# Unsupervised Features Ranking via Coalitional Game Theory for Categorical Data[⋆]

Chiara Balestra[1], Florian Huber[3], Andreas Mayr[2], and Emmanuel Müller[1]

[1] TU Dortmund, Germany
[2] Department of Medical Biometry, Informatics and Epidemiology,
University Hospital of Bonn, Germany
[3] University of Bonn, Germany

**Abstract.** Not all real-world data are labeled, and when labels are not available, it is often costly to obtain them. Moreover, as many algorithms suffer from the curse of dimensionality, reducing the features in the data to a smaller set is often of great utility. Unsupervised feature selection aims to reduce the number of features, often using feature importance scores to quantify the relevancy of single features to the task at hand. These scores can be based only on the distribution of variables and the quantification of their interactions. The previous literature, mainly investigating anomaly detection and clusters, fails to address the redundancy-elimination issue. We propose an evaluation of correlations among features to compute feature importance scores representing the contribution of single features in explaining the dataset's structure.
Based on Coalitional Game Theory, our feature importance scores include a notion of redundancy awareness making them a tool to achieve redundancy-free feature selection. We show that the deriving features' selection outperforms competing methods in lowering the redundancy rate while maximizing the information contained in the data. We also introduce an approximated version of the algorithm to reduce the complexity of Shapley values' computations.

**Keywords:** feature ranking · game theory · redundancy reduction

## 1 Introduction

In machine learning, both feature selection methods and reduction of dimensionality are often performed to increase interpretability and to reduce computational complexity. As an example, for unsupervised applications such as clustering [5] or anomaly detection [14], the curse of dimensionality poses a major challenge. Unsupervised feature selection enables the detection of data patterns, as well as the description of these patterns using a concise set of relevant features [20, 24]. The corresponding methods are mostly based on the analysis of multivariate data distributions, pairwise correlations, higher-order interactions among features, or pseudo-labels. The use of such complex measures implies that both the selection as well as the interpretation of why some
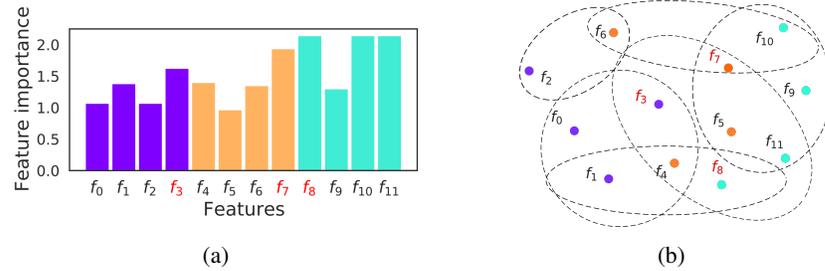
Fig. 1: (a) Unsupervised Shapley values-based feature importance scores. (b) Shapley values consider interactions within all possible subsets of features $f_i$s. In both figures, correlated subsets of features are color-coded and features selected by the proposed algorithm are marked in red.

features have been selected is challenging. On the one hand, selection requires basic measures to quantify the interaction within a set of features [5, 14, 16]. On the other hand, interpretation of higher-order interactions is non-straight-forward and requires the decomposition of complex non-linear, higher-order, and multivariate measures to feature importance scores.

In different application domains, raising understanding over the mechanisms of underlying machine learning techniques has become a crescent necessity. Assigning scores to features based on their contributions to the machine learning procedure plays a decisive role to this end. Feature importance scores are prevalent in supervised learning, e.g., random forests. At the same time, for unsupervised tasks, the literature is limited either to traditional scores [20, 24] not sensitive to higher-order interactions, or the scores are not easily interpretable higher-order correlation measures.

We propose new unsupervised feature importance scores decomposing the information contained in the data using axiomatic game-theoretic properties. In particular, Shapley values enable us to consider the interactions present in each possible subset of features (Figure 1(b)) and to assign importance scores to the single features accordingly. Our approach consists of two steps. In the first step, we introduce a game-theoretic solution to decompose the information contained in the dataset and assign importance scores to the single features. The scores obtained consider complex higher-order feature interactions, can be based on different correlation measures, and do not rely on specific notions of clustering or anomalies. In particular, we use Shapley values [18] to get the feature importance scores where features explaining the most information on the overall dataset obtain a higher score. In the second step, we take care of a mechanism to reduce the redundancy among features. To this end, feature importance scores are penalized through an information-theoretic measure of correlation to yield a redundancy-free feature selection. Figure 1(a) displays how correlated features are ranked similarly before applying the redundancy elimination step and how our method is capable of avoiding selecting highly correlated features. In the experimental results, we show that our ranking is achieving a redundancy free-ranking; the redundancy rate of the selected features is kept low both in synthetic and real datasets.

As a final remark, the scores flexibly rely on different correlation measures and are not bound to any clustering or anomaly detection goals. We choose to present the

| | versatile quality notion | feature ordering | iterative selection | redundancy awareness | higher-order interactions |
|---|---|---|---|---|---|
| UDFS [25] | ✗ | ✓ | ✗ | ✓ | ✓ |
| MCFS [2] | ✗ | ✓ | ✗ | ✓ | ✗ |
| NDFS [11] | ✗ | ✓ | ✗ | ✓ | ✗ |
| SPEC [26] | ✓ | ✓ | ✗ | ✗ | ✗ |
| LS [8] | ✓ | ✓ | ✗ | ✗ | ✗ |
| PFA [12] | ✓ | ✗ | ✗ | ✓ | ✗ |
| FSFC [27] | ✗ | ✗ | ✓ | ✓ | ✗ |
| **this paper** | ✓ | ✓ | ✓ | ✓ | ✓ |

Table 1: Summary table of the competing methods and this paper.

results obtained using the total correlation; hence, the presented experimental results are limited to discrete and categorical data. The procedure can be extended to mixed datasets replacing the total correlation or applying discretization on continuous data.

## 2  Related work

Dimensionality reduction helps avoid the curse of dimensionality and increases the interpretability of data and machine learning techniques. Different methods analyze the relationship among features, the class label, and the correlation among variables [23] and get feature importance scores in order to allow for a more aware use of machine learning by non-experts. Those scores are often not aware of correlations among variables, thus leading to a necessary integration of a redundancy awareness concept [19].

In 2007 game theory found application in supervised feature selection [6, 15] where the value function was defined as the accuracy or the generalization error of the trained model; to the best of our knowledge, the approaches proposed in the recent years are limited to labeled data. A recent paper [17] underlined how Shapley values spread through machine learning; in particular, they appear in several techniques to increase the overall interpretability of black-box models [13, 21] and new insights on Shapley values and their applications continue appearing in the literature [4].

As a downside, Shapley values are well known to be computationally expensive. Several approximations found place in the literature, e.g., [1, 3, 22] among others; the first attempt of a comprehensive survey of Shapley values' approximations is represented by Rozemberczki et al. [17]. To reduce the computational run-time, we implement Castro et al. [3] approximation, i.e., the most common Shapley values' approximation non relying on additional assumptions on the players.

As a parallel area of research, in recent years, unsupervised feature selection methods have raised strong interest in the community [10, 20]. We selected a representative sample within the vast number of unsupervised feature selection methods to compare the performance of our approach. Among them, UDFS [25] creates pseudo-labels to perform the feature selection in unlabelled data; MCFS and NDFS [2, 11] concentrate on keeping the clustering structure. LS [8] selects features by their local preserving power. PFA [12] tries to eliminate the downside of PCA while keeping the information within the data. Most of these algorithms tend to select features as a by-product of retaining a clustering structure in the data. Finally, FSFC [27] is meant to select only

non-redundant variables using a new definition of distance in the $k$-nearest neighbors. Table 1 illustrates a summary of the properties of the various methods analyzed in comparison with our paper.

## 3   Feature importance measures

Consider a $N$-dimensional dataset containing $D$ instances. We interpret each dimension as the realization set of a random variable, refer to the set of variables as $\mathcal{F} = \{X_1, \ldots, X_N\}$ and to each dimension $X_i$ as $i$th feature or variable. Feature selection methods often internally assign to subsets of features an importance score and output the subset maximizing the mentioned score. We propose to rank features considering their average contribution to all the possible subsets of features. The higher the average contribution of a feature is, the more convenient it is to keep it within the selected features. Additionally, we will also introduce redundancy awareness in these scores.

Given a function that assigns a value to each subset of features, assessing the *importance* of single features is not trivial as each feature belongs to $2^{N-1}$ subsets of features. In unsupervised contexts, we can assess the usefulness of a set of features measuring correlations or clustering properties. Throughout the manuscript, we stick to a value function that captures the maximal *information* contained in the data. Following this choice, the approach presented is restricted to categorical tabular data. We compute feature importance scores and obtain a ranking prioritizing features highly correlated with the rest of the dataset.

### 3.1   Feature importance score

We obtain feature importance scores using coalitional game theory. Each game is fully represented by the set of players $\mathcal{F}$ and a set function $v$ that maps each subset $\mathcal{A} \subseteq \mathcal{F}$ to $v(\mathcal{A}) \in \mathbb{R}$. $v$ is referred to as *value function* [18] and satisfies the following properties

1. $v(\emptyset) = 0$,
2. $v(\mathcal{A}) \geq 0$ for any $\mathcal{A} \subseteq \mathcal{F}$, and
3. $v(\mathcal{A}) \leq v(\mathcal{B})$ for any $\mathcal{A}, \mathcal{B} \subseteq \mathcal{F}$ such that $\mathcal{A} \subseteq \mathcal{B}$.

Working with unlabelled data, we can not rely on ground truth labels. Hence, we define value functions relying on intrinsic properties of the dataset; we opt for a value function measuring the independence of the features in $\mathcal{A} \subseteq \mathcal{F}$. One possible initialization for $v$ is the *total correlation* of $\mathcal{A}$.

**Definition 1.** *The* total correlation $C$ *of a set of variables* $\mathcal{A} \subseteq \mathcal{F}$ *is defined as*

$$C(\mathcal{A}) = \sum_{X \in \mathcal{A}} H(X) - H(\mathcal{A}). \tag{1}$$

$H(\mathcal{A})$ *is the Shannon entropy of the subset of discrete random variables* $\mathcal{A}$, *i.e.,*

$$H(\mathcal{A}) = -\sum_{\vec{x} \in \mathcal{A}} p_{\mathcal{A}}(\vec{x}) \log p_{\mathcal{A}}(\vec{x}) \tag{2}$$

*where* $p_{\mathcal{A}}(\cdot)$ *is the joint probability mass function of* $\mathcal{A}$.
$H(X)$ *is the Shannon entropy of* $X$, *i.e.,* $H(X) = -\sum_{x \in X} p_X(x) \log p_X(x)$.

We choose the total correlation as it satisfies properties (2) and (3), it has an intuitive meaning and can be easily extended such that it satisfies property (1).

Shannon entropy [7] measures the uncertainty contained in a random variable $X$ considering how uniform data are distributed: its value is close to zero when its probability mass function $p_X$ is highly skewed while, as the distribution approaches a uniform distribution, its value increases. Moreover, the Shannon entropy is a monotone non-negative function and can be extended such that $H(\emptyset) = 0$. We assume that all features in $\mathcal{F}$ are discrete as the extension of Shannon entropy to continuous variables is not monotone [9]. As a consequence of Shannon entropy's properties, the total correlation $C(\mathcal{A})$ is close to zero if the variables in $\mathcal{A}$ are independent, and it increases when they are correlated. To study the impact of adding a feature $Y$ to $\mathcal{A} \subseteq \mathcal{F}$, we compute the value function of the incremented subset $v(\mathcal{A} \cup Y)$ and compare it with $v(\mathcal{A})$: The difference $v(\mathcal{A} \cup Y) - v(\mathcal{A}) = H(\mathcal{A}) + H(Y) - H(\mathcal{A} \cup Y)$ is non-negative and measures how much $\mathcal{A}$ and $Y$ are correlated. We refer to $H(\mathcal{A}) + H(Y) - H(\mathcal{A} \cup Y)$ as *marginal contribution of $Y$ to $\mathcal{A}$*. If $\mathcal{A}$ and $Y$ are independent, then the marginal contribution of $Y$ to $\mathcal{A}$ equals zero. Vice versa, the marginal contribution grows the stronger the correlation between $Y$ and $\mathcal{A}$ is. As importance score, we assign to $X_i$ the average of its marginal contributions and we refer to it as $\phi(X_i)$, i.e.,

$$\phi(X_i) = \sum_{\mathcal{A} \subseteq \mathcal{F} \setminus X_i} \frac{1}{N\binom{N-1}{|\mathcal{A}|}} [H(\mathcal{A}) + H(X_i) - H(\mathcal{A} \cup X_i)] \qquad (3)$$

corresponding to the *Shapley value* of the player $X_i$ in the game $(\mathcal{F}, v)$ when $v$ is the total correlation. The general definition of Shapley values reads [18]:

**Definition 2.** *Given a coalitional game $(\mathcal{F}, v)$ and a player $X_i \in \mathcal{F}$, the* Shapley value *of $X_i$ is defined by*

$$\phi_v(X_i) = \sum_{\mathcal{A} \subseteq \mathcal{F} \setminus X_i} \frac{1}{N\binom{N-1}{|\mathcal{A}|}} [v(\mathcal{A} \cup X_i) - v(\mathcal{A})].$$

It can be proven that the Shapley value is the only function that satisfies the *Pareto optimality*, i.e., $\sum_{X_i \in \mathcal{F}} \phi_v(X_i) = v(\mathcal{F})$, the dummy, the symmetry and additive properties [18]. Moreover, Shapley values represent a fair assignment of resources to players based on their contributions to the game. We use the scores $\phi(X_i)$ to rank the features in the dataset $\mathcal{F}$. However, Shapley values do not consider redundancies, and linearly dependent features obtain equal Shapley values.

### 3.2   Importance scores of low correlated features

We use a dataset with three sets of correlated features (color-coded in Figure 1(a)), and we aim to select features from subsets with different colors; however, as we have already underlined, correlated features are characterized by similar Shapley values. In particular, the three highest Shapley values are obtained by correlated features in the blue-colored set. Before addressing the problem of redundancy-awareness inclusion in Shapley values, we show that the Shapley values rank features that are not correlated with the rest of the dataset in low positions.

---

**Algorithm 1** SVFS

---

1: **procedure** SVFS($\mathcal{F}, \epsilon$)
2:     $\mathcal{S} = \emptyset$
3:     **while** $\mathcal{F} \neq \emptyset$ **do**
4:         **while** $X \in \mathcal{F}$ **do**
5:             **if** $H(X) + H(\mathcal{S}) - H(\mathcal{S}, X) > \epsilon$ **then**
6:                 $\mathcal{F} = \mathcal{F} \setminus X$
7:             **else**
8:                 $\mathcal{F} = \mathcal{F}$
9:             $\mathcal{S} = \mathcal{S} \cup \arg\max_{X \in \mathcal{F}}\{\phi(X)\}$
10:            $\mathcal{F} = \mathcal{F} \setminus \mathcal{S}$
    **return** $\mathcal{S}$

---

**Algorithm 2** SVFR

---

1: **procedure** SVFR($\mathcal{F}$)
2:     $\mathcal{S} = \arg\max_{X \in \mathcal{F}}\{\phi(X)\}$
3:     ordered = [ ]
4:     ordered[0] = $\arg\max_{X \in \mathcal{F}}\{\phi(X)\}$,    j = 1
5:     $\mathcal{F} = \mathcal{F} \setminus \mathcal{S}$
6:     **while** $\mathcal{F} \neq \emptyset$ && $j < N$ **do**
7:         **for** $X \in \mathcal{F}$ **do**
8:             rk($X$) = $\phi(X) - H(X) - H(\mathcal{S}) + H(\mathcal{S}, X)$
9:             ordered[j] = $\arg\max_{X \in \mathcal{F}}\{\text{rk}(X)\}$
10:            $\mathcal{S} = \mathcal{S} \cup \arg\max_{X \in \mathcal{F}}\{\text{rk}(X)\}$
11:            $\mathcal{F} = \mathcal{F} \setminus \mathcal{S}$
12:            $j + +$
    **return** ordered

---

**Theorem 1.** *Given a subset of features $\mathcal{B} \subset \mathcal{F}$ that satisfies the following properties*

1. *for all $X_j \notin \mathcal{B}$ and for all $\mathcal{A} \subseteq \mathcal{F} \setminus \{X_j\}$, $H(\mathcal{A}) + H(X_j) = H(\mathcal{A} \cup X_j)$*
2. *for all $X_i \in \mathcal{B}$ and for all $\mathcal{A} \subseteq \mathcal{F} \setminus \{X_i\}$, $H(\mathcal{A}) + H(X_i) \geq H(X_i \cup \mathcal{A})$*

*then $\phi(X_i) \geq \phi(X_j)$ for all $X_i \in \mathcal{B}$ and $X_j \notin \mathcal{B}$.*

*Proof.* From (1) we know that, since the marginal contribution of $X_j \notin \mathcal{B}$ to any $\mathcal{A} \subseteq \mathcal{F} \setminus \{X_j\}$ is equal to zero, $\phi(X_j) = \sum_{\mathcal{A} \subseteq \mathcal{F} \setminus \{X_j\}} \frac{1}{N\binom{N-1}{|\mathcal{A}|}} \cdot 0 = 0$.

For any $X_i \in \mathcal{F}$ and $\mathcal{A} \subseteq \mathcal{F}$, we know that $H(\mathcal{A} \cup X_i) \leq H(\mathcal{A}) + H(X_i)$ from Shannon entropy's properties [7]. Hence, all marginal contributions are non-negative. Hence, $\phi(X_i) \geq 0 = \phi(X_j)$ for all $X_i \in \mathcal{B}$ and $X_j \notin \mathcal{B}$.

This concludes the proof.

Thus with total correlation as value function, Shapley values are non-negative and equal zero if and only if the feature is non-correlated with any subset of features. Moreover, features highly correlated with other subsets of features get high Shapley values.

## 4    Redundancy removal

We address the challenge of adding redundancy awareness to Shapley values. For this purpose, we develop a pruning criteria based on the total correlation and greedily rank features to get a redundancy-free ranking of features while still looking for features with high Shapley values. Feature selection based on this ranking selects the variables ranked first by Shapley values which show little dependencies.

We propose two algorithms. The Shapley Value Feature Selection (SVFS) needs a parameter $\epsilon$ representing the correlation among features that we are willing to accept; hence, SVFS requires some expert knowledge on the dataset to specify the parameter $\epsilon$ in an opportune interval. The Shapley Value Feature Ranking (SVFR) works automatically with an included notion of redundancy. We show that the two algorithms lead to consistent results in Section 6.5. At each step, both algorithms select the highest-ranked feature among the ones left.

We use a total correlation-based punishment; In particular, $H(\mathcal{A}) + H(X) - H(\mathcal{A} \cup X) \geq 0$ represents the strength of the correlation among $X$ and $\mathcal{A}$ and it is equal to zero if and only if $X$ and $\mathcal{A}$ are independent.

SVFS's inputs are the set of unordered features $\mathcal{F}$ and the parameter $\epsilon > 0$; $\epsilon$ plays the role of a stopping criterion and represents the maximum correlation that we are willing to accept within the set of selected features. Whenever $\epsilon$ is high, we end up with the ordering given by Shapley values alone; instead, for $\epsilon \approx 0$ the criterion can lead to the selection of the only features which are uncorrelated with the first one. The optimal range of $\epsilon$ highly depends on the dataset. We show that SVFS is robust w.r.t. the choice of $\epsilon$. At each iteration, SVFS excludes from the ranking the features $X$s that are correlated with the already ranked features $\mathcal{S} \subseteq \mathcal{F}$ more than $\epsilon$, i.e., $H(X) + H(\mathcal{S}) - H(\mathcal{S}, X) > \epsilon$, computes the Shapley values of all remaining features $X$ and adds to $\mathcal{S}$ the feature whose Shapley value is the highest. When there are no features left, it stops and returns $\mathcal{S}$.

SVFR takes as an input $\mathcal{F}$ and outputs a feature ranking without the need of any additional parameter. The ranking is aware of correlations as each of the Shapley values $\phi(X_i)$ is penalized using the correlation measure $H(X_i) + H(\mathcal{S}) - H(X_i \cup \mathcal{S})$ where $\mathcal{S}$ is the set of already ranked features, and $X_i$ is a new feature to be ranked. This algorithm provides a complete ranking of features and can be prematurely stopped including an upper bound of features we are willing to rank. The absence of the additional parameter $\epsilon$ is the main advantage of SVFR over SVFS.

## 5    Scalable algorithms

The size of $\mathcal{P}(\mathcal{F})$ being exponential in $N$, computing Shapley values involves $2^N$ evaluations of the value function. We use approximated Shapley values to obtain scalable versions of SVFR and SVFS. We implement three versions of the algorithms that differ only in the computations of Shapley values used:

– *full algorithm*: it uses the full computation of the Shapley values
– *bounded algorithm*: consider only subsets up to size $k$ fixed to compute the Shapley values

– *sampled algorithm*: it uses the approximation proposed by Castro et al. [3] based on $n$ random sampled subsets of features.

The time complexity for the sampled algorithm is $\mathcal{O}(D \cdot n)$, for the bounded algorithm is $\mathcal{O}(D \cdot N^k)$ while for the full algorithm is $\mathcal{O}(D \cdot 2^N)$ where $N$ is the number of features and $D$ the number of samples in the dataset.

## 6  Experiments

We show that our feature ranking method outperforms competing representative feature selection methods in terms of redundancy reduction. Metrics such as NMI, ACC, and redundancy rate are often used in the previous literature to evaluate unsupervised feature selection methods. NMI and ACC focus on the cluster structure in the data; therefore, as clustering is not the goal of our approach, we compare it with the competing methods using the redundancy rate. The redundancy rate of $\mathcal{S} \subseteq \mathcal{F}$ is defined in terms of pairwise Pearson correlations, i.e.,

$$\text{Red}(\mathcal{S}) = \frac{1}{2m(m-1)} \sum_{X,Y \in \mathcal{S}, X \neq Y} \rho_{X,Y} \tag{4}$$

where $\rho_{X,Y} \in [0,1]$ is the Pearson correlation of features $X$ and $Y$. It represents the averaged correlation among the pairs of features in $\mathcal{S}$ and varies in the interval $[0,1]$: a $\text{Red}(\mathcal{S})$ close to 1 shows that many selected features in $\mathcal{S}$ are strongly correlated while a value close to zero indicates that $\mathcal{S}$ contains little redundancy. In the experiments, we use the *redundancy rate* as evaluation criteria re-scaling it to the interval $[0, 100]$ via the maximum pair-wise correlation to facilitate the comparison among different datasets.

### 6.1  Datasets and competing methods

We show a comparison against *SPEC*, *MCFS*, *UDFS*, *NDFS*, *PFA*, *LS* and FSFC [2, 8, 11, 12, 25–27].

We use various synthetic and publicly available datasets: the *Breast Cancer dataset*, the *Big Five Personalities Test dataset*[4] and the *FIFA dataset*[5]. The datasets that we use throughout the paper are all categorical or discrete. We consider subsets of the full dataset in order to apply the full versions of the algorithms and investigate the performance of the approximations of SVFR and SVFS at the end of the section.

### 6.2  Redundancy awareness

We compare the feature selection results of our algorithm against the competitors by evaluating the redundancy rate in Table 2. For the FIFA dataset, we select 15 features

---

[4] The first 50 features in the Big Five dataset are the categorical answers to the personality test's questions and are divided into 5 personalities' traits (10 questions for each personality trait). To apply the full algorithm, we select questions from different personalities and restrict to 10000 instances.

[5] We restrict to the 5000 highest-rated players by the overall attribute.

| | Breast Cancer | B5_balanced | B5_unbalanced | FIFA | Synthetic |
|---|---|---|---|---|---|
| NDFS | 36.30 | 22.11 | 20.75 | 18.97 | **1.49** |
| MCFS | 20.26 | 23.59 | 18.79 | 20.63 | 3.74 |
| UDFS | 33.59 | 28.13 | 35.18 | 57.73 | 4.06 |
| SPEC | 13.89 | 39.09 | 21.46 | 42.14 | 29.4 |
| LS | 7.05 | 28.83 | 58.25 | 48.28 | 100.00 |
| PFA | **5.10** | 23.22 | 34.28 | 57.42 | 35.84 |
| FSFC | 8.74 | 22.64 | 20.99 | 36.45 | 2.12 |
| SVFR | 6.68 | **15.65** | **18.02** | **14.79** | 1.51 |

Table 2: Redundancy rate of the sets of three selected features using the competing algorithms and SVFR (highlighted in green color in the table) on different datasets. The lowest rates are represented in bold characters.

from the entire data which characterize the *agility*, *attacking* and *defending* skills of the football players; we keep the whole datasets for Breast Cancer and synthetic data; in the case of the Big Five Personality Traits dataset, we select respectively 5 questions from three different personality traits for the balanced dataset and 9 features from one trait and 3 from other two personality traits in the case of the unbalanced dataset. In order to avoid bias towards the random selection of personality traits and features in the Big Five data, we average the redundancy rate over 30 trials on randomly selected personalities and variables both in the case of the balanced and unbalanced setup.

In each column, bold characters highlight the lowest redundancy rate. We use SVFR for ranking the features and select the three highest-ranked features. We consequently specified the parameters of the competing methods in order to get a selection of features as close to three features as possible. For FCFS we set $k = 4$ for BC dataset, $k = 8$ for FIFA dataset, $k = 8$ for the synthetic data and for *Big Five* dataset we use different $k$ at each re-run such that the number of selected variables varies between $2$ and $5$ and then we average the redundancy rates; for NDFS, MCFS, UDFS and LS we used $k = 5$ ($k$ being the number of clusters in the data); for the other competitors, we specify the number of features to be selected. Table 2 illustrates that SVFR outperforms the competing methods in nearly all the cases. In particular, while SVFR achieves low redundancy rates in all datasets, the competing algorithms show big differences in performance in the various datasets. On the Breast Cancer data and the synthetic dataset respectively, PFA and NDFS slightly outperform SVFR. However, they do not keep an average low redundancy rate on the other datasets. For reproducibility, we make the code publicly available [6].

### 6.3   Relevance of unsupervised feature selection and effectiveness

In Figure 2(a), each plot corresponds to a different subset of features of the Big Five dataset, i.e., 10 features selected from three different personality traits. Running SVFS with $\epsilon = 0.3$ we detect correlated features and avoids selecting them together as shown in the plots. Using the scaled versions of our algorithms from Section 5 we can extend the approach towards the complete Big Five dataset.

---

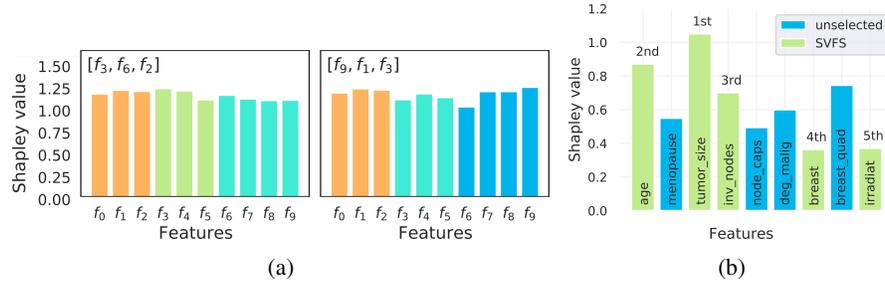[6] https://github.com/chiarabales/unsupervised_sv

Fig. 2: (a) Barplot of Shapley values and respective feature selections by SVFS ($\epsilon = 0.3$) in Big Five dataset restricted to 10 features. Different personalities traits are color-coded in each plot. (b) Barplot of Shapley values for Breast cancer data; in green, the ordering of features' selection by SVFS when $\epsilon = 0.5$.

|  | Big Five | Synthetic Data | Breast Cancer |
|---|---|---|---|
| $\epsilon = 0.2$ | $[11, 0, 5]$ | $[8, 7, 0]$ | $[2, 0, 8]$ |
| $\epsilon = 0.3$ | $[11, 0, 10]$ | $[8, 7, 2]$ | $[2, 0, 4, 6]$ |
| $\epsilon = 0.4$ | $[11, 0, 14]$ | $[8, 7, 3]$ | $[2, 0, 4, 8, 6]$ |
| $\epsilon = 0.5$ | $[11, 0, 14, 9]$ | $[8, 7, 3]$ | $[2, 0, 3, 8, 6]$ |
| $\epsilon = 0.6$ | $[11, 0, 14, 5]$ | $[8, 7, 3]$ | $[2, 0, 3, 8, 6]$ |
| $\epsilon = 0.7$ | $[11, 0, 14, 13]$ | $[8, 7, 3]$ | $[2, 0, 3, 4, 8, 6]$ |
| $\epsilon = 0.8$ | $[11, 0, 14, 13]$ | $[8, 7, 3, 0]$ | $[2, 0, 3, 4, 5, 8, 6]$ |
| SVFR | $[11, 0, 5, 10, 12, 8, 6, 2]$ | $[8, 7, 3, 0, 6, 5, 2, 10]$ | $[2, 0, 4, 6, 8, 5, 1, 3]$ |

Table 3: Orderings of selection given by SVFS for various $\epsilon$ and first 8 ranked features by SVFR. Features are color-coded in order to simplify the visualization.

Figure 1(a) represents the Shapley values of features in a 12 dimensional synthetic dataset where subsets of correlated features are color-coded. We measure the ability of the algorithm in selecting features from different subsets of correlated features; SVFS selects one feature from each subset of correlated features. In particular, when $\epsilon = 1$, SVFS achieves this goal by selecting $\{f_8, f_7, f_3\}$ while the ranking given by the Shapley values alone is $\{f_8, f_{10}, f_{11}\}$ which belong to the same subset of correlated features. This nicely underlines the inability of Shapley values to detect correlations and the necessity of integrating correlation-awareness to perform a feature selection.

Our unsupervised feature selection allow to construct more efficient psychological tests avoiding redundancies and reducing the number of questions that need to be answered without losing too much information.

### 6.4   Interpretation of feature ranking

We apply SVFS when $\epsilon = 0.5$ to the Breast Cancer dataset. In Figure 2(b), the resulting Shapley values and the ordering of selected features are displayed. The selection resulting from SVFS shows a low redundancy rate while the selected features (e.g., the size of the tumour, age, and the number of involved lymph nodes) are clearly in line with domain knowledge on risk factors for disease progression (label). Furthermore, the comparison with the ranking without redundancy awareness nicely highlights the importance of our approach to avoid redundancies when possible.

|      |         | $k = 1$ | $k = 3$ | $k = 5$ |
|------|---------|---------|---------|---------|
| BIG5 | random  | 0.04    | 0.19    | 0.33    |
|      | sampled | 0.04    | 0.37    | 0.49    |
|      | bounded | **0.08** | **0.56** | **0.55** |
| FIFA | random  | 0.06    | 0.24    | 0.35    |
|      | sampled | 0.00    | 0.33    | 0.40    |
|      | bounded | **1.00** | **0.67** | **0.80** |

Table 4: $recall@k$ for $k \in \{1, 3, 5\}$ comparing a random ranking and the rankings given by SVFR using the sampled and bounded algorithms to the full SVFR ranking. We show results for FIFA and Big Five datasets restricted to 15 features randomly chosen. Bold text highlights the best approximation.
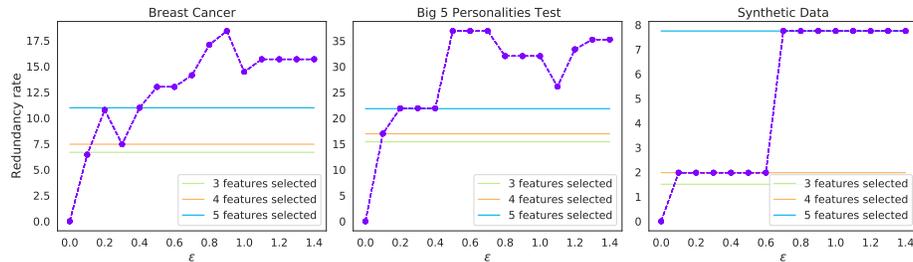


Fig. 3: Redundancy rates of the selected features' sets as a function of $\epsilon$ for SVFS (bullets points connected by the dashed line) and for $3, 4$, and $5$ selected features when using SVFR.

### 6.5   Comparison among the proposed algorithms

In Figure 3, we plot a comparison among SVFS and SVFR w.r.t. the redundancy rate on three datasets with different values of $\epsilon$. As benchmarks, we use for SVFR the selection of $3, 4$ and $5$ features respectively while for SVFS, $\epsilon$ varies in the interval $[0, 1.4]$ with steps of size $0.1$.

Using the number of features as a stopping criterion in SVFR would produce consistent results to SVFS: as an example, using the breast cancer data the ranking given by SVFR, i.e., $[2, 0, 4, 6, 8, 5, 1, 3]$, is consistent with the selection given by SVFS respectively using $\epsilon = 0.2$ and $\epsilon = 0.6$, i.e., $[2, 0, 8]$ and $[2, 0, 3, 8, 6]$. Table 3 shows a full comparison among the SVFR and SVFS on three different representative datasets. We recommend applying SVFS when no previous knowledge of the data is available and it is hard to establish an optimal range for $\epsilon$. Vice versa, one could apply SVFR when the expertise in the dataset domain allows determining a reasonable number of features as stopping criterion or the observation of the ranking given can provide insights to the non-expert on which features to keep and which can be discarded for further analysis.

### 6.6   Run-time analysis

As a consequence of the full computation of Shapley values, the run-time of SVFR and SVFS increases exponentially with the number of features as shown by Figure 4. Using
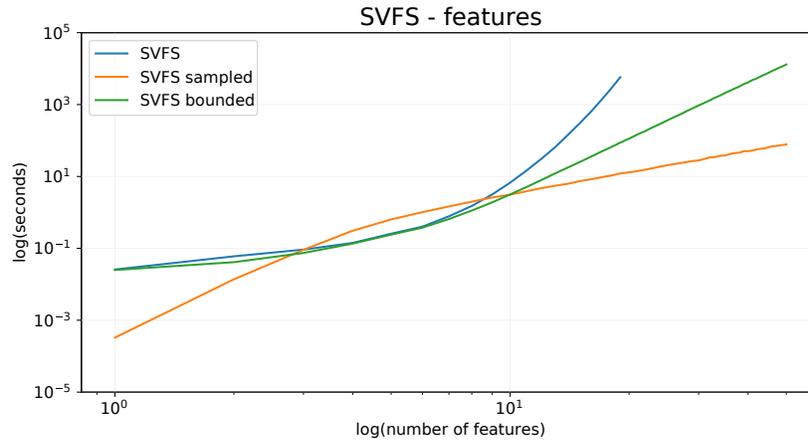
Fig. 4: Log-log plots of the run-time as a function of the number of features for the approximated and full SVFS ($\epsilon = 0.5$, $D = 1000$). The full SVFS is stopped with 20 features.

the approximated algorithms, this growth turns out to be slower. In particular, when using the sampled algorithm, the run-time increases only linearly with the number of features while the growth of the bounded algorithm's run-time is polynomial in the number of features. In the additional material, we show the log-log plot of the run-time for increased number of samples in the dataset. For each algorithm, we use random subsets of the Big Five dataset and average over 10 trails.

We further compare the rankings of the approximated and full algorithms using the $recall@k$ metric interpreting rankings of the full version of SVFR as ground truth. We use the Big Five dataset, randomly selecting 5 questions from 3 different personalities and average the scores over 100 trails (see Table 4). Overall, the results for the approximated algorithms clearly outperform random ordering - but still deviate often from the full versions. It is worth to note that the bounded algorithm using subsets up to size 5 performs better than the sampled version.

## 7   Conclusions

In the paper, we develop a new method to assess feature importance scores in unsupervised learning, bridging the gap between unsupervised feature selection and cooperative game theory. We integrate Shapley values with redundancy awareness making use of an entropy-based function to get feature importance scores.

We present two algorithms: SVFS implements feature selection using a redundancy aware criterion while SVFR assigns a ranking to each feature while being aware of correlations with previously ranked features. We show how the results of the two algorithms are consistent and state-of-the-art regarding their application. Our feature selection methods outperform previously proposed algorithms w.r.t. the redundancy rate. We additionally introduce approximated versions of the algorithms that are scalable to higher dimensions.

# References

1. M. A. BURGESS AND A. C. CHAPMAN, *Approximating the Shapley Value Using Stratified Empirical Bernstein Sampling*, in Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, Montreal, Canada, Aug. 2021, International Joint Conferences on Artificial Intelligence Organization, pp. 73–81.

2. D. CAI, C. ZHANG, AND X. HE, *Unsupervised feature selection for multi-cluster data*, in KDD, 2010.

3. J. CASTRO, D. GÓMEZ, AND J. TEJADA, *Polynomial calculation of the shapley value based on sampling*, Computers & Operations Research, (2009).

4. A. CATAV, B. FU, Y. ZOABI, A. L. W. MEILIK, N. SHOMRON, J. ERNST, S. SANKARARAMAN, AND R. GILAD-BACHRACH, *Marginal contribution feature importance - an axiomatic approach for explaining data*, in Proceedings of the 38th International Conference on Machine Learning, M. Meila and T. Zhang, eds., vol. 139 of Proceedings of Machine Learning Research, PMLR, 18–24 Jul 2021, pp. 1324–1335.

5. C.-H. CHENG, A. FU, AND F. ZHANG, *Entropy-based subspace clustering for mining numerical data*, KDD, (1999).

6. S. COHEN, G. DROR, AND E. RUPPIN, *Feature selection via coalitional game theory*, Neural computation, (2007).

7. T. M. COVER AND J. A. THOMAS, *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*, Wiley-Interscience, USA, 2006.

8. X. HE, D. CAI, AND P. NIYOGI, *Laplacian score for feature selection*, in NIPS, 2006.

9. A. V. LAZO AND P. RATHIE, *On the entropy of continuous probability distributions*, IEEE, (1978).

10. J. LI, K. CHENG, S. WANG, F. MORSTATTER, R. P. TREVINO, J. TANG, AND H. LIU, *Feature selection: A data perspective*, ACM Comput. Surv., (2017).

11. Z. LI, Y. YANG, J. LIU, X. ZHOU, AND H. LU, *Unsupervised feature selection using non-negative spectral analysis*, in AAAI, 2012.

12. Y. LU, I. COHEN, X. S. ZHOU, AND Q. TIAN, *Feature selection using principal feature analysis*, in MM, 2007.

13. S. M. LUNDBERG AND S.-I. LEE, *A unified approach to interpreting model predictions*, in NIPS, 2017.

14. H. NGUYEN, E. MÜLLER, J. VREEKEN, F. KELLER, AND K. BÖHM, *Cmi: An information-theoretic contrast measure for enhancing subspace cluster and outlier detection*, SDM, (2013).

15. K. PFANNSCHMIDT, E. HÜLLERMEIER, S. HELD, AND R. NEIGER, *Evaluating tests in medical diagnosis: Combining machine learning with game-theoretical concepts*, in IPMU, 2016.

16. D. RESHEF, Y. RESHEF, H. FINUCANE, S. GROSSMAN, G. MCVEAN, P. TURNBAUGH, E. LANDER, M. MITZENMACHER, AND P. SABETI, *Detecting novel associations in large data sets*, Science, (2011).

17. B. ROZEMBERCZKI, L. WATSON, P. BAYER, H.-T. YANG, O. KISS, S. NILSSON, AND R. SARKAR, *The shapley value in machine learning*, 2022.

18. L. S. SHAPLEY, *A value for n-person games*, Contributions to the Theory of Games, (1953).

19. A. SHEKAR, T. BOCKLISCH, P. SÁNCHEZ, C. STRAEHLE, AND E. MÜLLER, *Including multi-feature interactions and redundancy for feature ranking in mixed datasets*, in ECML PKDD, 2017.

20. S. SOLORIO-FERNÁNDEZ, J. CARRASCO-OCHOA, AND J. F. MARTÍNEZ-TRINIDAD, *A review of unsupervised feature selection methods*, Artificial Intelligence Review, (2019).

21. E. STRUMBELJ AND I. KONONENKO, *An efficient explanation of individual classifications using game theory*, J Mach Learn Res, (2010).
22. T. VAN CAMPEN, H. HAMERS, B. HUSSLAGE, AND R. LINDELAUF, *A new approximation method for the shapley value applied to the wtc 9/11 terrorist attack*, Social Network Analysis and Mining, (2018).
23. J. VERGARA AND P. ESTEVEZ, *A review of feature selection methods based on mutual information*, Neural. Comput. Appl., (2014).
24. S. WANG, J. TANG, AND H. LIU, *Embedded unsupervised feature selection*, in AAAI, 2015.
25. Y. YANG, H. SHEN, Z. MA, Z. HUANG, AND X. ZHOU, $l_{2,1}$-*norm regularized discriminative feature selection for unsupervised learning*, in IJCAI, 2011.
26. Z. ZHAO AND H. LIU, *Spectral feature selection for supervised and unsupervised learning*, in ICML, 2007.
27. X. ZHU, Y. WANG, Y. LI, Y. TAN, G. WANG, AND Q. SONG, *A new unsupervised feature selection algorithm using similarity-based feature clustering*, Computational Intelligence, (2019).
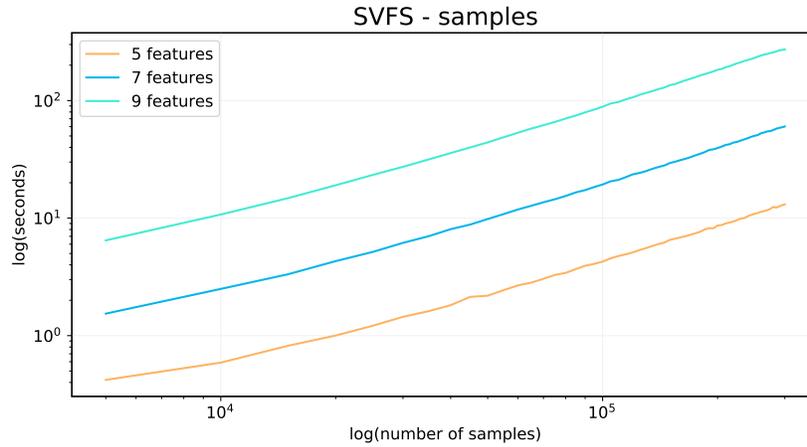
## Additional material



Fig. 5: Log-log plot of the run-time for the full SVFS with $\epsilon = 0.5$ as a function of the number of the samples $D$ and fixed number of features.

|  | features | samples |
|---|---|---|
| Breast Cancer dataset | 9 | 286 |
| Big Five dataset | 50 | 1013558 |
| FIFA20 dataset | 46 | 15257 |
| synthetic dataset | 12 | 10000 |

Table 5: Summary of the datasets' structures.