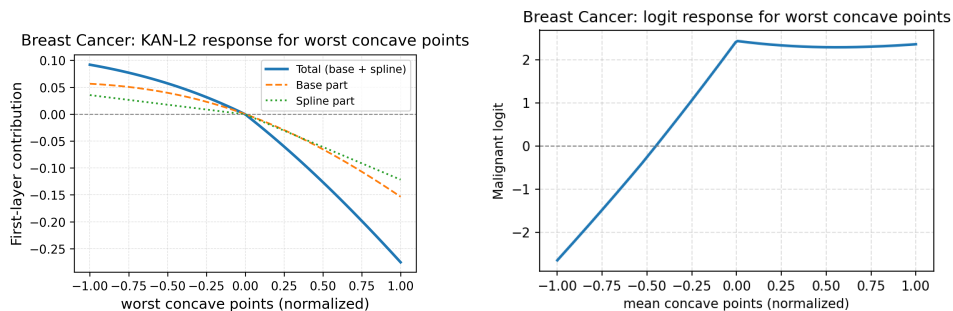


817
818
819
820
821
822
823
824
825
826
827 (a) KAN first-layer responses of *mean concave* (b) Malignant-class logit as a function of *mean*
 828 *points* *concave points* (normalized), obtained by varying
 829 only this feature around a reference input.

830 Figure 3: Breast Cancer dataset interpretability visualizations using mean concave points.

832 C.2 WORST CONCAVE POINTS INTERPRETABILITY STUDY

833
834 For *worst concave points*, the KAN layer response is monotonically decreasing (consistently with the base and
 835 spline variation), so increasing concavity systematically pushes the representation along a single direction, in
 836 line with the clinical view that strongly concave, spiculated margins are characteristic of malignant lesions. To
 837 relate this internal behaviour to the actual prediction, we also examine the malignant-class logit as a function
 838 of *worst concave points* (Figure 4). As worst concavity increases, the malignant logit rises sharply before
 839 saturating at extreme values, showing that the internal direction induced by high concavity is ultimately mapped
 840 to a higher malignancy score. Together, the first-layer responses and the logit curve indicate that KAN encodes
 841 highly concave, irregular tumour margins as a strong, monotone risk factor for malignancy, consistent with
 842 established radiological criteria.



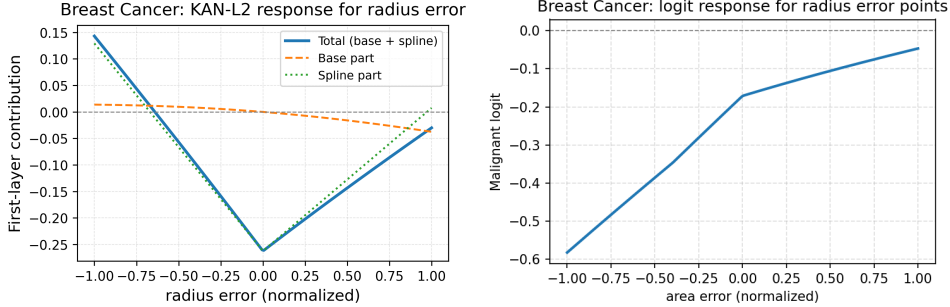
844
845
846
847
848
849
850
851
852
853 (a) KAN first-layer responses of *worst concave* (b) Malignant-class logit as a function of *worst*
 854 *points* *concave points* (normalized), obtained by varying
 855 only this feature around a reference input.

856
857 Figure 4: Breast Cancer dataset interpretability visualizations using KAN-L2 feature choice.

860 C.3 RADIUS ERROR POINTS INTERPRETABILITY STUDY

861 For *radius error*, the first-layer KAN response is approximately V-shaped (Figure 5a), indicating that the net-
 862 work uses a nonlinear code in which moderate values drive the hidden activation most strongly negative, while
 863 very small or very large errors have a weaker effect. When this representation is propagated through the second
 KAN layer, the resulting malignant-class logit becomes roughly monotone in *radius error* (Figure 5b): higher

radius error leads to a larger malignant logit with a mild saturation effect. This suggests that the model combines the V-shaped hidden code with other directions in feature space so that, at the prediction level, increasingly irregular tumour borders are treated as a stronger risk factor for malignancy, consistent with radiological practice.



(a) KAN first-layer responses of *radius error* (b) Malignant-class logit as a function of *radius error points* (normalized), obtained by varying only this feature around a reference input.

Figure 5: Breast Cancer dataset interpretability visualizations. (a) KAN first-layer responses for the top three features. (b) Logit sensitivity curve with respect to *worst concave points*.

D STABILITY AND REDUNDANCY ANALYSIS

To further assess the reliability of the proposed KAN-L2 feature selection method, we conducted two targeted diagnostic checks on representative datasets.

D.1 STABILITY

Evaluating the reproducibility of a feature selection method is a central principle in modern feature selection, formalized in stability selection frameworks (Meinshausen & Bühlmann, 2010). To quantify how consistently each method selects the same features under small perturbations of the dataset, we measure stability across $K = 5$ cross-validation folds. For a dataset with feature matrix $X \in \mathbb{R}^{n \times d}$, let $s^{(f)} = (s_1^{(f)}, \dots, s_d^{(f)})$ denote the feature-importance vector computed in fold f , and let $\alpha = 0.4$ be the retention fraction. We define the selected feature set in fold f as

$$A^{(f)} = \text{TopK}(s^{(f)}, k), \quad k = \lfloor \alpha d \rfloor.$$

To measure agreement between folds i and j , we use the Jaccard similarity, a standard stability index in feature selection studies (Nogueira et al., 2018),

$$J(A^{(i)}, A^{(j)}) = \frac{|A^{(i)} \cap A^{(j)}|}{|A^{(i)} \cup A^{(j)}|}, \quad J \in [0, 1].$$

The overall stability of a method M is given by the mean Jaccard similarity across all $\binom{K}{2}$ fold pairs:

$$\mu_M = \frac{2}{K(K-1)} \sum_{1 \leq i < j \leq K} J(A^{(i)}, A^{(j)}),$$

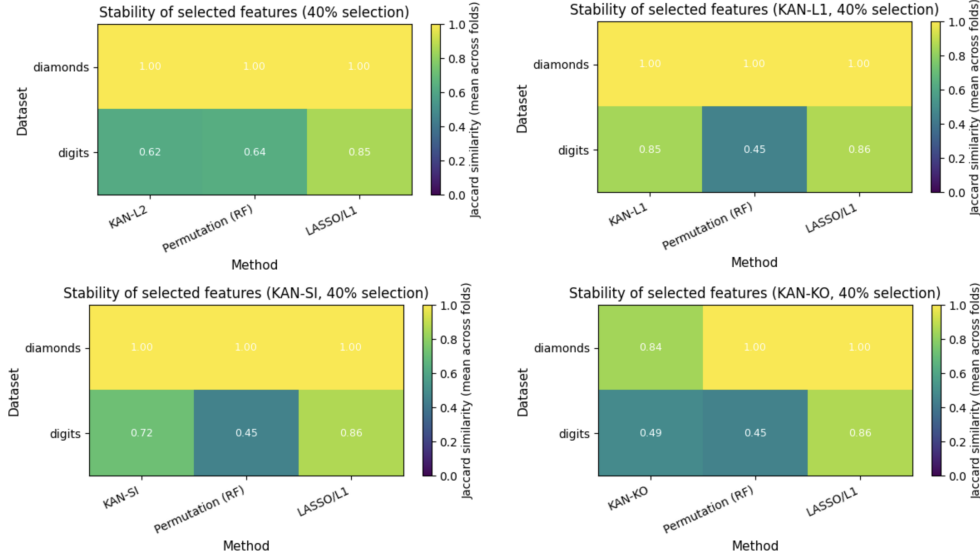
with variability quantified by the standard deviation

$$\sigma_M = \sqrt{\frac{2}{K(K-1)} \sum_{1 \leq i < j \leq K} (J(A^{(i)}, A^{(j)}) - \mu_M)^2}.$$

Higher values of μ_M indicate greater reproducibility of the selected features, while lower σ_M reflects reduced sensitivity to sampling noise.

Across all four variants (KAN-L1, KAN-KO, KAN-SI, and KAN-L2)(See Figure 6, the stability heatmaps show a consistent pattern: feature selection is perfectly stable on the diamonds dataset (Jaccard ≈ 1.00 for every method), indicating that all selectors repeatedly identify essentially the same subset of features across

918 folds. On the more challenging digits dataset, stability decreases as expected because the feature space is larger
 919 and more redundant. Still, the KAN-based methods remain competitive with classical baselines: KAN-L1
 920 (0.85), KAN-KO (0.49), KAN-SI (0.72), and KAN-L2 (0.62) all match or exceed the stability of Random
 921 Forest permutation importance (≈ 0.45), and are generally close to LASSO/L1 (≈ 0.86). Overall, these results
 922 show that the proposed KAN-based feature scores produce stable and reproducible feature subsets, especially
 923 when compared to existing selectors, and remain reliable even in high-dimensional settings



943 Figure 6: Stability analysis

944 D.2 REDUNDANCY

947 To determine whether a method selects mutually correlated (and therefore potentially redundant) predictors,
 948 we compare the correlation structure of the full feature space with that of the selected subset. This follows the
 949 classical redundancy perspective underlying mRMR-style criteria (Peng et al., 2005). Let $C = \text{corr}(X) \in$
 950 $\mathbb{R}^{d \times d}$ be the Pearson correlation matrix of all features. The baseline redundancy level is defined as the average
 951 absolute correlation across all off-diagonal pairs:

952
$$\overline{|\rho|}_{\text{all}} = \frac{2}{d(d-1)} \sum_{1 \leq i < j \leq d} |C_{ij}|.$$

955 For the selected subset $S = A^{(f)}$ of size k , the redundancy among selected features is computed from the
 956 submatrix $C_S = C[S, S]$:

957
$$\overline{|\rho|}_{\text{sel}} = \frac{2}{k(k-1)} \sum_{\substack{i, j \in S \\ i < j}} |C_{ij}|.$$

959 A method is said to reduce redundancy when

960
$$\overline{|\rho|}_{\text{sel}} < \overline{|\rho|}_{\text{all}},$$

962 indicating that the selected predictors are less correlated than the dataset as a whole. Comparable values sug-
 963 gest that the selector does not inflate redundant structure, even without an explicit sparsity or decorrelation
 964 constraint. KAN variants offer complementary advantages: KAN-L2 provides the cleanest, least-redundant
 965 subsets; KAN-L1 yields sparse, efficient selections; KAN-KO balances structure and redundancy; and KAN-
 966 SI captures the strongest predictive signals even when features are correlated. Together, they offer flexible,
 967 task-adapted feature selection capabilities

968 Across both datasets in Figure 7, we determined the average feature correlation (all feature vs 40% selected
 969 features). The redundancy analysis shows that the KAN feature selectors behave differently depending on the
 970 selector. For diamonds, all methods start from the same baseline average correlation (≈ 0.267), but their se-
 971 lected subsets diverge significantly: KAN-L2 achieves the strongest redundancy reduction (down to 0.054),
 KAN-L1 and KAN-KO reduce redundancy moderately (0.054 – 0.097), whereas KAN-SI increases redun-
 dancy (0.328), meaning its top-ranked features are more mutually correlated. For digits, where the baseline

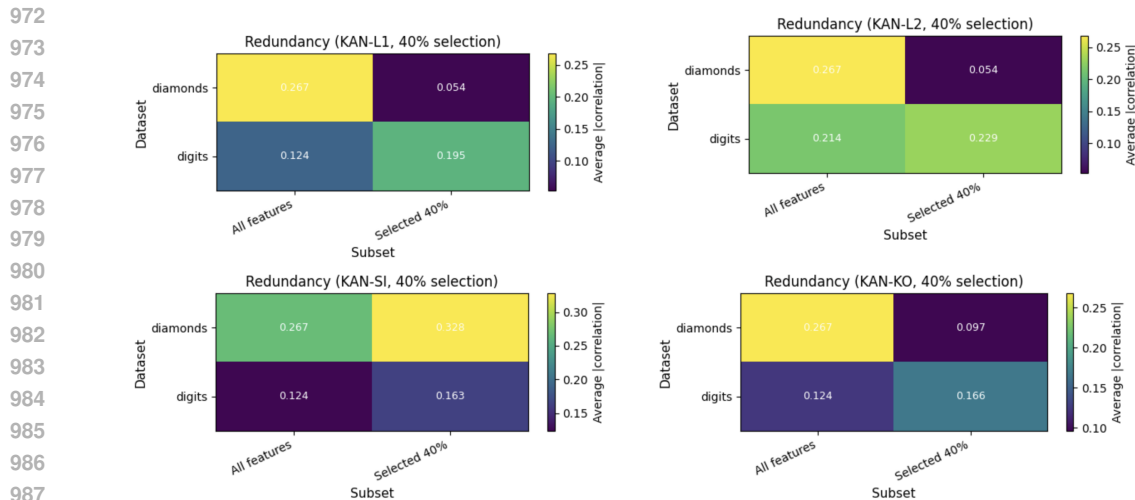


Figure 7: Redundancy analysis

redundancy is lower (≈ 0.124), KAN-L2 again maintains controlled redundancy (0.229), KAN-KO and KAN-SI show modest increases (0.166 – 0.195), while KAN-SI increases correlation the most. Overall, the results show that KAN-L2 consistently provides the least redundant subsets, KAN-L1 and KAN-KO give moderate redundancy behaviour, and KAN-SI tends to prioritise stronger, more correlated signals, which may help accuracy but does not enforce redundancy minimisation.

E PREDICTION PERFORMANCES PER FEATURE RETENTION LEVELS

We provided a predictive analysis in this part for different retention levels.

E.1 AVERAGE PREDICTORS PERFORMANCES AT DIFFERENT RETENTION LEVELS

Figure equation 8 shows that Across both the Breast Cancer and Digits datasets, the heatmaps show that most selectors achieve near-saturated predictive performance once a moderate fraction of features is retained. For Breast Cancer, accuracy becomes very stable around 94-96% from roughly 30-40% feature retention onward, indicating substantial redundancy in the original predictors: using only one third of the features is almost as good as using them all. Differences between methods are most visible under aggressive pruning (10-20% retention). In this regime, LASSO/L1, Mutual Information, SVM-RFE, and the KAN-based selectors (especially KAN-L1 and KAN-SI) remain comparatively robust, whereas KAN-KO is noticeably unstable at 10% retention.

E.2 INSIGHT AT 60% RETENTION LEVEL

In this section, we provide a comprehensive overview of the performance results of our models under different feature selection methods across all datasets considered in this study. The appendix is organized into two main subsections: the first subsection focuses on the **classification datasets**, while the second subsection presents the results for the **regression datasets**. For each dataset, we report model performance at 60% feature retention, highlighting the impact of feature selection techniques on the predictive accuracy.

E.2.1 CLASSIFICATION DATA SETS

This subsection presents the detailed results for all classification datasets, including breast cancer, digits, make_classification, and wine. Each table reports the macro F1 score of different predictor models using the full features data as well as reduced (60% reduction) selected feature sets obtained via KAN-based selection, LASSO, Mutual Information, Permutation importance, and SVM-RFE. The results illustrate how feature selection affects model performance and help identify the most effective selection strategies for each dataset. *Breast cancer.* (Table 2) Tree ensembles (GB, RF, XGB) remain competitive or improve with KAN-based selectors. In particular, KAN-L1 attains the top or tied-top accuracy for RF and XGB, and KAN-KO is best for GB. Logistic regression (LogReg) favors sparsity-oriented embedded/wrapper methods (SVM-RFE, LASSO), consistent with its linear inductive bias.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

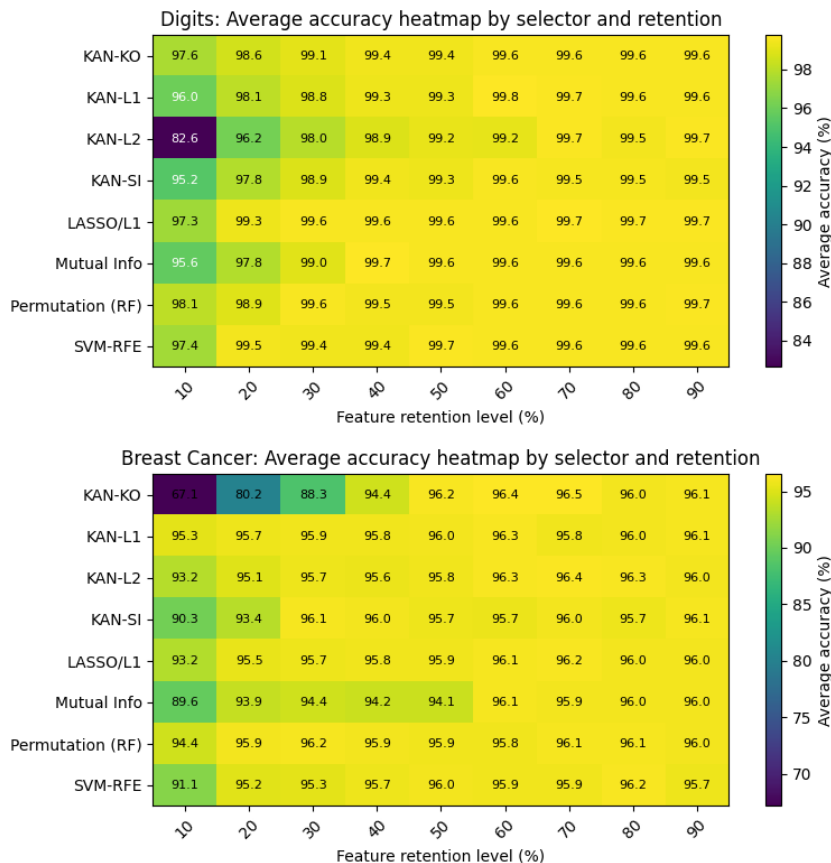


Figure 8: Average Accuracy Score (averaged across XGB, Random Forest and Gradient Boosted Trees, and Logistic Regression Classification) for 10/20/30/40/50/60/70/80/90% retained features) for the Digits and Breast cancer data sets

Table 2: 60% retention level for the *breast cancer* dataset

Models	All Features	KAN-KO	KAN-L1	KAN-L2	KAN-SI	LASSO/L1	Mutual Info	Perm. (RF)	SVM-RFE
GB	0.9624	0.9698	0.9603	0.9603	0.9623	0.9603	0.9623	0.9605	0.9587
LogReg	0.9487	0.9431	0.9433	0.9489	0.9374	0.9488	0.9508	0.9376	0.9526
RF	0.9642	0.9678	0.9698	0.9661	0.9585	0.9678	0.9606	0.9584	0.9547
XGB	0.9716	0.9753	0.9773	0.9754	0.9698	0.9678	0.9718	0.9773	0.9681

Digits. (Table 3) Performance is near a ceiling (LogReg reaches 1.00), so differences are small. Still, KAN-L1 (and KAN-SI for GB) match or set the best scores for GB/RF/XGB, showing that KAN selectors can remove 40% of inputs without harming multiclass performance on image-like digits.

Synthetic make_classification. (Table 4) Moderate gains appear for RF (KAN-KO/L2) and GB (KAN-SI), suggesting that removing weak/noisy variables helps tree ensembles; LogReg again prefers sparse selectors (all strong and tied within noise).

Wine. (Table 5) Mixed but consistent story: KAN-L1/L2 and Permutation (RF) yield the best or tied-best for GB/XGB, while LogReg peaks with LASSO (as expected). RF is very strong overall and slightly benefits from KAN-KO.

Takeaways. KAN-L1/L2 are reliable selectors for nonlinear learners (GB, RF, XGB), often matching or exceeding the full-feature baseline at 60% retention. KAN-KO is competitive when measured against the end-to-end loss (notably GB and RF), validating a perturbation view of importance. Linear LogReg benefits most from coefficient-based sparsity (LASSO, SVM-RFE). The ability to prune $\approx 40\%$ of features with no loss (and

Table 3: 60% retention level for the *digits* dataset

Models	All Features	KAN-KO	KAN-L1	KAN-L2	KAN-SI	LASSO/L1	Mutual Info	Perm. (RF)	SVM-RFE
GB	0.9869	0.9907	0.9925	0.9813	0.9925	0.9906	0.9907	0.9925	0.9907
LogReg	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	0.9981	1.0000
RF	0.9944	0.9981	1.0000	0.9925	0.9981	0.9962	0.9962	0.9962	0.9963
XGB	0.9963	0.9963	0.9981	0.9926	0.9944	0.9981	0.9981	0.9963	0.9963

Table 4: 60% retention level for the *make_classification* dataset

Models	All Features	KAN-KO	KAN-L1	KAN-L2	KAN-SI	LASSO/L1	Mutual Info	Permutation (RF)	SVM-RFE
GB	0.8838	0.8879	0.8879	0.8879	0.8919	0.8617	0.8819	0.8879	0.8819
LogReg	0.8577	0.8596	0.8596	0.8596	0.8596	0.8477	0.8577	0.8596	0.8577
RF	0.8999	0.9139	0.9059	0.9139	0.9059	0.8819	0.8879	0.9098	0.8919
XGB	0.8959	0.8979	0.9019	0.9039	0.8999	0.8779	0.8959	0.8979	0.8959

Table 5: 60% retention level for the *wine* dataset

Models	All Features	KAN-KO	KAN-L1	KAN-L2	KAN-SI	LASSO/L1	Mutual Info	Permutation (RF)	SVM-RFE
GB	0.9452	0.9449	0.9566	0.9566	0.9458	0.9505	0.9507	0.9566	0.9450
LogReg	0.9549	0.9663	0.9611	0.9611	0.9610	0.9774	0.9663	0.9611	0.9611
RF	0.9832	0.9837	0.9676	0.9676	0.9784	0.9781	0.9729	0.9728	0.9732
XGB	0.9724	0.9674	0.9615	0.9615	0.9671	0.9667	0.9557	0.9670	0.9728

occasional gains) indicates substantial redundancy and supports using KAN-based selection as a practical dimensionality reduction step for tabular classification.

E.2.2 REGRESSION DATA SETS

This subsection presents the results for all regression datasets, including California housing, diabetes, and *make_regression*. Each table reports the predictive performance of predictor models, measured in terms of R^2 , under various feature selection methods using a data resulting in 60% feature retention, as well as the data having all features. The analysis emphasizes the robustness of models to feature reduction and highlights which selection methods preserve or enhance model performance. On California dataset Table 6, ensembles remain resilient at 60% retention. GB peaks with RF permutation; KAN-SI and KAN-KO track the all-features baseline closely. RF improves with KAN-L1, indicating redundancy removal helps tree splits. Ridge is fragile: all-features wins while KAN-L1/L2 over-prune complementary linear signal. XGB prefers all-features, with KAN-SI and Mutual Information nearly matching, SVM-RFE trailing slightly. Across selectors, differences remain modest.

Diabetes dataset Table 7 exemplifies small-n, noisy regression. GB favors LASSO; RF performs best with all features. Ridge is exceptionally steady near 0.41, tying All and KAN-KO, and barely changing under other selectors. XGB prefers RF permutation. KAN-L1 consistently underperforms, suggesting over-sparsification removes weak yet useful signal. Overall, cautious selection helps, but aggressive pruning hurts most predictors. Mutual Information lags across models slightly.

Diamonds dataset 8 shows a strong signal, near-ceiling performance. Ensembles (RF, GB, XGB) remain virtually unchanged; tiny gains appear with KAN-KO/L2/SI. Ridge benefits modestly from LASSO, suggesting linear redundancy. Overall, selection is largely neutral: most approaches match the all-features baseline within the third decimal. Thus, when structure and signal to noise ratio are high, compact subsets neither help nor harm meaningfully. KAN methods remain competitive On synthetic *make_regression* (Table. 9), signal to noise is high. Ridge sits essentially at one across most selectors, confirming linear recoverability. RF improves modestly with KAN-L2/SI; GB and XGB notch their best with SVM-RFE, though gaps are tiny. KAN variants are stable and competitive, neither overshooting nor collapsing. Overall, many selectors converge, reflecting informative features and limited redundancy. Permutation, LASSO, Mutual Information perform similarly

Table 6: 60% retention level for the *California* dataset

Models	All Features	KAN-KO	KAN-L1	KAN-L2	KAN-SI	LASSO/L1	Mutual Info	Perm. (RF)	SVM-RFE
GB	0.7858	0.7825	0.7416	0.7471	0.7833	0.7565	0.7832	0.7909	0.7563
RF	0.8103	0.8117	0.8266	0.8230	0.8180	0.8257	0.8179	0.8109	0.8257
Ridge	0.6037	0.5886	0.3068	0.3933	0.5854	0.5933	0.5854	0.5946	0.5933
XGB	0.8431	0.8317	0.8318	0.8295	0.8409	0.8284	0.8395	0.8382	0.8261

Table 7: 60% retention level for the *diabetes* dataset

Models	All Features	KAN-KO	KAN-L1	KAN-L2	KAN-SI	LASSO/L1	Mutual Info	Perm. (RF)	SVM-RFE
GB	0.4208	0.4157	0.2405	0.4204	0.4021	0.4326	0.3933	0.4041	0.4205
RF	0.4266	0.4032	0.2601	0.4015	0.4089	0.4132	0.3916	0.4091	0.4186
Ridge	0.4205	0.4209	0.2769	0.4027	0.4127	0.4130	0.4137	0.4127	0.4132
XGB	0.3591	0.3413	0.1716	0.3636	0.3432	0.3554	0.3255	0.3366	0.3682

Takeaways. KAN-SI and KAN-KO are stable selectors across regression tasks, often matching the all-feature baseline (California, `make_regression`). KAN-L1/L2 can benefit tree ensembles (RF, GB, XGB) by removing redundancy, but risk over-pruning in noisy, small- n data (Diabetes). Ridge regression is fragile to sparsification, performing best with all features unless signal-to-noise is very high. In high-signal datasets (Diamonds, `make_regression`), feature selection has negligible effect, with all methods converging near the ceiling. The ability to drop $\approx 40\%$ of features without loss, and sometimes small gains, demonstrates redundancy and supports KAN-based feature selection as a practical dimensionality reduction tool for regression. For the Digits dataset, feature redundancy is even more pronounced. Many selectors already achieve 95-98% average accuracy at 10% retention, and by 30-40% retention essentially all methods converge to $\approx 99-99.7\%$ accuracy, making them practically indistinguishable. The only clear weakness appears for KAN-L2 at the most extreme pruning level (10%), where performance drops to about 83%, but it quickly recovers once more features are kept. Overall, these results suggest that: (i) strong dimensionality reduction is possible without meaningful loss in predictive performance, especially for highly redundant datasets such as Digits; and (ii) KAN-based selectors are competitive with, and often comparable to, classical baselines such as LASSO, Mutual Information, and SVM-RFE, except for a few edge cases under very aggressive pruning.

F RUNTIME PROFILING OF FEATURE SELECTORS

Table 10 reports wall-clock times for computing feature scores on two representative datasets. Once a KAN has been trained, the associated selectors are relatively cheap: coefficient-based criteria (*Selector-L1*, *Selector-L2*) are as fast as or faster than LASSO, and the gradient-based score (*Selector-SI*) remains competitive with mutual information. The knock-out score (*Selector-KO*) is more expensive, as it requires repeated forward passes with ablated features, but is still considerably lighter than wrapper-style methods such as permutation importance or SVM-RFE, whose runtimes are dominated by repeated model retraining or resampling. The main overhead of our approach lies in the initial KAN training (*KAN-Train*), which is substantially slower than fitting a single linear or tree model. However, this cost is incurred once per dataset and can then be amortized across multiple selectors, retention levels, and interpretability analyses, whereas wrapper methods must pay a similar price every time a new feature subset is evaluated.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

Table 8: 60% retention level for the *diamonds* dataset

Models	All Features	KAN-KO	KAN-L1	KAN-L2	KAN-SI	LASSO/L1	Mutual Info	Perm. (RF)	SVM-RFE
GB	0.9710	0.9710	0.9708	0.9710	0.9710	0.9704	0.9707	0.9707	0.9707
RF	0.9800	0.9801	0.9792	0.9801	0.9801	0.9800	0.9800	0.9800	0.9800
Ridge	0.8805	0.8806	0.8824	0.8806	0.8806	0.8847	0.8797	0.8797	0.8797
XGB	0.9816	0.9817	0.9810	0.9817	0.9817	0.9810	0.9811	0.9812	0.9811

Table 9: 60% retention level for the *make_regression* dataset

Models	All Features	KAN-KO	KAN-L1	KAN-L2	KAN-SI	LASSO/L1	Mutual Info	Permutation (RF)	SVM-RFE
GB	0.9456	0.9418	0.9418	0.9501	0.9501	0.9500	0.9499	0.9500	0.9503
RF	0.8819	0.8922	0.8922	0.8976	0.8976	0.8972	0.8972	0.8972	0.8971
Ridge	0.99999	0.9885	0.9885	0.99999	0.99999	0.99999	0.99999	0.99999	0.99999
XGB	0.9233	0.9234	0.9234	0.9325	0.9325	0.9336	0.9341	0.9336	0.9344

Table 10: Selector Runtime Comparison on Diamonds and Wine Datasets

Selector	Diamonds (sec)	Wine (sec)
Selector-L2	0.022246	0.001013
Selector-L1	0.039464	0.001047
Selector-SI	0.125190	0.003804
LASSO/L1	0.292232	0.001218
Selector-KO	0.426020	0.023628
Mutual Info	2.349357	0.029882
Permutation (RF)	51.709088	0.799606
SVM-RFE	80.940045	0.021224
KAN-Train	527.862040	2.441519