



# Grouped feature importance and combined features effect plot

Quay Au<sup>1</sup> · Julia Herbinger<sup>1</sup>  · Clemens Stachl<sup>2</sup> · Bernd Bischl<sup>1</sup> · Giuseppe Casalicchio<sup>1</sup>

Received: 23 April 2021 / Accepted: 18 April 2022  
© The Author(s) 2022

## Abstract

Interpretable machine learning has become a very active area of research due to the rising popularity of machine learning algorithms and their inherently challenging interpretability. Most work in this area has been focused on the interpretation of single features in a model. However, for researchers and practitioners, it is often equally important to quantify the importance or visualize the effect of feature groups. To address this research gap, we provide a comprehensive overview of how existing model-agnostic techniques can be defined for feature groups to assess the grouped feature importance, focusing on permutation-based, refitting, and Shapley-based methods. We also introduce an importance-based sequential procedure that identifies a stable and well-performing combination of features in the grouped feature space. Furthermore, we introduce the combined features effect plot, which is a technique to visualize the effect of a group of features based on a sparse, interpretable linear com-

---

Responsible editor: Martin Atzmueller, Johannes Fürnkranz, Tomáš Kliegr and Ute Schmid.

---

Quay Au and Julia Herbinger have contributed equally to this work.

---

✉ Julia Herbinger  
julia.herbinger@stat.uni-muenchen.de

Quay Au  
quayau@gmail.com

Clemens Stachl  
clemens.stachl@unisg.ch

Bernd Bischl  
bernd.bischl@stat.uni-muenchen.de

Giuseppe Casalicchio  
giuseppe.casalicchio@stat.uni-muenchen.de

<sup>1</sup> Department of Statistics, Ludwig-Maximilians-University Munich, 80539 Munich, Germany

<sup>2</sup> Institute of Behavioral Science and Technology, University of St. Gallen, 9000 St. Gallen, Switzerland

bination of features. We used simulation studies and real data examples to analyze, compare, and discuss these methods.

**Keywords** Grouped feature importance · Combined features effects · Dimension reduction · Interpretable machine learning

## 1 Introduction

Machine learning (ML) algorithms are nowadays used in many diverse fields e.g. in medicine (Shipp et al. 2002), criminology (Berk et al. 2009), and increasingly in the social sciences (Stachl et al. 2020b; Yarkoni and Westfall 2017). Interpretable models are paramount in many high-stakes settings, such as medical and juridical applications (Lipton 2018). However, well-performing ML models often bear a lack of interpretability. In the context of interpretable ML (IML) research, several model-agnostic methods to produce explanations for single features have been developed (Molnar 2019). Examples include the permutation feature importance (PFI; Fisher et al. 2019), leave-one-covariate out (LOCO) importance (Lei et al. 2018), SHAP values (Lundberg and Lee 2017), or partial dependence plots (PDP; Friedman 2001).

In many applications, it can be more informative to produce explanations for the importance or effect of a group of features (which we refer to as grouped interpretations) rather than for single features. It is important to note that the meaning of grouped interpretations, in general, differs from single feature interpretations, and resulting interpretations are usually not directly comparable (e.g., as Gregorutti et al. (2015) shows for the permutation feature importance). Hence, our aim is not to challenge single feature interpretations as both single and grouped feature interpretation methods measure different things and are useful on their own.

Grouped interpretations might be especially interesting for high-dimensional settings with hundreds or thousands of features. In particular, when analyzing the influence of these features visually (e.g., by plotting the marginal effect of a feature on the target) on a single feature level, this might result in an information overload which might not provide a comprehensive understanding of the learned effects (Molnar et al. 2020b). Furthermore, the runtime of some interpretation methods—such as Shapley values—does not scale linearly in the number of features. Hence, calculating them on a single feature level might not be computationally feasible for high-dimensional settings, making grouped computations a feasible remedy (Lundberg and Lee 2017; Covert et al. 2020; Molnar et al. 2020b).

From a use case perspective, the concept of grouped interpretations is particularly useful when the feature grouping is available *a priori* based on the application context. In that sense, features that either belong to the same semantic area (e.g., behaviors in psychology or biomarkers in medicine) or are generated by the same mechanism or device (e.g., fMRI, EEG, smartphones) can be grouped together to assess their joint effect or importance. For example, in our application in Sect. 7, we use a real-world use case from psychology that studies how the human behavior on smartphone app usage is associated to different personality traits (Stachl et al. 2020a). Features were extracted from longitudinal data collected from smartphones of 624 participants,



**Fig. 1** A possible process from group definition to grouped interpretations. First, the feature groups must be defined. A model is then fitted, typically on the feature space where the information of the pre-defined grouping might be used (e.g., if the fitting process is combined with a feature selection procedure) or ignored. When the best model is found, model-agnostic grouped interpretation methods are applied on the previously defined feature groups. A commonly used approach is to first obtain an overview of which groups are most important for achieving a good model performance (grouped feature importance) to subsequently analyze how the most important feature groups influence the model's prediction (grouped feature effect) (Color figure online)

and can be grouped into different behavioral classes (i.e., communication and social activity, app-usage, music consumption, overall phone activity, mobility). Another example is applications with sensor data (Chakraborty and Pal 2008), where multiple features measured by a single sensor naturally belong together, and hence grouped interpretations on sensor-level might be more informative.

There are also situations where the interpretation of single features might be misleading and where grouped interpretations can provide a remedy. Examples include datasets with time-lagged or categorical features (e.g., dummy or one-hot encoded categories) and the presence of feature interactions (Gregorutti et al. 2015). A concrete example for dummy encoded categorical features is shown in Appendix A.

Even in situations where feature groups are not naturally given in advance, it still might be beneficial to define groups in a data-driven manner and apply interpretation methods on groups of features (for examples, see Sect. 1.2).

Hence, compared to single feature interpretation methods, the grouping structure must be defined beforehand. A possible process—from group membership definition to modeling up to post-hoc interpretations—is illustrated in Fig. 1. Since defining the underlying group structure is a relevant step in this process, we discuss some applied techniques on how to find groups of features in Sect. 1.2. However, in this paper, we focus on the interpretation component once the groups are known (the green part in Fig. 1).

Although the grouped feature perspective is relevant in many applications, most IML research has focused on methods that attempt to provide explanations on a single-feature level. Model-agnostic methods for feature groups are rare and not well-studied.

## 1.1 Real data use cases with grouped features

In the following we summarize further exemplary predictive tasks with pre-specified feature groupings. These tasks will also be used in Sect. 3.4 for further empirical analysis. For more details on features and associated groups see Table 1.

*Heat value of fossil fuels* In this small scale regression task ( $n = 129$ ), the objective is to predict the heat value of fossil fuels from spectral data (Fuchs et al. 2015). In addition to one scalar feature (humidity), the dataset contains two groups of curve data, the first from the ultraviolet-visible spectrum (UVVIS) and the second from the near infrared spectrum (NIR).

**Table 1** Real world datasets with grouped features and their pre-specified group memberships

Dataset	Single features	Group membership	Description
<i>Birthweight</i>	age1, age2, age3	Age	Mother's age represented by 3 orthogonal polynomials
	lwt1, lwt2, lwt3	lwt	Mother's weight represented by 3 orthogonal polynomials
	White, black	Race	Mother's race (indicator functions)
	Smoke	Smoke	Smoking status (indicator function)
	ptl1, ptl2m	ptl	One, or two or more previous premature labors
	ht	ht	History of hypertension (indicator function)
	ui	ui	Presence of uterine irritability (indicator function)
	ftv1, ftv2, ftv3m	ftv	One, two, or three or more physician visits during first trimester
<i>Colon</i>	x1, ..., x5	Gene1	Gene expression data for gene 1
	⋮	⋮	⋮
<i>Fuelsubset</i>	x96, ..., x100	Gene20	Gene expression data for gene 20
	H20	H20	Humidity in percent
	UVVIS1, ..., UVVIS134	UVVIS	Data from the ultraviolet-visible spectrum (134 wavelength points)
	NIR1, ..., NIR231	NIR	Data from the near infrared spectrum (231 wavelength points)

*Birthweight* The *birthweight* dataset has data on 189 births at the Baystate Medical Centre in Massachusetts during 1986 (Venables and Ripley 2002). The objective is to predict the birth weight in kilograms from a set of 16 features, some of which are grouped (e.g., dummy encoded categorical features).

*Colon cancer* The *colon* dataset contains gene expression data of 20 genes (5 basis B-Splines each) for 62 samples from microarray experiments of colon tissue (Alon et al. 1999). The task is to predict cancerous tissue from the resulting 100 predictors.

## 1.2 Grouping procedures

Following the definitions of He and Yu (2010), we provide a brief overview of different procedures to define feature groups in a knowledge-driven and data-driven manner. In data-driven grouping, an algorithmic approach such as clustering or density estimation is used to define groups of features. Knowledge-driven grouping, on the other hand, uses domain knowledge to define the grouping structure of features. Throughout our

paper, we mainly assume a user-defined grouping structure. However, all methods introduced in this paper should also be compatible with an appropriate data-driven method if the defined groups have a meaningful interpretation.

### ***Data-driven grouping***

One method to group features in a data-driven manner is to use clustering approaches such as hierarchical clustering (Park et al. 2006; Toloşi and Lengauer 2011; Rapaport et al. 2008) or fuzzy clustering (Jaeger et al. 2003). These approaches often work well in highly correlated feature spaces, such as in genomics or medicine, where correlated features are grouped together so that no relevant information is discarded (Toloşi and Lengauer 2011). For instance, Jaeger et al. (2003) tackles a feature selection problem for a high-dimensional and intercorrelated feature space when working with microarray data. To simultaneously select informative and distinct genes, they first apply fuzzy clustering to obtain groups of similar genes from microarray data. Next, the informative representatives of each group are selected based on a suitable test statistic. The disadvantage of data-driven grouping is that groups depend only on the statistical similarity between features, which might not coincide with domain-specific interpretations (Chakraborty and Pal 2008).

### ***Knowledge-driven grouping***

Knowledge-driven group formation has the advantage that the dimensionality reduction might lead to better interpretability than the data-driven path. Gregorutti et al. (2015) apply a knowledge-driven approach in the context of multiple functional data analysis, where they then select groups for subsequent modeling based on their group importance values. Chakraborty and Pal (2008) also select groups of features, where data from one sensor (e.g., to capture satellite images in different spectral bands) represents a group. Hence, features are grouped based on their topical character (e.g., measurement device) rather than their shared statistical properties. Another use case of knowledge-driven grouping is described in Lozano et al. (2009), who group time-lagged features of the same time series for gene expression data. They use the given grouping structure in a group feature selection procedure and apply group LASSO as well as a boosting method.

## **1.3 Related work**

A well-known model that handles groups of features is the *group LASSO* (Yuan and Lin 2006), which extends the LASSO (Tibshirani 1996) for feature selection based on groups. Moreover, other extensions—e.g., to obtain sparse groups of features (Friedman et al. 2010), to support classification tasks (Meier et al. 2008) or non-linear effects (Gregorova et al. 2018)—also exist. However, group LASSO is a modeling technique that focuses on selecting groups in the feature space rather than quantifying their importance.

A large body of research already exists regarding the importance of individual features (see, e.g., Fisher et al. 2019; Hooker and Mentch 2019; Scholbeck et al. 2020). Hooker and Mentch (2019) distinguish between two loss-based feature importance approaches, namely permutation methods and refitting methods. Permutation meth-

ods measure the increase in expected loss (or error) after permuting a feature while the model remains untouched. Refitting methods measure the increase in expected loss after leaving out the feature of interest completely and hence require refitting the model (Lei et al. 2018). Since the model remains untouched in the former approach, interpretations refer to a specific fitted model, while interpretations for refitting methods refer to the underlying ML algorithm. Gregorutti et al. (2015) introduced a model-specific, grouped PFI score for random forests and applied this approach to functional data analysis. Valentin et al. (2020) introduced a model-agnostic grouped version of the model reliance score (Fisher et al. 2019). However, they focus more on the application and omit a detailed theoretical foundation. Recently, a general refitting framework to measure the importance of (groups of) features was introduced by Williamson et al. (2020). In their approach, the feature importance measurement is detached from the model level and defined by an algorithm-agnostic version to measure the intrinsic importance of features. The importance score is defined as the difference between the performance of the full model and the performance based on all features *except* the group of interest.

Permutation methods can be computed much faster than refitting methods. However, the PFI, for example, has issues when features are correlated and interact in the model due to extrapolation in regions without any or just a few observations (Hooker and Mentch 2019). Hence, interpretations in these regions might be misleading. To avoid this problem, alternatives based on conditional distributions or refitting have been suggested (e.g., Strobl et al. 2008; Nicodemus et al. 2010; Hooker and Mentch 2019; Watson and Wright 2019; Molnar et al. 2020a). Although the so-called conditional PFI provides a solution to this problem, its interpretation is different and “must be interpreted as the additional, unique contribution of a feature given all remaining features we condition on were known” (Molnar et al. 2020a). This property complicates the comparison with non-conditional interpretation methods. Therefore, we do not consider any conditional variants in this paper.

A third class of importance measures is based on Shapley values (Shapley 1953), a theoretical concept of game theory. The SHAP (Lundberg and Lee 2017) approach quantifies the contribution of each feature to the predicted outcome and is a permutation-based method. It has the advantage that contributions of interactions are distributed fairly between features. Besides being computationally more expensive, SHAP itself is based on the model’s predicted outcome rather than the model’s performance (e.g., measured by the model’s expected loss). Casalicchio et al. (2019) extended the concept of Shapley values to fairly distribute the model’s performance among features and called it Shapley Feature IMPortance (SFIMP). A similar approach called SAGE has also been proposed by Covert et al. (2020), who showed the benefits of the method on various simulation studies. One approach that uses Shapley values to explain grouped features was introduced by de Mijolla et al. (2020). However, instead of directly computing Shapley importance on the original feature space, they first apply a semantically-meaningful latent representation (e.g. by projecting the original feature space into a lower dimensional latent variable space using disentangled representations) and compute the Shapley importance on the resulting latent variables. Williamson and Feng (2020) mention that their feature importance method based on Shapley values can also be extended to groups of features. Additionally,

Amoukou et al. (2021) investigated grouping approaches for Shapley values in the case of encoded categorical features and subset selection of important features for tree-based methods. The calculation of Shapley values on groups of features based on performance values has only been applied with regard to feature subset selection methods and not for interpretation purposes (Cohen et al. 2005; Tripathi et al. 2020).<sup>1</sup>

After identifying which groups of features are important, the user is often interested in how they (especially the important groups) influence the model's prediction. Several techniques to visualize single-feature effects exist. These include partial dependence plots (PDP) (Friedman 2001), individual conditional expectation (ICE) curves (Goldstein et al. 2013), SHAP dependence plots (Lundberg et al. 2018), and accumulated local effects (ALE) plots (Apley and Zhu 2019). However, in the case of high-dimensional feature spaces, it is often not feasible to compute, visualize, and interpret single-feature plots for all (important) features. If features are grouped, visualization techniques become computationally more complex, and it may become even harder to visualize the results in an easily interpretable way. In the case of low-dimensional feature spaces, this might still be feasible, for example by using two-feature PDPs or ALE plots. Recently, effect plots that visualize the combined effect of multiple features have been introduced by Seedorff and Brown (2021) and Brenning (2021). They use principal component analysis (PCA) to reduce the dimension of the feature space and calculate marginal effect curves for the principal components. However, the employed dimension reduction method does not include information about the target variable and lacks sparsity (and hence, interpretability).

## 1.4 Contribution

Our contributions can be summarized as follows: We extend the permutation-based and refitting-based grouped feature importance methods introduced by Valentin et al. (2020) and Williamson et al. (2020) by comparing these methods to not only the full model (i.e., taking into account all features), but also to a null model (i.e., ignoring all features). Hence, we can quantify to what extent a group itself contributes to the prediction of a model without the presence of other groups. Furthermore, we introduce Shapley importance for feature groups and describe how these scores can be decomposed into single-feature importance scores of the respective groups. Our main contributions are: (1) We define a new algorithm to sequentially add groups of features depending on their importance, thereby enabling identification of well-performing combinations of groups. (2) We compare all grouped feature importance methods with respect to the main challenges that arise when applying these methods by creating small simulation examples. Subsequently, we provide recommendations for using and interpreting the respective methods correctly. (3) We introduce a model-agnostic method to visualize the joint effect of a group of features. To that end, we use a suitable dimension reduction technique and the conceptual idea of PDPs to calculate and plot the mean prediction of a sparse group of features with regard to their linear

---

<sup>1</sup> Feature subset selection methods usually aim to find sparse, well-performing feature combinations. Hence, the intended purpose of employing these methods is not to produce interpretability, but rather to generate a sufficient performance with fewer features.

combination. This novel method finally enables the user to visualize effects for groups of features. Finally, we showcase the usefulness of all these methods in real data examples.

The structure of this paper is as follows: First, we provide some general notation and definitions in Sect. 2. We formally define the grouped feature importance methods and introduce the sequential grouped feature importance procedure in Sects. 3 and 4, respectively. We compare these methods for different scenarios in Sect. 5. In Sect. 6, we introduce the combined features effect plot (CFEP) to visualize the effects of feature groups based on a supervised dimension reduction technique. Moreover, we also show the suitability of this technique compared to its unsupervised counterpart in a simulation study. Finally, in Sect. 7, all methods are applied to a real data example before summarizing and offering an outlook for future research in Sect. 8.

## 2 Background and notation

Analogous to Casalicchio et al. (2019), we use the term *feature importance* to refer to the influence of features on a model's predictive performance, which we measure by the expected loss when we perturb these features in a permutation approach or remove these features in a refitting approach.

### 2.1 General notation

Consider a  $p$ -dimensional feature space  $\mathcal{X} = (\mathcal{X}_1 \times \dots \times \mathcal{X}_p)$  and a one-dimensional target space  $\mathcal{Y}$ . The corresponding random variables that are generated from these spaces are denoted by  $X = (X_1, \dots, X_p)$  and  $Y$ . We denote a ML prediction function that maps the  $p$ -dimensional feature space to a one-dimensional target space by  $\hat{f} : \mathcal{X} \rightarrow \mathbb{R}$  for regression tasks.<sup>2</sup> ML algorithms try to learn this functional relationship using  $n \in \mathbb{N}$  i.i.d. observations drawn from the joint space  $\mathcal{X} \times \mathcal{Y}$  with unknown probability distribution  $\mathcal{P}$ . The resulting dataset is denoted by  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$ , where the vector  $\mathbf{x}^{(i)} = (x_1^{(i)}, \dots, x_p^{(i)})^\top \in \mathcal{X}$  is the  $i$ -th observation associated with the target variable  $y^{(i)} \in \mathcal{Y}$ . The  $j$ -th feature is denoted by  $\mathbf{x}_j = (x_j^{(1)}, \dots, x_j^{(n)})^\top$ , for  $j = 1, \dots, p$ . The dataset  $\mathcal{D}$  can also be written in matrix form:

$$\begin{pmatrix} x_1^{(1)} & \dots & x_p^{(1)} & y^{(1)} \\ \vdots & \ddots & \vdots & \vdots \\ x_1^{(n)} & \dots & x_p^{(n)} & y^{(n)} \end{pmatrix} = (\mathbf{X}, \mathbf{Y}), \text{ with } \mathbf{X} = \begin{pmatrix} x_1^{(1)} & \dots & x_p^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(n)} & \dots & x_p^{(n)} \end{pmatrix}, \mathbf{Y} = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{pmatrix}. \quad (1)$$

The general error measure  $\rho(\hat{f}, \mathcal{P}) = \mathbb{E}(L(\hat{f}(X), Y))$  of a learned model  $\hat{f}$  is measured by a loss function  $L$  on test data drawn independently from  $\mathcal{P}$  and can be

<sup>2</sup> The target space is defined by  $\mathbb{R}^g$  in the case of scoring classifiers with  $g$  classes.



estimated using unseen test data  $\mathcal{D}_{\text{test}}$  by

$$\hat{\rho}(\hat{f}, \mathcal{D}_{\text{test}}) = \frac{1}{|\mathcal{D}_{\text{test}}|} \sum_{(\mathbf{x}, y) \in \mathcal{D}_{\text{test}}} L(\hat{f}(\mathbf{x}), y). \tag{2}$$

The application of an ML algorithm (or *learner*)  $\mathcal{I}$  to a given dataset  $\mathcal{D}$  results in a fitted model  $\mathcal{I}(\mathcal{D}) = \hat{f}_{\mathcal{D}}$ . The *expected generalization error* of a learner  $\mathcal{I}$  takes into account the variability introduced by sampling different datasets  $\mathcal{D}$  of equal size  $n$  from  $\mathcal{P}$  and is defined by

$$GE(\mathcal{I}, \mathcal{P}, n) = \mathbb{E}_{|\mathcal{D}|=n}(\rho(\mathcal{I}(\mathcal{D}), \mathcal{P})). \tag{3}$$

In practice, resampling techniques such as cross-validation or bootstrapping on the available dataset  $\mathcal{D}$  are used to estimate Eq. (3). Resampling techniques usually split the dataset  $\mathcal{D}$  into  $k \in \mathbb{N}$  training datasets  $\mathcal{D}_{\text{train}}^i, i = 1, \dots, k$ , of roughly the same size  $n_{\text{train}} < n$ . Eq. (3) can be estimated by

$$\widehat{GE}(\mathcal{I}, \mathcal{D}, n_{\text{train}}) = \frac{1}{k} \sum_{i=1}^k \hat{\rho}(f_{\mathcal{D}_{\text{train}}^i}, \mathcal{D}_{\text{test}}^i). \tag{4}$$

In the following, we often associate the set of numbers  $\{1, \dots, p\}$  in a one-to-one manner with the features  $\mathbf{x}_1, \dots, \mathbf{x}_p$  by referring a number  $j \in \{1, \dots, p\}$  as feature  $x_j$ . We call  $G \subset \{1, \dots, p\}$  a *group of features*.

### 2.2 Permutation feature importance (PFI)

Fisher et al. (2019) proposed a model-agnostic version of the PFI measure used in random forests (Breiman 2001). The PFI score of the  $j$ -th feature of a fitted model  $\hat{f}$  is defined as the increase in expected loss after permuting feature  $X_j$ :

$$\text{PFI}_j(\hat{f}) = \mathbb{E}(L(\hat{f}(X_{[j]}), Y)) - \mathbb{E}(L(\hat{f}(X), Y)). \tag{5}$$

Here,  $X_{[j]} = (X_1, \dots, X_{j-1}, \tilde{X}_j, X_{j+1}, \dots, X_p)$  is the  $p$ -dimensional random variable vector of features, where  $\tilde{X}_j$  is an independent replication of  $X_j$  following the same distribution. The idea behind this method is to break the association between the  $j$ -th feature and the target variable by permuting its feature values. If a feature is not useful for predicting an outcome, changing its values by permutation will not increase the expected loss.<sup>3</sup> For an accurate estimation of Eq. (5), we would need to calculate all possible permutation vectors over the index set  $\{1, \dots, n\}$  (see Casalicchio et al. (2019) for an in-depth discussion on this topic). However, Eq. (5) can be approximated on a dataset  $\mathcal{D}$  with  $n$  observations by Monte Carlo integration using  $m$

<sup>3</sup> We consider the case of loss functions that are to be minimized. Hence, the larger PFI<sub>*j*</sub>, the more substantial the increase in expected loss and the more important the  $j$ -th feature.

random permutations:

$$\widehat{\text{PFI}}_j(\hat{f}, \mathcal{D}) = \frac{1}{nm} \sum_{i=1}^n \sum_{k=1}^m \left( L\left(\hat{f}((x_1^{(i)}, \dots, x_j^{(\tau_k^{(i)})}, \dots, x_p^{(i)})\right), y^{(i)})\right) - L(\hat{f}(\mathbf{x}^{(i)}, y^{(i)})) \right), \quad (6)$$

where  $\tau_k$  is a random permutation vector of the index set  $\{1, \dots, n\}$  for  $k = 1, \dots, m$  permutations.<sup>4</sup>

Equation (6) could also be embedded into a resampling technique, where the permutation is always applied on the held-out test set of each resampling iteration (Fisher et al. 2019). However, this leads to refits and is computationally more expensive. The resulting resampling-based PFI of a learner  $\mathcal{I}$  is estimated by

$$\widehat{\text{PFI}}_j^{\text{res}}(\mathcal{I}, \mathcal{D}, n_{\text{train}}) = \frac{1}{k} \sum_{i=1}^k \widehat{\text{PFI}}_j(\hat{f}_{\mathcal{D}_{\text{train}}^i}, \mathcal{D}_{\text{test}}^i), \quad (7)$$

where the permutation strategy is applied on the test sets  $\mathcal{D}_{\text{test}}^i$ .

### 3 Feature importance for groups

In our first minor contribution, we provide a general notation and formal definitions for grouped permutation and refitting methods and explain them by answering the following questions:

- To what extent does a group of features contribute to the model's performance in the presence of other groups?
- To what extent does a group itself increase the expected loss if it is added to a null model like the mean prediction of the target for refitting methods?
- How can we fairly distribute the expected loss among all groups and all features within a group?

The definitions of all grouped feature importance scores are based on loss functions. They are defined in such a way that important groups will yield positive grouped feature importance scores. The question of how to interpret the differing results of these methods is addressed in Sect. 5.

#### 3.1 Permutation methods

Here, we extend the existing definition of PFI to groups of features and introduce the GPFI (Grouped Permutation Feature Importance) and GOPFI (Group Only Permutation Feature Importance) scores. For ease of notation, we will only define these scores for a fitted model  $\hat{f}$  (see Eq. 5).

<sup>4</sup> An example for  $n = 3$  would be  $\tau_1 = (1, 3, 2)^\top$  with  $\tau_1^{(i)}$  being the  $i$ -th entry of that vector.

### 3.1.1 Grouped permutation feature importance (GPFI)

For the definition of GPFI—which is based on the definitions of Gregorutti et al. (2015) and Valentin et al. (2020)—let  $G \subset \{1, \dots, p\}$  be a group of features. Let  $\tilde{X}_G = (\tilde{X}_j)_{j \in G}$  be a  $|G|$ -dimensional random vector of features, which is an independent replication of  $X_G = (X_j)_{j \in G}$  following the same joint distribution. This random vector is independent of both the target variable and the random vector of the remaining features, which we define by  $X_{-G} := (X_j)_{j \in \{1, \dots, p\} \setminus G}$ . With slight abuse of notation to index the feature groups included in  $G$ , we define the grouped permutation feature importance of  $G$  as

$$\text{GPFI}_G = \mathbb{E}(L(\hat{f}(\tilde{X}_G, X_{-G}), Y)) - \mathbb{E}(L(\hat{f}(X), Y)). \tag{8}$$

Equation (8) extends Eq. (5) to groups of features so that the interpretation of GPFI scores always refers to the importance when the feature values of the group defined by  $G$  are permuted jointly (i.e., without destroying the dependencies of the features within the group). Similar to Eq. (7), the grouped permutation feature importance can be estimated by Monte Carlo integration:

$$\widehat{\text{GPFI}}_G = \frac{1}{nm} \sum_{i=1}^n \sum_{k=1}^m \left( L(\hat{f}(\mathbf{x}_G^{(\tau_k^{(i)})}, \mathbf{x}_{-G}^{(i)}), y^{(i)}) - L(\hat{f}(\mathbf{x}^{(i)}, y^{(i)})) \right). \tag{9}$$

The GPFI measures the contribution of one group to the model’s performance if all other groups are present in the model (see (a) from Sect. 3).

### 3.1.2 Group only permutation feature importance (GOPFI)

To evaluate the extent to which a group itself contributes to a model’s performance (see (b) from Sect. 3), one can also use a slightly different measure. As an alternative to Eq. 9, we can compare the expected loss after permuting all features jointly with the expected loss after permuting all features except the considered group. We define this GOPFI for a group  $G \subset \{1, \dots, p\}$  as

$$\text{GOPFI}_G = \mathbb{E}(L(\hat{f}(\tilde{X}), Y)) - \mathbb{E}(L(\hat{f}(X_G, \tilde{X}_{-G}), Y)), \tag{10}$$

which can be approximated by

$$\widehat{\text{GOPFI}}_G = \frac{1}{nm} \sum_{j=1}^n \sum_{k=1}^m \left( L(\hat{f}(\mathbf{x}^{(\tau_k^{(j)})}, y^{(j)})) - L(\hat{f}(\mathbf{x}_G^{(j)}, \mathbf{x}_{-G}^{(\tau_k^{(j)})}, y^{(j)})) \right). \tag{11}$$

While the relevance of GOPFI as an importance measure might be limited, it is technically useful for the grouped Shapley importance (see Eq. 14).

### 3.2 Refitting methods

Here, we introduce two refitting-based methods for groups of features. The first definition is similar to the one introduced in Williamson et al. (2020).

#### 3.2.1 Leave-one-group-out importance (LOGO)

For a subset  $G \subset \{1, \dots, p\}$ , we define the reduced dataset  $\tilde{\mathcal{D}} := \{(\mathbf{x}_{-G}^{(i)}, y^{(i)})\}_{i=1}^n$ . Given a learner  $\mathcal{I}$ , which generates models  $\mathcal{I}(\mathcal{D}) = \hat{f}_{\mathcal{D}}$  and  $\mathcal{I}(\tilde{\mathcal{D}}) = \hat{f}_{\tilde{\mathcal{D}}}$ , we define the Leave-One-Group-Out Importance (LOGO) as

$$LOGO(G) = \mathbb{E}(L(\hat{f}_{\tilde{\mathcal{D}}}(X_{-G}), Y)) - \mathbb{E}(L(\hat{f}_{\mathcal{D}}(X), Y)). \quad (12)$$

The LOGO can be estimated by using a learner  $\mathcal{I}$  on  $\tilde{\mathcal{D}}$  and should be embedded in a resampling technique:

$$\begin{aligned} \widehat{LOGO}(G) &= \widehat{GE}(\mathcal{I}, \tilde{\mathcal{D}}, n_{\text{train}}) - \widehat{GE}(\mathcal{I}, \mathcal{D}, n_{\text{train}}) \\ &= \frac{1}{k} \sum_{i=1}^k \hat{\rho}(\hat{f}_{\tilde{\mathcal{D}}_{\text{train}}^i}, \tilde{\mathcal{D}}_{\text{test}}^i) - \frac{1}{k} \sum_{i=1}^k \hat{\rho}(\hat{f}_{\mathcal{D}_{\text{train}}^i}, \mathcal{D}_{\text{test}}^i). \end{aligned}$$

Consequently, we compare the increase in expected loss compared to the full model's expected loss when leaving out a group of features and performing a refit (see (a) from Sect. 3).

While GPFI can be calculated with a resampling-based strategy by using refits to receive the algorithm-based instead of model-based GPFI, the meaning still varies from LOGO. For the algorithm-based GPFI, we calculate for each fitted model the importance score by permuting the regarded group and predicting with the same model. Then we average over all models from our resampling strategy and receive an importance score, which tells us how important a group of features is for some learner  $\mathcal{I}$  when we break the association between this group and all other groups and the target. LOGO, on the other hand, leaves the group out and then performs the refit to calculate the importance of the group, and hence, it addresses the question: Can we remove this group from our dataset without reducing our model's performance? This is not answered by permutation-based methods.

#### 3.2.2 Leave-one-group-in importance (LOGI)

While it may be too limiting to estimate the performance of a model based on one feature only, it can be informative to determine the extent to which a group of features (e.g., all measurements from a specific medical device) can reduce the expected loss in contrast to a null model (see (b) from Sect. 3). The Leave-One-Group-In (LOGI) method could be particularly helpful in settings where information on additional groups of measures will induce significant costs (e.g., adding functional imaging

data for a diagnosis) and/or limited resources are available (e.g., in order to be cost-covering, only one group of measures can be acquired). The LOGI method can also be useful for theory development in the natural and social sciences (e.g., which group of behaviors is most predictive by itself).

Let  $\mathcal{I}_{\text{null}}$  be a null algorithm, which results in a null model  $\hat{f}_{\text{null}}$  that only guesses the mean (or majority class for classification) of the target variable for any dataset. We additionally define a learner  $\mathcal{I}$ , which generates a model  $\mathcal{I}(\hat{\mathcal{D}}) = \hat{f}_{\hat{\mathcal{D}}}$  for a dataset  $\hat{\mathcal{D}} := \{(\mathbf{x}_G^{(i)}, y^{(i)})\}_{i=1}^n$ , which only contains features defined by  $G \subset \{1, \dots, p\}$ . We define the LOGI of a group  $G$  as

$$LOGI(G) = \mathbb{E}(L(\hat{f}_{\text{null}}, Y)) - \mathbb{E}(L(\hat{f}_{\hat{\mathcal{D}}}(X_G), Y)). \tag{13}$$

The LOGI can be estimated by using a learner  $\mathcal{I}$  on  $\hat{\mathcal{D}} = \{(\mathbf{x}_G^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^n$  and should be embedded in a resampling technique:

$$\begin{aligned} \widehat{LOGI}(G) &= \widehat{GE}(\mathcal{I}_{\text{null}}, \mathcal{D}, n_{\text{train}}) - \widehat{GE}(\mathcal{I}, \hat{\mathcal{D}}, n_{\text{train}}) \\ &= \frac{1}{k} \sum_{i=1}^k \hat{\rho}(\hat{f}_{\text{null}}, \mathcal{D}_{\text{test}}^i) - \frac{1}{k} \sum_{i=1}^k \hat{\rho}(\hat{f}_{\hat{\mathcal{D}}_{\text{train}}^i}, \hat{\mathcal{D}}_{\text{test}}^i). \end{aligned}$$

### 3.3 Grouped Shapley importance (GSI)

The importance measures defined above either exclude (or permute) individual groups of features from the total set of features or consider only the importance of groups by omitting (or permuting) all other features. The grouped importance scores are usually not affected if interactions within the groups are present. However, they can be affected if features from different groups interact, since permuting a group of features jointly destroys any interactions with other features outside the considered group. Therefore, we define the grouped Shapley importance (GSI) based on Shapley values (Shapley 1953). GSI scores account for feature interactions, as they measure the average contribution of a given group to all possible combinations of groups and fairly distribute the importance value caused by interactions among all groups (see (c) from Sect. 3).

We assume a set of distinct groups  $\mathcal{G} = \{G_1, \dots, G_l\}$ , with  $G_i \subset \{1, \dots, p\}$ , for  $i = 1, \dots, l$ . In our grouped feature context, the value function  $v : \mathcal{P}(\mathcal{G}) \rightarrow \mathbb{R}$  assigns a “payout” to each possible group or combination of groups included in  $\mathcal{G}$ . With slight abuse of notation, we define the value function for a subset  $S \subset \mathcal{G}$  as

$$v(S) := v(\cup_{G_i \in S} G_i).$$

We define the value function for a group  $G \in \mathcal{G}$  calculated by a refitting or a permutation method by

$$v_{\text{refit}}(G) = LOGI(G) \quad \text{or} \quad v_{\text{perm}}(G) = GOPFI(G), \tag{14}$$

respectively. The marginal contribution of a group  $G \in \mathcal{G}$ , with  $S \subset \mathcal{G}$  is

$$\Delta_G(S) = v(S \cup G) - v(S).$$

The GSI of the feature group  $G$  is then defined as

$$\phi(G) = \sum_{S \subset \mathcal{G} \setminus G} \frac{(|\mathcal{G}| - 1 - |S|)! \cdot |S|!}{|\mathcal{G}|!} \Delta_G(S), \quad (15)$$

which is a weighted average of marginal contributions to all possible combinations of groups.

The GSI cannot always be calculated in a time-efficient way, because the number of coalitions  $S \subset \mathcal{G} \setminus G$  can become large very quickly. In practice, the Shapley value is often approximated (Casalicchio et al. 2019; Covert et al. 2020) by drawing  $M \leq |\mathcal{G}|!$  different coalitions  $S \subset \mathcal{G} \setminus G$  and averaging the marginal, weighted contributions:

$$\hat{\phi}_M(G) = \frac{1}{M} \sum_{m=1}^M (|\mathcal{G}| - 1 - |S_m|)! \cdot |S_m|! \cdot \Delta_G(S_m), \quad (16)$$

with  $S_m \subset \mathcal{G} \setminus G$ , for all  $m = 1, \dots, M$ .

The GSI can in general not be exactly decomposed into the sum of the Shapley importances for single features of the regarded group. In Appendix B, we show that the remainder term  $R = \phi(G) - \sum_{i \in G} \phi(x_i)$  depends only on higher-order interaction effects between features of the regarded group and features of other groups. Hence, if one is interested in which features contributed most within a group, the Shapley importances for single features can be calculated, which provide a fair distribution of feature interactions within the group but not necessarily of feature interactions across groups. However, the remainder term can be used as a quantification of learned higher-order interaction effects between features of different groups.

While the GSI can be calculated with permutation- as well as refitting-based approaches, we will only apply the permutation-based approach in the upcoming simulation studies and the real-world example.

### 3.4 Real world use cases

For each dataset from Sect. 1.1, we fitted a random forest and summarized the three most important groups according to different grouped feature importance methods. For the importance scores of LOGI and LOGO, we used a 10-fold cross-validation (Table 2).

For the *birthweight* task, the feature **lwt** (mother's weight) was the most important group to predict the birthweight for all grouped feature importance methods except for LOGI. While all methods except LOGI also agree on the second most important group **ui** (presence of uterine irritability), feature groups differ for the third rank. However, this may also be due to statistical variability, as the importance values become very

**Table 2** Best 3 groups for each grouped feature importance score

Dataset	GPMI	GPMI	GSI	LOGI	LOGO
<i>Birthweight</i>	lwt (0.067)	lwt (0.056)	lwt (0.062)	ui (0.041)	lwt (0.036)
	ui (0.056)	ui (0.047)	ui (0.046)	Race (0.017)	ui (0.029)
	Smoke (0.009)	Race (0.045)	ptl (0.019)	ptl (0.015)	Race (0.005)
<i>Colon</i>	Gene14 (0.143)	Gene14 (0.174)	Gene14 (0.125)	Gene14 (0.128)	Gene14 (0.131)
	Gene10 (0.007)	Gene16 (0.087)	Gene16 (0.042)	Gene20 (0.045)	Gene17 (0.036)
	Gene7 (0.001)	Gene12 (0.057)	Gene13 (0.019)	Gene13 (0.028)	Gene18 (0.033)
<i>Fuelsubset</i>	NIR (30.51)	NIR (42.20)	NIR (36.21)	NIR (27.35)	NIR (8.34)
	UVVIS (2.85)	UVVIS (14.38)	UVVIS (7.99)	UVVIS (15.74)	H <sub>2</sub> O (0.14)
	H <sub>2</sub> O (0.01)	H <sub>2</sub> O (1.26)	H <sub>2</sub> O (0.24)	H <sub>2</sub> O (− 12.17)	UVVIS (− 2.14)

For the classification task (*colon*) the scores were calculated as differences in classification accuracy. For the other two regression tasks the scores result from differences in MSE

small. It is interesting that **lwt**, despite being the most important group for all other scores, is not very important in terms of LOGI. Thus, **lwt** is less important as a stand-alone group, but appears important if the other feature groups are included in the model.

In the *colon* task, the feature group *gene14* is by far the most important group to predict cancerous tissue for all grouped feature important methods. However, there are variations in the second and third most important groups.

For the *fuelsubset* task, the permutation-based grouped importance methods (GPMI, GOPMI and GSI) show the same importance ranking for the three most important feature groups. However, for the refitting-based grouped importance methods (LOGI and LOGO), we can observe interesting differences. The features from the *UVVIS* group are important as a stand-alone group as can be seen by their positive LOGI score. However, the negative LOGO score of the *UVVIS* group indicates that the algorithm seems to perform better with only the *NIR* and *H<sub>2</sub>O* groups.

GPMI, GOPMI and GSI provide importance scores for feature groups of a given trained model without the necessity to refit the model. In contrast, LOGI and LOGO provide grouped importance scores based on the underlying algorithm and should always be considered together.

## 4 Sequential grouped feature importance

In general, feature groups do not necessarily have to be distinct or independent of each other. When groups partly contain the same or highly correlated features, we may obtain high grouped feature importance scores for similar groups. This can lead to misleading conclusions regarding the importance of groups. Quantifying the importance of different combinations of groups is especially relevant in applications where extra costs are associated with using additional features from other data sources. In this case, one might be interested in the sparsest, yet most important combination of groups or in understanding the interplay of different combinations of groups. Hence,

in practical settings, it is often important to decide which additional group of features to make available (e.g., buy or implement) for modeling and how groups should be prioritized under economic considerations.

Gregorutti et al. (2015) introduced a method called *grouped variable selection*, which is an adaptation of the recursive feature elimination algorithm from Guyon et al. (2002) and uses permutation-based grouped feature importance scores for the selection of feature groups. In Algorithm 1, we introduce a sequential procedure that is based on the idea of stability selection (Meinshausen and Bühlmann 2010). The procedure primarily aims at understanding the interplay of different combinations of groups by analyzing how the importance scores change after including other groups in a sequential manner. The feature groups must be pre-specified by the user. We prefer a refitting-based over a permutation-based grouped feature importance score when the secondary goal is to find well-performing combinations of groups. Here, the fundamental idea is to start with an empty set of features and to sequentially add the next best group in terms of LOGI until no further substantial improvement can be achieved. Our sequential procedure is based on a greedy forward search and creates an implicit ranking by showing the order in which feature groups are added to the model. To account for the variability introduced by the model, we propose to use repeated subsampling or bootstrap with sufficient repetitions (e.g., 100 repetitions).

To better understand Algorithm 1, we will demonstrate it with a small example with four groups  $\mathcal{G} = \{G_1, G_2, G_3, G_4\}$  here. As a reminder, each group is a subset of  $\{1, \dots, p\}$ , and we want to find a subset  $B \subset \{1, \dots, p\}$ , which consists of the union of groups in  $\mathcal{G}$ . The subset  $B$  is found by our sequential grouped feature importance procedure. To account for variability, the whole dataset is split into two sets (training and test set) repeatedly so that the train-test splits are different in each repetition of the resampling strategy (bootstrap or subsampling). For each training set, Algorithm 1 starts with an empty set  $B = \emptyset$  (line 2, Algorithm 1). In line 5 of Algorithm 1, the candidate set  $\mathcal{B} \subset \mathcal{P}(\mathcal{G})$  is defined as all subsets of the power set with cardinality 1. These are all individual groups  $\mathcal{B} = \{\{G_1\}, \{G_2\}, \{G_3\}, \{G_4\}\}$ . The LOGI score of each single group is then calculated. In our example, let  $G_1$  have the highest LOGI score, which also exceeds the threshold  $\delta$ . The desired combination  $B$  is preliminarily defined as  $G_1$  (line 8), and for the comparison in the next step, the LOGI score of  $G_1$  is defined as  $L_0$  (line 9). Then, a new candidate set  $\mathcal{B}$  is defined (line 11), which consists of all subsets of the power set of  $\mathcal{G}$  of size  $i$  (at this step, we have  $i = 2$ ), where  $B = G_1$  is also a subset of  $\mathcal{B}$ . Hence,  $\mathcal{B} := \{\{G_1, G_2\}, \{G_1, G_3\}, \{G_1, G_4\}\}$ . The LOGI score of elements of  $\mathcal{B}$  is calculated as the LOGI score of the union of all subsets. Now, let  $\widehat{LOGI}(G_1 \cup G_3)$  have the highest score. This score is compared to the LOGI score of the previous iteration  $L_0$  (line 13). Let the difference exceed the threshold  $\delta$  for our example. In line 14 and 15, the desired combination  $B$  is now defined as  $G_1 \cup G_3$  and the LOGI score is again defined as  $L_1$ . Algorithm 1 now jumps to line 10 again with  $i = 3$ . The candidate set is now  $\mathcal{B} = \{\{G_1, G_3, G_2\}, \{G_1, G_3, G_4\}\}$  (line 11). The LOGI scores are now calculated again for each element of  $\mathcal{B}$ . Let no LOGI score exceed  $L_0$  by the threshold  $\delta$  (line 13). Algorithm 1 now ends for this dataset split and returns  $B = G_1 \cup G_3$  as the best combination. This procedure is repeated for each train-test split in each repetition.



**Algorithm 1:** Sequential Grouped Feature Importance

```

input : Set of groups  $\mathcal{G} = \{G_1, \dots, G_k\}$ .
          Improvement threshold  $\delta > 0$ .
          Number of repetitions for the data splitting.
output: For every data split: a combination  $B \subset \{1, \dots, p\}$  and the order in which feature groups
          were added.
1 for Every outer data split do
2   Let  $B = \emptyset$  for  $i = 1, \dots, k$  do
3     if  $i = 1$  then
4       Define candidate set  $\tilde{B} := \{\tilde{G} \in \mathcal{P}(\mathcal{G}) \mid |\tilde{G}| = 1\}$ 
5       Find best single group  $G^* = \underset{\tilde{G} \in \tilde{B}}{\arg \max} (\widehat{LOGI}(\tilde{G}))$ 
6       if  $\widehat{LOGI}(G^*) > \delta$  then
7          $B = G^*$ 
8          $L_{i-1} = \widehat{LOGI}(B)$ 
9       if  $i > 1$  and  $B \neq \emptyset$  then
10        Define candidate set  $\tilde{B} := \{\tilde{G} \in \mathcal{P}(\mathcal{G}) \mid |\tilde{G}| = i \text{ and } B \subset \tilde{G}\}$ 
11        Find best combination  $G^* = \underset{\tilde{G} \in \tilde{B}}{\arg \max} (\widehat{LOGI}(\bigcup_{G' \in \tilde{G}} G'))$ 
12        if  $\widehat{LOGI}(\bigcup_{G' \in G^*} G') - L_{i-1} > \delta$  then
13           $B = \bigcup_{G' \in G^*} G'$ 
14           $L_{i-1} = \widehat{LOGI}(B)$ 
15        else
16          break for loop
17
18

```

Since the order in which feature groups are added is also known, alluvial charts (Allaire et al. 2017) can be created for visualization purposes (see Figs. 2 and 10). In these charts, we included the number of times feature groups were added as well as the performance on the test datasets. These charts show how frequently a group was selected given that another group was already included and thereby highlight robust combinations of groups.

## 5 Comparison of grouped feature importance methods

After introducing the methodological background of the different loss-based grouped feature importance measures in Sect. 3, we will now compare them in different simulation settings. We analyze the impact on all methods for settings where (1) groups are dependent, (2) correlations within groups vary, and (3) group sizes differ.

### 5.1 Dependencies between groups and sparsity

In this section, we compare refitting- and permutation-based grouped feature importance methods and show how different dependencies between groups can influence the importance scores. We demonstrate the benefits of the sequential grouped feature importance procedure and conclude with a recommendation of when to use refitting or permutation-based methods depending on the use-case.

We simulate a data matrix  $\mathbf{X}$  with  $n = 1000$  instances and 3 groups  $G_1, G_2, G_3$ , with each of them containing 10 normally distributed features. Features are simulated in such a way that features within each group are highly correlated. However, features in  $G_3$  are independent of features in  $G_1$  and  $G_2$ , while features in  $G_1$  and  $G_2$  are also highly correlated with each other. To generate normally distributed features with such correlation patterns, we follow the approach of Toloşi and Lengauer (2011) and use prototype vectors in the following way: (1) We draw  $n$  instances of the prototype vector  $\mathbf{U} \sim \mathcal{N}(0, 1)$ . (2) We generate features in  $G_1$  by adding a normally distributed error term  $\epsilon \sim \mathcal{N}(0, 0.5)$  to 10% of the instances of the prototype vector  $\mathbf{U}$ . (3) Features in  $G_2$  are generated by copying features of  $G_1$  and adding a small normally distributed error term  $\epsilon \sim \mathcal{N}(0, 0.01)$  to the copied features. It follows that features within  $G_1$  and  $G_2$  as well as features between the two groups are highly correlated. (4) We generate a new prototype vector  $\mathbf{V}$ , which is independent of  $\mathbf{U}$ . (5) We generate features for  $G_3$  in the same way as done for  $G_1$  in step (2) but with the prototype vector  $\mathbf{V}$ .

The target vector  $\mathbf{Y}$  is generated by  $\mathbf{Y} = 2\mathbf{U} + 1\mathbf{V} + \epsilon$ , with  $\epsilon \sim \mathcal{N}(0, 0.1)$ . We fitted a support vector machine with a radial basis function kernel<sup>5</sup>, as an example of a black-box algorithm.

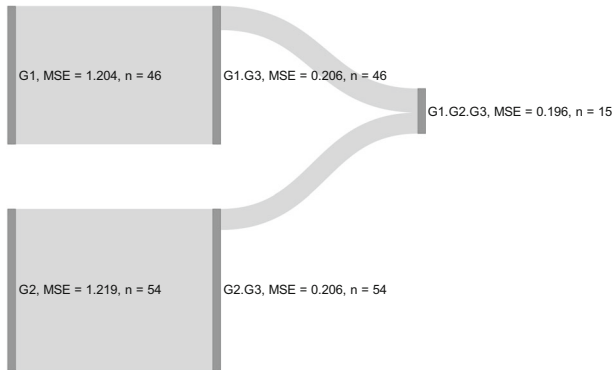
The results in Table 3 show that there can be major differences depending on how the grouped feature importance is calculated. Permutation methods (GOPFI & GPFI & GSI) reflect the importance of the groups based on a model trained on a fixed dataset. In contrast, refitting methods (LOGI & LOGO) retrain the model on a reduced dataset and can therefore learn new relationships. Looking at the results from the permutation methods, we can see that the groups  $G_1$  and  $G_2$  are approximately equally important while both being more important than  $G_3$ . However, the results from the refitting methods can reveal some interesting relationships between the groups. The refitting methods highlight that  $G_1$  and  $G_2$  are more or less interchangeable if we only consider a performance-based interpretation (which might not coincide with a domain-specific

<sup>5</sup> Epsilon regression,  $\epsilon = 0.1, C = 1$  with heuristically chosen kernel width according to (Caputo et al. 2002) (here:  $\sigma = 0.079$ ).

**Table 3** Results of different feature importance calculations of the simulation

Group	GOPFI	GPMFI	GSI	LOGI	LOGO
$G_1$	6.04 ( $\pm 0.37$ )	2.64 ( $\pm 0.07$ )	4.12 ( $\pm 0.45$ )	3.93 ( $\pm 0.75$ )	-0.01 ( $\pm 0.02$ )
$G_2$	5.90 ( $\pm 0.35$ )	2.57 ( $\pm 0.09$ )	4.01 ( $\pm 0.47$ )	3.93 ( $\pm 0.76$ )	-0.00 ( $\pm 0.02$ )
$G_3$	1.76 ( $\pm 0.39$ )	1.75 ( $\pm 0.05$ )	1.54 ( $\pm 0.39$ )	0.58 ( $\pm 1.01$ )	1.01 ( $\pm 0.22$ )

GSI scores were calculated without approximation, with  $v_{perm}$  as value function (see Eq. 14). All results were averaged by a 10-fold cross-validation scheme, with standard deviations reported in parentheses



**Fig. 2** Sequential grouped feature importance for the simulation in Sect. 5.1. 100 times repeated subsampling. Improvement threshold  $\delta = 0.001$ . Vertical bars show one step of the sequential procedure (left to right). Height of the vertical bars represent the number of subsampling iterations that a combination of groups was chosen. *MSE* scores show predictive performance. Streams represent the addition of a group

perspective)<sup>6</sup>. Hence, the two groups do not complement each other. This is reflected by the near-zero LOGO scores, which indicate that leaving each group out of the full model does not considerably change the model’s expected loss.

Figure 2 illustrates the results of the sequential procedure introduced in Algorithm 1. We see that across 100 subsampling iterations,  $G_1$  was chosen 46 times as the most important first group, and  $G_2$  was chosen 54 times with similar predictive performance for both groups, while  $G_3$  was never chosen as the first most important group. Hence, similar to LOGI, we can see that if only one group can be chosen, it would either be  $G_1$  or  $G_2$  with approximately the same probability. In the second step, the group  $G_3$  was added in all cases to either  $G_1$  or  $G_2$  (depending on which group had been chosen in the first step). This step resulted in an on-average drop in the MSE score from 1.2 to 0.2. In only a few cases (15 out of 100), the final addition of either  $G_1$  or  $G_2$  to a full model in step 3 exceeded the very low chosen threshold of  $\delta = 0.001$ . This rather unlikely improvement is represented by the proportionally narrower band that connects the second and the third step (dark gray bars) in the chart in Fig. 2. This reveals that these two groups are—from a performance or loss perspective—rather interchangeable and do not benefit from one another.

<sup>6</sup> It is possible that adding a group of features to the model might not lead to a better model performance, but the group may still be relevant due to the domain-specific context. However, this depends on the regarded use case. All our interpretations here are purely statistical.

The choice between using permutation-based or refitting-based grouped feature importance methods might depend on the number of groups and correlation strength between the different groups. If feature groups are distinct and features between the groups are almost uncorrelated, we might prefer permutation over refitting methods due to lower computation time. In cases where groups are correlated with each other (e.g., because some features belong to multiple groups), refitting methods might be preferable, as they are not misleading in correlated settings. Since the number of groups is usually smaller than the number of features in a dataset, refitting methods for groups of features could become a viable choice. Furthermore, with the sequential grouped feature importance procedure, it is possible to find sparse and well performing combinations of groups in an interpretable manner. Thus, this approach helps to better understand which groups of features were important (e.g., as they were more frequently selected) given that certain groups were already selected.

## 5.2 Varying correlations within groups

In many use cases, it is quite common to group similar (and therefore, often correlated) features together, while groups of features may be almost independent of each other. However, compared to Sect. 5.1, correlations of features within groups might differ. We created a data matrix  $\mathbf{X}$  with  $n = 1000$  instances and 4 groups  $G_1$ ,  $G_2$ ,  $G_3$ , and  $G_4$ , with each of these groups containing 10 normally distributed features. Using fivefold cross-validation, we fitted a random forest with 2000 trees and a support vector regression with a radial basis function kernel.<sup>7</sup> The univariate target vector  $\mathbf{Y}$  is defined as follows:

$$\mathbf{Z}_j = 3\mathbf{X}_{G_j,3}^2 - 4\mathbf{X}_{G_j,5} - 6\mathbf{X}_{G_j,7} + 5\mathbf{X}_{G_j,9} \cdot d_j, \quad j \in \{1, 2, 3\}$$

$$\mathbf{Y} = \sum_{j=1}^3 \mathbf{Z}_j + \epsilon$$

with

$$d_j = \begin{cases} 1, & \text{if } \text{mean}(\mathbf{X}_{G_j,8}) > 0 \\ 0, & \text{otherwise} \end{cases}$$

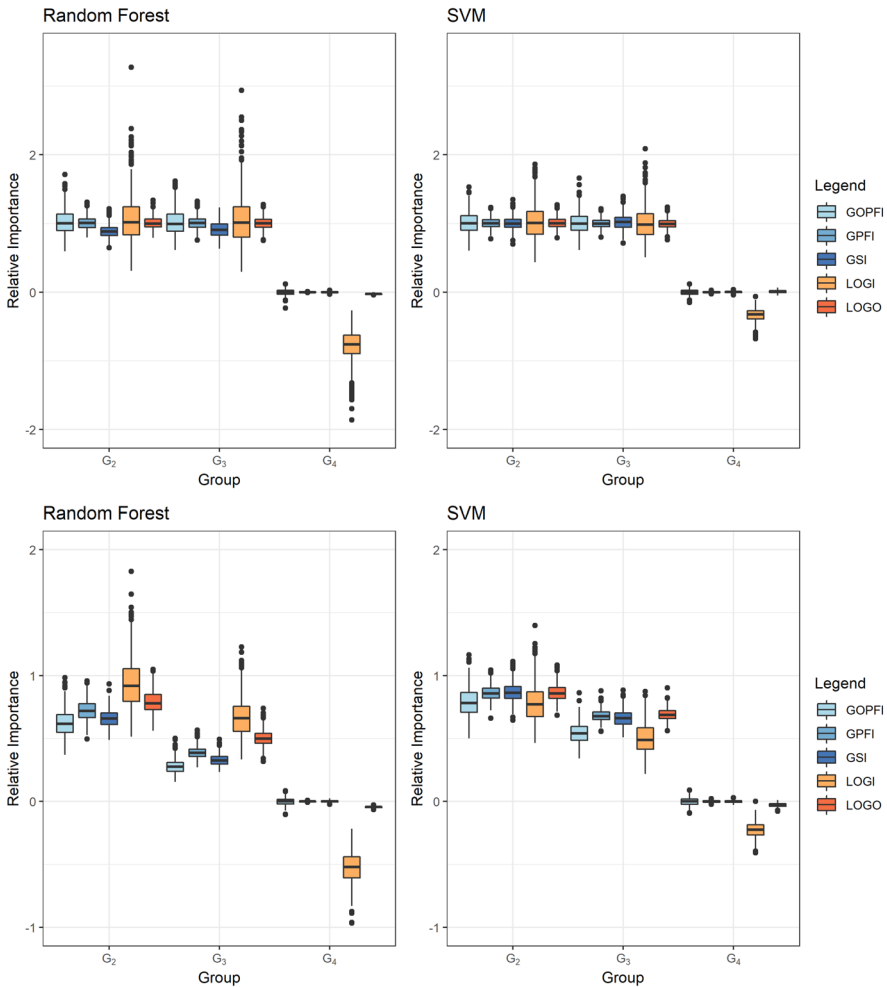
and  $\epsilon \stackrel{iid}{\sim} N(0, 1)$ . The  $i$ -th feature of the  $j$ -th group is denoted by  $\mathbf{X}_{G_j,i}$ . We repeated the simulation 500 times.

It follows that  $G_1$ ,  $G_2$ , and  $G_3$  have the same influence on the target variable, while  $G_4$  has no influence on  $\mathbf{Y}$ . We generate the feature space  $\mathbf{X}$ —similar to the approach in Sect. 5.1—as follows: (1) For each feature group  $j$ , we generate a prototype vector  $\mathbf{U}_j \sim \mathcal{N}(0, 1)$  with  $n$  instances. (2) We generate the features of a group  $G_j$  by altering a proportion  $\alpha$  with  $0 \leq \alpha \leq 1$  of the  $n$  instances of  $\mathbf{U}_j$ . We alter these instances by taking

<sup>7</sup> We used a cost parameter of  $C = 1$  and estimate the kernel width based on the heuristic introduced by Caputo et al. (2002)

a weighted average between the respective values of  $U_j$  (20%) and a standard normally distributed random variable  $W_i$  (80%). For the results shown in Fig. 3, we set  $\alpha$  to 0.1 for all features within the same group. Hence, correlations within groups are the same (around 90%) for all groups, while groups themselves are independent of each other. The plots show that all methods correctly attribute the same importance to the first three groups, while the fourth group is not important for predicting  $Y$ . The lower plots in Fig. 3, on the other hand, correlations within groups vary across groups. The altering proportion parameter  $\alpha$  is set to 0.1 for features of  $G_1$  and  $G_4$ , to 0.3 for features of  $G_2$ , and to 0.6 for features of  $G_3$ . Hence, features in  $G_1$  and  $G_4$  are highly correlated within the respective group, while features within  $G_2$  and  $G_3$  show a medium and small correlation, respectively. While  $G_4$  is still recognized to be unimportant, the relative importance of groups 1 to 3 drops with decreasing within-group correlation. This artifact seems—at least, in this simulation setting—to be even more severe for the random forest compared to the support vector machine. For example,  $G_3$  is on average less than half as important as  $G_1$  for permutation-based methods. Thus, none of the methods reflect the true importance of the different groups of the underlying data generating process. A possible reason for this artifact is that the regarded model learned effects different from those given by the underlying true relationship. Especially for the random forest, this has already been studied extensively in the presence of different correlation patterns in the feature space (Strobl et al. 2008; Nicodemus et al. 2010). Additionally, Hooker and Mentch (2019) showed that permutation-based methods are more sensitive in this case than refitting methods, which is also visible for both models in Fig. 3. Since the model is learned on the original feature space and group structures are not considered in the modelling process, we can also observe this effect when applying grouped feature importance methods. This is due to the fact that we can only quantify which groups are important for the model or algorithm performance but not for the underlying data generating process, which is usually unknown. Another approach to quantify feature importance when using random forests is to extract the information on how often a feature has been used as a splitting variable for the different trees. The feature chosen for the first split has the most influence within each tree. Hence, we calculated for each repetition the percentage of how often a feature is chosen as the first splitting feature. The distribution over all repetitions is displayed in Fig. 4. Each of the features of  $G_1$  is on average chosen more often as the first splitting feature than all features of the other groups, no matter if it has an influence on the target or not. The influential features of  $G_3$  (which has the lowest within-group correlation) are rarely chosen as the first splitting feature. This observation confirms the results of the grouped importance methods in Fig. 3, since all of them rank  $G_3$  as least important from the influential feature groups.

Note that while GPFI and LOGO are calculated with reference to the full model's performance—which on average leads to higher absolute values than the two counter-methods based on the null model's performance—GOPFI and LOGI might lead to less robust results, as the newly learned effects as well as the approximation of the permutation effect underlie a higher uncertainty. This effect might increase when relative values instead of absolute values are considered due to smaller absolute importance scores of GOPFI and LOGI. However, the methods are only comparable on a relative scale. This effect is also visible in the boxplots of Fig. 3. Furthermore, LOGI can also

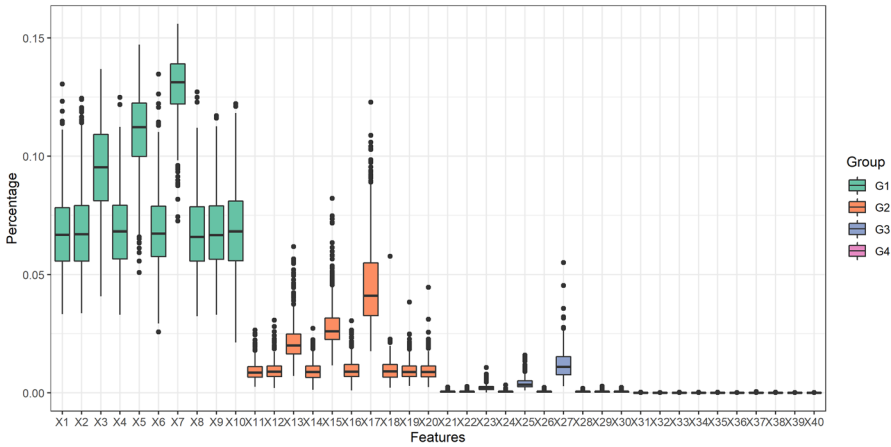


**Fig. 3** Upper (lower) plots: Grouped relative importance scores in the case of *equally sized* (*varying sizes*) within-group correlations for random forest (left) and SVM (right). Relative importance is calculated by dividing each of the absolute group importance scores by the importance score of  $G_2$ . Hence, the relative importance of  $G_1$  is 1. Boxplots illustrate the variation between different repetitions

take negative values in the case of  $G_4$ , as the feature group does not affect the target in the underlying data generating process, and hence it might be counterproductive to only include  $G_4$  compared to the null model.

### 5.3 Varying sizes of groups

Another factor to consider when calculating grouped rather than individual feature importance scores is that differing group sizes might influence the ranking of the scores. Groups with more features might often have higher grouped importance scores and

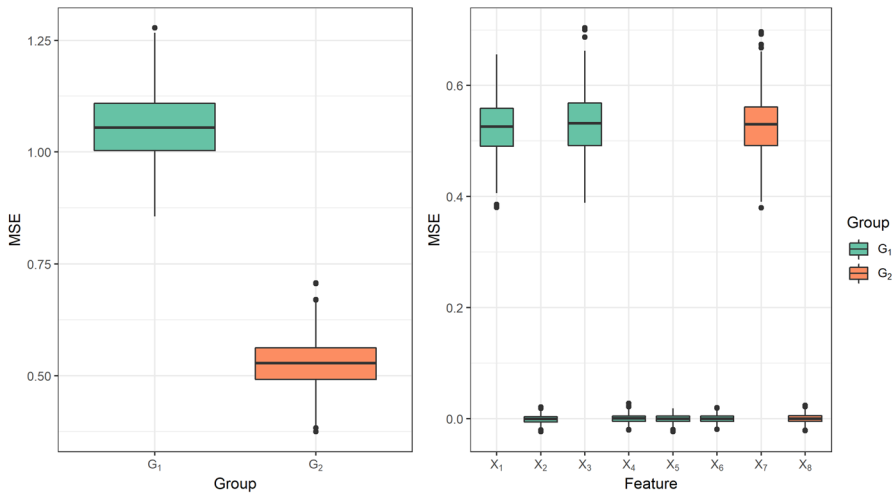


**Fig. 4** Percentage of how often each feature is chosen as the first splitting feature within the trained random forests. Results have been averaged over the cross-validation folds for each repetition. Boxplots show the distribution over all 500 repetitions

might contain more noise features than smaller groups. Therefore, Gregorutti et al. (2015) argue that in case one must decide between two groups that have an equal importance score, one would prefer the group with fewer features. Following from that, they normalize the grouped feature importance scores regarding the group size with the factor  $|G|^{-1}$ . This is also used in the default definition of the grouped model reliance score in Valentin et al. (2020). However, the usefulness of normalization highly depends on the question the user would like to answer. This is illustrated in a simulation example in Fig. 5. We created a data matrix  $\mathbf{X}$  with  $n = 2000$  instances and 2 groups, with  $G_1$  containing  $\{x_1, \dots, x_6\}$  and  $G_2$  containing  $\{x_7, x_8\}$  i.i.d. uniformly distributed features on the interval  $[0, 1]$ . The univariate target variable  $\mathbf{Y}$  is defined as follows:

$$\mathbf{Y} = 2\mathbf{X}_1 + 2\mathbf{X}_3 + 2\mathbf{X}_7 + \epsilon, \quad \text{with } \epsilon \stackrel{iid}{\sim} N(0, 1).$$

We used 1000 observations for fitting a random forest with 2000 trees and 1000 observations for prediction and calculating the GSI as defined in Sect. 3.3 with a permutation-based value function. This was repeated 500 times. Figure 5 shows that  $G_1$  is about twice as important as  $G_2$ . As shown in Sect. 3.3 and Appendix B, we can compare the GSI with the Shapley importance on feature level. In case there are no higher-order interaction terms between groups modeled by the random forest, the single feature importance scores will approximately sum up to the grouped importance score, as shown in this example. This provides a more detailed view of how many and which features are important within each group. In this case, there are two equally important features in  $G_1$  and one equally important feature in  $G_2$ . If we use the normalization constant in this example, we would divide the grouped importance score of  $G_1$  (which is on average approximately 1.1) by 6 and the one of  $G_2$  (which is on average approximately 0.55) by 2. Consequently,  $G_2$  with a normalized score of



**Fig. 5** Shapley importance on group (left) and on feature level (right). Boxplots show the variation between the 500 repetitions of the experiment

approximately 0.27 would be regarded as more important than  $G_1$  with a normalized score of approximately 0.18. It follows that if we must decide between two groups, we would choose  $G_2$  when we follow the approach of Gregorutti et al. (2015). However, since  $G_1$  contains two features with the same importance as the one important feature of  $G_2$ , and hence  $G_1$  contains more information from a statistical perspective, the user might prefer  $G_1$ . Furthermore, breaking down the GSI to the single-feature Shapley importance scores puts the user in the position of defining sparser groups by excluding non-influential features.

Finally, Table 4 presents a summary of the key takeaways regarding all discussed grouped feature importance methods.

## 6 Feature effects for groups

Feature effect methods quantify or visualize the influence of features on the model's prediction. For a linear regression model, we can easily summarize the feature effect in one number, thus making interpretation very simple: If we change feature  $x_1$  by one unit, our prediction will change by the corresponding coefficient estimate  $\hat{\beta}_1$  (positively or negatively depending on the sign of the coefficient). For more complex non-linear models like generalized additive models, such a simplified summary of the feature effect is not adequate, as the magnitude and sign of the effect might change over the feature's value range. Hence, it is more common to visualize the marginal effect of the feature of interest on the predicted outcome. Since ML models are often complex non-linear models, different visualization techniques for the feature effect have been introduced in recent years. Common methods are PDP, ICE curves or ALE (Friedman 2001; Goldstein et al. 2013; Apley and Zhu 2019), which show how changes in the feature values affect the predictions of the model. However, these are



**Table 4** Overview of pros and cons of the grouped feature importance methods

Criteria	GOPFI & GPFI	GSI	LOGI & LOGO
Time efficient	Yes (in comparison to alternatives)	Depends on number of groups	Depends on number of groups
Dependencies between groups (Sect. 5.1)	No full picture	No full picture	More insights than permutation-based if regarded together
Identify well performing combinations of groups (Sect. 5.1)	Not in general	Not in general	Only LOGI within Algorithm 1
Correlations within groups but independence between groups (Sect. 5.2)	Depends on learned effects of the model, less problematic if within group correlations do not differ strongly between groups	Depends on learned effects of the model, less problematic if within group correlations do not differ strongly between groups	More robust than permutation-based methods but still dependent on learned effects
Drilldown of grouped importance score on feature level (Sect. 5.3)	No	Yes (approximately depending on the influence of higher-order interactions)	No

While GOPFI is less relevant on its own, LOGI can provide insightful interpretations, e.g., if feature groups are correlated with each other or when used within the sequential procedure introduced in Sect. 4. The sequential procedure is the only method that can identify well performing and sparse combination of groups. Note that GSI is only evaluated w.r.t. a permutation-based calculation

usually only defined for a maximum of two features. For larger groups of features, this becomes more challenging, since it is difficult to visualize the influence of several features simultaneously. The approach described in this section aims to create effect plots for a predefined group of features that have an interpretation similar to that of the single-feature PDP. To achieve this, we transform the high-dimensional space of the feature group into a low-dimensional space by using a supervised dimension reduction method, which is discussed in Sect. 6.1. We want to find a few underlying factors that are attributed to a sparse and interpretable combination of features that explain the effect of the regarded group on the model's expected loss. We provide a detailed description of this method in Sect. 6.3 and introduce the resulting combined features effect plot (CFEP). In Sect. 6.4, we illustrate the advantages of applying a supervised rather than an unsupervised dimension reduction method and compare our method to the main competitor, which is the totalvis effect plot introduced in Seedorff and Brown (2021).

## 6.1 Choice of dimension reduction method

The most prominent dimension reduction technique is arguably PCA (Jolliffe 1986). PCA is restricted to explaining most of the variance of the feature space, and the identified projections are not related to the target variable (for more details see Appendix

C.1). Because we want to visualize the mean prediction of combined features as a result of the dimension reduction process, we prefer supervised procedures that maximize dependencies between the projected data  $\mathbf{XV}$ —with  $\mathbf{V}$  being a projection  $\mathbf{V} \in \mathbb{R}^{p \times p}$ —and the target vector  $\mathbf{Y}$  (as we show in Sect. 6.4). Many methods for supervised PCA have been established. For example, see Bair et al. (2006), who used a subset of features that were selected based on their linear correlation with the target variable. Another very popular method that maximizes the covariance between features and the target variable is partial least squares (PLS) (Wold et al. 1984). The main difference between these methods and the supervised PCA (SPCA) introduced by Barshan et al. (2011) is that the SPCA is based on a more general measure of dependence, called the Hilbert-Schmidt Independence Criterion (HSIC). This independence measure is constructed to be zero, if and only if any bounded continuous function between the feature and target space is uncorrelated. In practice, an empirical version of the HSIC criterion is calculated with kernel matrices. It follows that while this SPCA technique can cover a variety of linear and non-linear dependencies between  $\mathbf{X}$  and  $\mathbf{Y}$  by choosing an appropriate kernel, the other suggested methods are only able to model linear dependencies between the features and the target variable. The approach that is probably best suited for our application of finding *interpretable* sets of features in a high-dimensional dataset is the method called sparse SPCA, described in Sharifzadeh et al. (2017). Similar to the SPCA method from Barshan et al. (2011), sparse SPCA not only uses the HSIC criterion to maximize the dependency between projected data  $\mathbf{XV}$  and the target  $\mathbf{Y}$ , but also incorporates an  $L_1$  penalty of the projection  $\mathbf{V}$  for sparsity. The sparse SPCA problem can be solved with a *penalized matrix decomposition* (Witten et al. 2009). More theoretical details on the sparse SPCA, including the HSIC criterion and how it can be calculated empirically, and the choice of kernels and hyperparameters can be found in Appendix C.

## 6.2 Totalvis effect plot

Seedorff and Brown (2021) recently introduced a method that aims to plot the combined effect of multiple features by using PCA. Their approach can be described as follows: First, they apply PCA on the regarded feature space to receive the principal components matrix after rotation. For the principal component of interest, they create an equidistant grid. Second, for each grid value, they replace all values of the selected principal component with this grid value and transform the matrix back to the original feature space. Third, The ML model is applied on these feature values and a mean prediction for the grid point of the regarded principal component is calculated. Steps 2 and 3 are repeated for all grid points.

Hence, with this method, combined effect plots for up to  $p$  principal components can be created. Thus, Seedorff and Brown (2021) do not focus on explaining groups of features explicitly. Furthermore, they use PCA for unsupervised dimension reduction, and thus, projections might not be related to the target. Due to using PCA and not sparse PCA, the results might be difficult to interpret, as many or all features might have an influence on the principal component. Lastly, with the back-transformation from the principal component matrix to the original feature space, all feature values

change and might not be meaningful anymore. For example, in the case of integer features, the back-transformation might lead to real feature values. We illustrate the drawbacks of the method compared to the CFEP in Sect. 6.4.

### 6.3 Combined features effect plot (CFEP)

The CFEP picks up the idea of PDPs (Friedman 2001) and extends it to groups of features. The partial dependence function is defined as

$$f_S^{PD}(\mathbf{x}_S) = \mathbb{E}_{X_C}[\hat{f}(\mathbf{x}_S, X_C)] \tag{17}$$

with  $S \subset \{1, \dots, p\}$  and  $C = \{1, \dots, p\} \setminus S$ . Since the joint distribution of  $X_C$  is usually unknown, the Monte Carlo method is used to estimate  $f_S^{PD}(\mathbf{x}_S)$ :

$$\hat{f}_S^{PD}(\mathbf{x}_S) = \frac{1}{n} \sum_{i=1}^n \hat{f}(\mathbf{x}_S, \mathbf{x}_C^{(i)}) \tag{18}$$

Hence, we marginalize over all features in  $C$  and with that we obtain the average marginal effect for the feature subset in  $S$ . The PDP usually visualizes this average marginal effect for  $|S| \leq 2$  by plotting  $(\mathbf{x}_S^{(k)}, \hat{f}_S^{PD}(\mathbf{x}_S^{(k)}))$  for some pre-specified grid points  $k = \{1, \dots, m\}$ .<sup>8</sup> However, this is usually only possible for  $|S| \leq 2$  and thus not directly applicable to visualize the combined effect of feature groups. To obtain a visualization in the case of  $|S| > 2$ , we need to reduce the dimensions and therefore define the CFEP of a certain group of features  $G$  as follows:

- (1) We first apply a suitable (preferably supervised) dimension reduction method (e.g., here we use the sparse SPCA, however, the CFEP follows a modular approach and hence the dimension reduction method is exchangeable) on features in  $G \subset \{1, \dots, p\}$  to obtain a low dimensional representation of the feature group  $G$ . We denote these principle component functions—which are ordered according to relevance<sup>9</sup> and which possibly depend on a reduced set of features<sup>10</sup>  $S_j \subseteq G$  with  $j \in \{1, \dots, |G|\}$ —by  $g_j : \mathcal{X}_{S_j} \rightarrow \mathbb{R}$ .
- (2) For visualization purposes, we choose from all possible  $g_j$  with  $j \in \{1, \dots, |G|\}$  a principle component function

$$g : \mathcal{X}_S \rightarrow \mathbb{R} \tag{19}$$

(with  $S$  being its reduced set of features) which serves as a proxy for the feature group  $G$ . We usually only consider the first few principle components.

<sup>8</sup> For example, by using an equidistant grid or a random sample of values of  $\mathbf{x}_S$ .

<sup>9</sup> The relevance is defined by the objective that is optimized by the dimension reduction method. For sparse SPCA this is the HSIC criterion (see also Appendix C) and for PCA it is the explained variance.

<sup>10</sup> If a dimension reduction method which results in a sparse solution (e.g., sparse SPCA) is applied, then  $S_j$  is only a subset of  $G$  and might differ for different principal components.

- (3) We calculate the average marginal effect  $\hat{f}_S^{PD}(\mathbf{x}_S)$  of the feature set  $S$  exactly as in Eq. (18).
- (4) We visualize the CFEP by plotting  $(g(\mathbf{x}_S^{(i)}), \hat{f}_S^{PD}(\mathbf{x}_S^{(i)}))$  for each observation in the dataset.

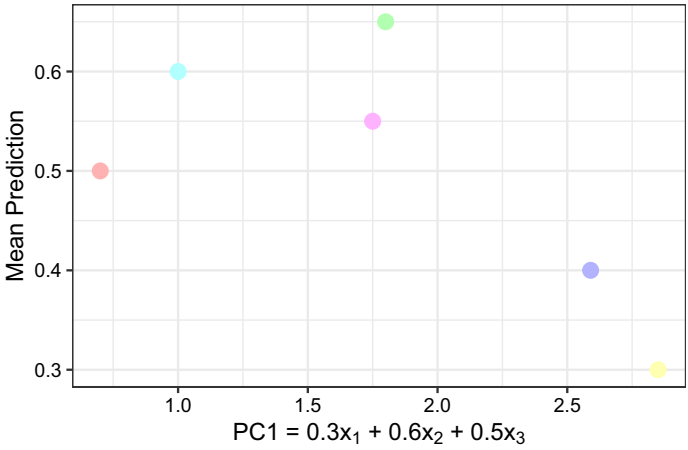
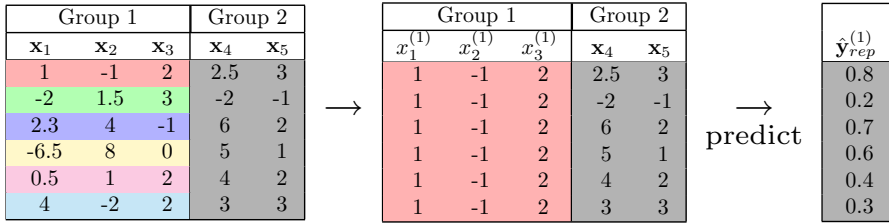
Hence, the CFEP visualizes the average marginal effect of features in  $S$  against the combinations of features received by the dimension reduction method (e.g., a linear combination of a principal component in the case of sparse SPCA) and thus shows how different values of  $g(\mathbf{x}_S)$  affect the predictions of a given model. For a feature group, several principle components  $g_j$  and hence several CFEPs may be of interest.

The CFEP is defined in Algorithm 2, but we will demonstrate the procedure of constructing a CFEP with the illustrative example in Fig. 6. In this example, we have two predefined groups of features, where the first group contains  $x_1, x_2,$  and  $x_3,$  and the second group contains features  $x_4$  and  $x_5.$  The sparse SPCA on the first group yields a first principal component ( $g_1$ ) with the loadings 0.3 for  $x_1,$  0.6 for  $x_2$  and 0.5 for  $x_3$  (step 1 to 3 of Algorithm 2). It follows that  $S = \{1, 2, 3\}$  and that the low dimensional representation of interest is  $g_1.$  For the construction of a CFEP for  $g_1,$  mean predictions for the principal component are calculated for each observation. To calculate the mean prediction of the first observation (shown in red), we replace the values of features with non-zero loadings of  $g_1$  of each instance in the dataset by the feature values of the first observation (step 6 in Algorithm 2). A prediction vector  $\hat{\mathbf{y}}_{rep}^{(1)}$  is then calculated with the previously trained model (step 7 in Algorithm 2). The value on the y-axis for the red point in the graph below corresponds to the mean over all predictions for the first observation:  $\bar{y}_{rep}^{(1)} = (0.8 + 0.2 + 0.7 + 0.6 + 0.4 + 0.3)/6 = 0.5.$  The value on the x-axis is the linear projection of the first observation for the regarded principal component (step 8 and 9 in Algorithm 2). Hence, it is calculated by the weighted sum of feature values  $x_1^{(i)}$  to  $x_3^{(i)},$  where the weights are defined by the loadings of the respective principal component that we receive with sparse SPCA.

In contrast to PDP or totalvis effect plots, CFEP produces a point cloud instead of a curve. The CFEP is, mathematically speaking, not a function, since points on the x-axis correspond to linear projections of features within a group. A point  $z$  on the x-axis can have multiple combinations of features, which lead to  $z$  and have different mean predictions on the y-axis. However, we now have the possibility to interpret the shape of the point cloud and can draw conclusions about the behavior of the mean prediction of the model regarding a linear combination of features of interest.

## 6.4 Experiments on supervised versus unsupervised dimension reduction

As discussed in Sect. 6.1, PCA might be the most popular dimension reduction method. However, since PCA is unsupervised, it does not account for the dependency between the feature space and the target variable. To evaluate the degree to which this drawback influences CFEP, we examine two regression problems on simulated data. The first is defined by a single underlying factor depending on a sparse set of features, which can be represented by a single principal component. The linear combination of this feature set is also linearly correlated with the target variable. The second regression problem



**Fig. 6** Explanation of estimating and visualizing CFEP; the x-coordinate reflects the linear combination of features with non-zero loadings for  $g_1$ , and the y-coordinate reflects the mean predictions  $\hat{y}_{rep}^{(i)}$  for each observation  $i$ . The substitution of values for each observation is only done for features with non-zero loadings

contains two underlying factors that depend on two sparse sets of features. While the linear combination of the first feature set is also linearly correlated with the target, the second factor has a quadratic effect on  $\mathbf{Y}$ . In both cases, we compare the usage of sparse supervised and unsupervised PCA (sparse SPCA and sparse PCA) as dimension reduction methods within CFEP and compare them to the totalvis effect plot. Here, we investigate if the respective dimension reduction method does correctly identify the sparse set of features for each group. Additionally, we determine how accurately we can predict the true underlying relationship between the linear combination of these features and the target variable. Since we simulated the data, we know the number of underlying factors (principal components).

**6.4.1 One factor**

In this example, we created a data matrix  $\mathbf{X}$  with 500 instances of 50 standard normally distributed features with decreasing correlations. Therefore, all features are generated as done in Sect. 5.2. The altering proportion  $\alpha$  is set to 0.2 for the first 10 features, to 0.4 for the next 10 features, and to 1 for the last 10 features. Thus, while the first 10 features are highly correlated with each other, the last 10 features are approximately

**Algorithm 2:** Combined Features Effect Plot

---

**input** : Dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$ ,  
group  $G \subset \{1, \dots, p\}$ ,  
model  $\hat{f}$  trained on  $\mathcal{D}$ .

**output**: Combined Features Effect Plot

- 1 Perform sparse SPCA on  $\hat{\mathcal{D}} := \{(\mathbf{x}_G^{(i)}, y^{(i)})\}_{i=1}^n$ ;
- 2 Choose a principle component function of interest  $g$ ;
- 3 Let  $S \subseteq G$  be the sparse set of features of  $g$ ;
- 4 **for**  $i \in \{1, \dots, n\}$  **do**
- 5     get feature values  $\mathbf{x}_S^{(i)}$ ;
- 6     create  $\mathcal{D}_{rep}^{(i)}$  by replacing feature values from  $S$  of every observation with  $\mathbf{x}_S^{(i)}$ ;
- 7     predict vector  $\hat{\mathbf{y}}_{rep}^{(i)}$  by applying  $\hat{f}$  on  $\mathcal{D}_{rep}^{(i)}$  row-wise;
- 8     calculate the mean prediction  $\tilde{y}_{rep}^{(i)}$  of  $\hat{\mathbf{y}}_{rep}^{(i)}$ ;
- 9     save  $g(\mathbf{x}_S^{(i)})$  as x-coordinate and  $\tilde{y}_{rep}^{(i)}$  as y-coordinate of observation  $i$  for the CFEP (see Eq. (19));

---

The CFEP can be used as a descriptive method to better understand the effect of a group of features on the target variable. The dimension reduction method in step 1 is exchangeable.

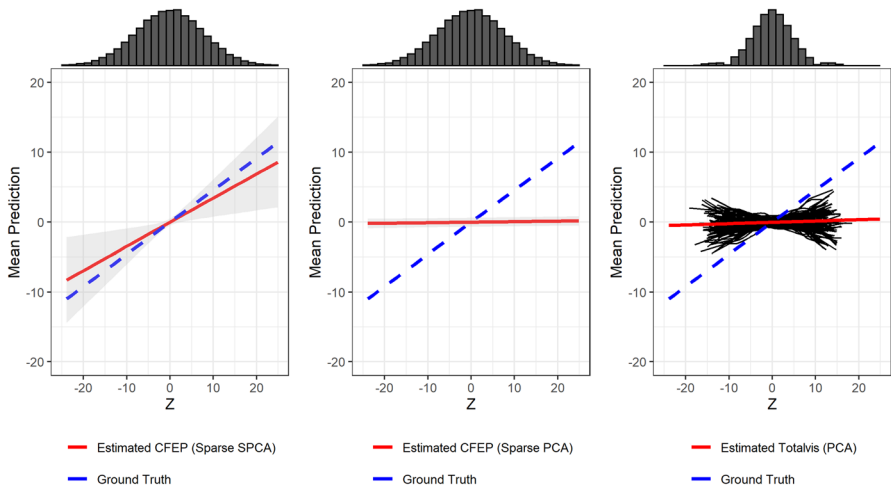
uncorrelated with each other. The sparse subgroup defined by the variable  $\mathbf{Z}$  is a linear combination of 5 features from  $\mathbf{X}$  and has itself a linear effect on the target variable  $\mathbf{Y}$ :

$$\mathbf{Z} = \mathbf{X}_5 - 2\mathbf{X}_8 - 4\mathbf{X}_{25} + 8\mathbf{X}_{47} + 4\mathbf{X}_{49}$$

$$\mathbf{Y} = \mathbf{Z} + \epsilon, \quad \text{with } \epsilon \stackrel{iid}{\sim} N(0, 1).$$

Hence, according to our notation,  $G_{\mathbf{Z}}$  is defined by  $G_{\mathbf{Z}} = \{5, 8, 25, 47, 49\}$ , and thus,  $X_{G_{\mathbf{Z}}}$  is the related subset of all features. We drew 100 samples and fitted a random forest with 2000 trees with each sample drawing. We used the 10-fold cross-validated results to perform sparse SPCA. For each dimension reduction method, we estimate  $\hat{\mathbf{Z}}$  by summing up the (sparse) loading vector (estimated by the dimension reduction method) multiplied by the feature matrix  $\mathbf{X}$ . Therefore,  $\mathbf{X}_{G_{\hat{\mathbf{Z}}}}$  is defined by the received sparse feature set. The mean prediction  $\tilde{\mathbf{Y}}_{rep}$  for the CFEP is calculated as described in Sect. 6.3.

The impact of choosing a supervised over an unsupervised sparse PCA approach is shown in Fig. 7, which also shows the average linear trend and 95% confidence bands of CFEP for the simulation results. To evaluate how well the estimated mean prediction  $\hat{\mathbf{Y}}_{rep}$  approximates the underlying trend, we assume that we know that  $\mathbf{Z}$  has a linear influence on the target. Thus, we fit a linear model on each simulation result. To compare the received regression lines, we evaluate each of them on a predefined grid and average over all 100 samples (represented by the red line). The confidence bands are then calculated by taking the standard deviation over all estimated regression lines on grid level and calculating the 2.5% and 97.5% quantiles using the standard normal approximation. The associated calculation steps for each of the 100 samples can be summarized as follows:



**Fig. 7** Average linear trend and confidence bands of CFEP over all samples using sparse SPCA (left) and sparse PCA (middle) compared to estimated totalvis effect curves over all 100 samples for the first principal component (black) and the average linear trend (red) (right) (Color figure online)

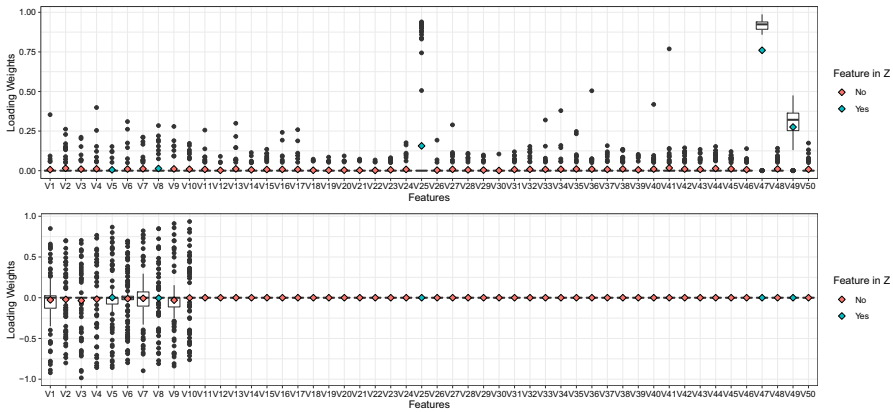
- (1) Estimate a linear model  $\hat{f}(\mathbf{X}_{G_{\mathbf{Z}}}) \sim \mathbf{Z}$ .
- (2) Define an equidistant grid of length 50 within the range of  $\mathbf{Z}$ .
- (3) Apply the linear model estimated in 1) on the grid defined in 2).
- (4) Repeat steps 1 to 3 for  $\hat{f}(\mathbf{X}_{G_{\mathbf{Z}}})$  by using the true underlying features of  $\mathbf{Z}$  to calculate the combined features dependencies that we call the ground truth.

The left plot in Fig. 7 shows a similar linear trend of the estimated CFEP compared to the average ground truth (represented by the blue line), while the red line in the right plot varies around 0. By using sparse SPCA, the underlying feature set  $\mathbf{X}_{G_{\mathbf{Z}}}$  is better approximated than with sparse PCA, which is reflected in the MSE between  $\mathbf{Z}$  and  $\hat{\mathbf{Z}}$  of 0.7 for sparse SPCA and 1.9 for sparse PCA. Figure 8 provides an explanation for those differences. While sparse SPCA (on average) more strongly weights features that have a large influence on the target, impactful loading weights for sparse PCA are solely distributed over highly correlated features in  $\mathbf{X}$  that explain the most variance in the feature space. Thus, including the relationship between the target and  $\mathbf{X}$  in the dimension reduction method may have a huge influence on correctly approximating the underlying factor and, hence, also on the CFEP.

Similar to using sparse PCA as a dimension reduction method within CFEP, on average, the totalvis effect curves based on PCA do not show a clear positive linear trend (see Fig. 7). For almost half of the samples, we even receive a negative instead of a positive trend for the underlying factor. The interpretation is opposite to the actual effect and, hence, is misleading.

### 6.4.2 Two factors

In real-world data settings are often more complex by containing non-linear relationships and the target variable is described by more than one underlying factor. Hence,



**Fig. 8** Distribution of feature loadings in sparse SPCA (top) and sparse PCA (bottom) over all samples. Rhombuses denote the mean values, with the blue rhombuses indicating the features that have an influence on the target in the underlying model formula (Color figure online)

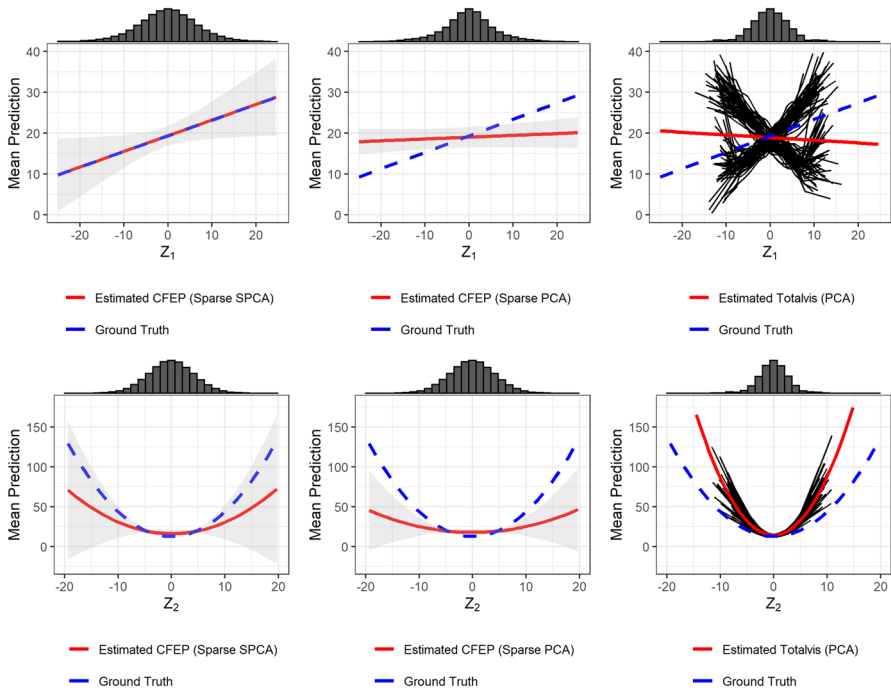
we now examine a more complex simulation setting to assess if we can observe the same behavior that we observed for the simple case. To that end, we simulated a data matrix  $\mathbf{X}$  with 500 instances for two feature sets, each containing 20 standard normally distributed features. The data for each feature set is generated as described in Sect. 5.2 but with an altering proportion of 0.15 and 0.35 for the features in the first set and 0.55 and 0.85 in the second set. Hence, within each set, the first ten features show a higher correlation among each other than the last ten features. Additionally, all features of the first set are on average more highly correlated than all features of the second set. Features between the two sets are uncorrelated. The first factor  $\mathbf{Z}_1$  is a linear combination of four features from the first set and  $\mathbf{Z}_2$  of two features from the second set.  $\mathbf{Z}_1$  has a linear and  $\mathbf{Z}_2$  a quadratic effect on  $\mathbf{Y}$ .

$$\begin{aligned} \mathbf{Z}_1 &= 3\mathbf{X}_3 - 2\mathbf{X}_8 - 4\mathbf{X}_{13} + 8\mathbf{X}_{18} \\ \mathbf{Z}_2 &= 2\mathbf{X}_{25} + 4\mathbf{X}_{35} \\ \mathbf{Y} &= \mathbf{Z}_1 + \mathbf{Z}_2^2 + \epsilon, \quad \text{with } \epsilon \stackrel{iid}{\sim} N(0, 1). \end{aligned}$$

Again, we drew 100 samples and fitted a random forest with 2000 trees with each sample drawing. The approach is almost the same as described for one factor, with the difference being that we use the first two principal components (as we want to find two sparse feature sets instead of one).

In Fig. 9, the average linear and quadratic trend of the underlying CFEPs of  $\mathbf{Z}_1$  and  $\mathbf{Z}_2$  are depicted for both dimension reduction methods. While the average linear regression line of sparse SPCA matches the average ground truth almost perfectly for  $\mathbf{Z}_1$ , the associated line of sparse PCA shows only a slightly positive trend and differs substantially from the ground truth. Regarding  $\mathbf{Z}_2$ , a similar propensity can be observed for the quadratic shape. Again, this behavior results from sparse SPCA (on





**Fig. 9** Top ( $Z_1$ ): Average linear trend and confidence bands of CFEP over all samples using sparse SPCA (left) and sparse PCA (middle) compared to estimated totalvis effect curves over all 100 samples for first principal component (black) and the average linear trend (red) (right). Bottom ( $Z_2$ ): Same structure as for  $Z_1$ , but showing the quadratic trend of  $Z_2$  (Color figure online)

average) more strongly weighting features that have a large effect on the target, while the unsupervised version focuses on features that explain the most variance in  $\mathbf{X}$ .

The estimated linear trend of the totalvis effect curves for the first principal component is negative instead of positive. Thus, for most of the samples and on average, these results are completely misleading (see Fig. 9). The quadratic shape of the second component is (on average and for almost all samples) steeper than the average ground truth. Additionally, the deviation is higher here than for CFEP with sparse SPCA.

## 7 Real data example: smartphone sensor data

Smartphones and other consumer electronics have increasingly been used to collect data for research (Miller 2012; Raento et al. 2009). The emerging popularity of these devices for data collection is grounded in their connectivity, the number of built-in sensors, and their widespread use. Moreover, smartphones enable users to perform a wide variety of activities (e.g., communication, shopping, dating, banking, navigation, listening to music) and thus provide an ideal means to study human behavior in naturalistic contexts, over extended periods of time, and at fine granularity (Harari et al. 2015, 2016, 2017). In this regard, smartphone data has been used to investigate

individual differences in personality traits (Stachl et al. 2017; Harari et al. 2019), in human emotion and well-being (Servia-Rodríguez et al. 2017; Rachuri et al. 2010; Saeb et al. 2016; Thomée 2018; Onnela and Rauch 2016; Kolenik and Gams 2021), and in daytime and nighttime activity patterns (Schoedel et al. 2020).

We use a dataset on human behavior, collected with smartphones, to illustrate methods for group-based feature importance. The PhoneStudy dataset was consolidated from three separate datasets (Stachl et al. 2017; Schuwerk et al. 2019; Schoedel et al. 2018). It consists of 1821 features on smartphone-sensed behavior and 35 target variables on self-reported Big Five personality trait dimensions (domains) and subdimensions (facets). The dataset has been published online and is openly available.<sup>11</sup> The Big Five personality trait taxonomy is the most widely used conceptualization of stable individual differences in human patterns of thoughts, feelings, and behavior (Goldberg 1990). In their original study, Stachl et al. (2020a) used the behavioral variables to predict self-reported Big Five personality trait scores (five dimensions and 30 subdimensions) and used grouped feature importance measures to explore which classes of behaviors were most predictive for each personality trait dimension. The groups in this study were created based on theoretical considerations from past work.

The personality prediction task is challenging because (1) the dataset contains many variables on similar behaviors, (2) these variables are often correlated, and (3) effects with the targets are interactive, very small, and partially non-linear. Many variables in the dataset can be manually grouped into classes of behavior (e.g., communication and social activity, app-usage, music consumption, overall phone activity, mobility).

We use this dataset to illustrate the idea of grouped feature importance with regard to the prediction of personality trait scores for the dimension of conscientiousness (Table 5). Conscientiousness is a personality trait dimension that globally describes people's propensity to be reliable, dutiful, orderly, ambitious, and cautious (Jackson et al. 2010). We chose this personality trait because it has high practical relevance due to its ability to predict important life outcomes and behaviors (Ozer and Benet-Martínez 2006). Here, we (1) fit a random forest model to predict the personality dimension of conscientiousness, (2) compute the introduced methods for grouped feature importance (GOPFI, GPFI, GSI, LOGI, LOGO), (3) use the proposed sequential grouped feature importance procedure to investigate which groups are most important in combination, and (4) visualize the effect of different groups with CFEPs. Once the importance of individual groups has been quantified, CFEPs can be helpful to further explore the variables in these groups with regard to the criterion variable of interest (i.e., conscientiousness) to generate new hypotheses for future research.

In Fig. 10, we show a sequential procedure for our personality prediction example. The figure shows that the groups *overall phone usage* and *app usage* lead to the best model performance if used alone and, in many cases, lead to even better performances if combined. The results also suggests that if only one group can be selected, the initial selection of the feature group app usage more often leads to the smallest expected loss (mean MSE = 0.519). For a practical application, this would indicate that if only one type of feature may be collected from smartphones to predict the personality trait conscientiousness, features on app usage should be used. If two groups of data can

<sup>11</sup> <https://osf.io/kqjhr/>.

**Table 5** Grouped feature importance values for predicting the personality trait conscientiousness based on MSE

Group	GOPFI	GPMFI	GSI	LOGI	LOGO
Mobility (Mo)	-0.002 (± 0.011)	-0.002 (± 0.001)	0.000 (± 0.003)	-0.011 (± 0.075)	0.000 (± 0.006)
Music (Mu)	-0.001 (± 0.011)	0.002 (± 0.002)	0.001 (± 0.006)	-0.019 (± 0.074)	0.001 (± 0.012)
Communication and social (C)	0.000 (± 0.008)	0.001 (± 0.003)	0.004 (± 0.006)	0.008 (± 0.070)	0.001 (± 0.010)
Overall phone usage (O)	0.007 (± 0.011)	0.009 (± 0.003)	0.012 (± 0.008)	0.032 (± 0.080)	0.009 (± 0.014)
App usage (A)	0.032 (± 0.009)	0.028 (± 0.005)	0.031 (± 0.012)	0.041 (± 0.069)	0.011 (± 0.019)

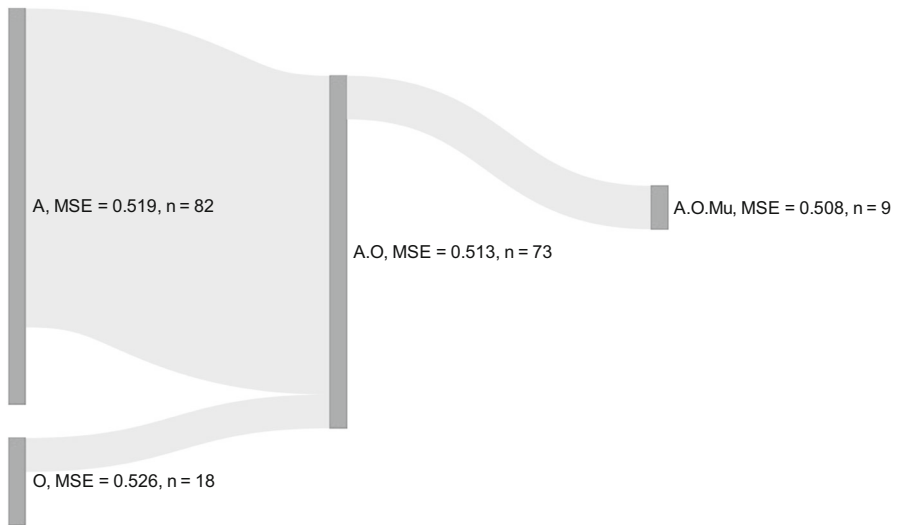
All values were calculated using a resampling method (10-times cross-validation)

be collected, overall phone usage should also be added (mean MSE = 0.513). Finally, the plot indicates that in some cases ( $n = 9$ ), the additional consideration of music listening behaviors in the model could lead to additional, small improvements of the expected loss (mean MSE= 0.508). If a feature group is not added, this means that it did not make a significant contribution in this iteration of the data split. Interestingly, the feature group *music* alone shows very low (or even negative) grouped feature importance scores. This would mean that music features are only predictive in the presence of other features.

To additionally explore meaningful and predictive directions in the feature space of the app usage group, we use CFEPs for the visualization. Subplot (a) in Fig. 11 shows that combinations of higher values in features on weather app usage on average lead to higher mean values in the personality trait conscientiousness. The increased frequency in weather app usage could signify the propensity of conscientious people to be prepared for future eventualities (e.g., bad weather; Jackson et al. 2010). Subplot (b) shows an interesting non-monotonic relationship between the number of different apps used each day and the mean value in conscientiousness. Subplot (c) shows that the combinations of higher values in overall phone activities lead to lower mean values in conscientiousness. Finally, plot (d) shows a similar, negative effect pattern with regard to music listening behaviors.

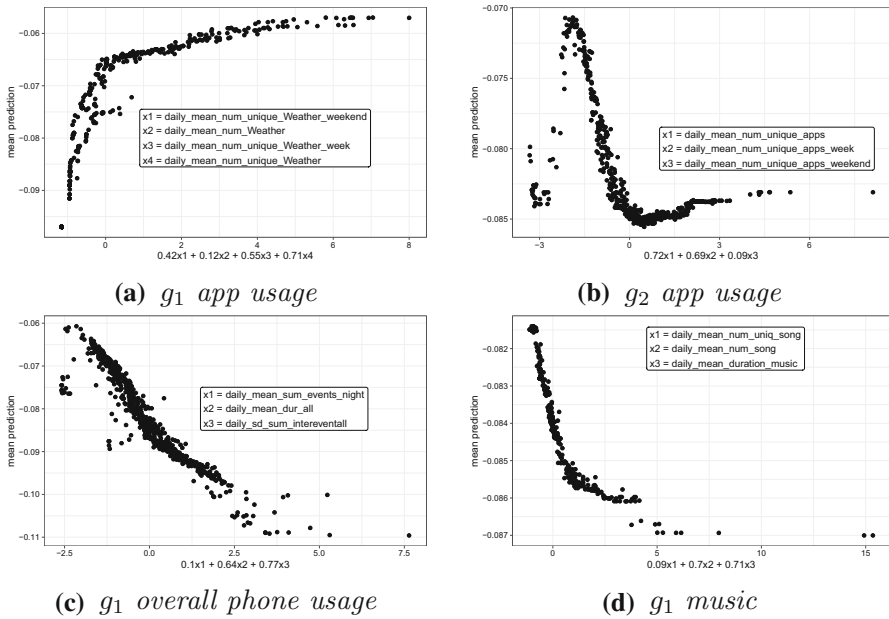
## 8 Conclusion

We introduced various techniques to analyze the importance and effect of user-defined feature groups on predictions of ML models. We provided formal definitions and dis-



**Fig. 10** Sequential grouped feature importance procedure for smartphone sensor data predicting *conscientiousness*. 100 times repeated subsampling. Inner resampling strategy: 10-fold cross-validation. Improvement threshold  $\delta = 0.01$ . Abbreviations: app-usage (A), communication & social (C), music (Mu), overall phone activity (O), mobility (Mo). Vertical bars show one step in the greedy forward search algorithm. Height of the vertical bars represent the number of subsampling iterations in which a combination of groups was chosen (for example, out of 100 subsampling iterations the group app-usage (A) was chosen 82 times as the best first group). Streams indicate the proportion of iterations that additionally benefited from a consequent step. Only streams containing at least 5 iterations and better mean performance at the end are displayed

inction criteria for grouped feature importance methods and distinguished between permutation- and refitting-based methods. For both approaches, we defined two calculation strategies that either start with a null model or with the full model. Based on these two definitions, we introduced Shapley importance scores for groups, which we defined for permutation as well as refitting methods. Moreover, we introduced as our first main contribution a sequential grouped feature importance procedure to find good and stable combinations of feature groups. To contrast the newly proposed methods with existing ones, we compared them for different scenarios. The key recommendations for the user can be summarized for four scenarios: (1) If high correlations between groups are present, refitting methods should be preferred over permutation methods, since they often deliver more meaningful results in these scenarios. Moreover, if the number of groups is reasonably small, refitting methods become computationally feasible. (2) If a sparse set of feature groups is of interest (e.g., due to data availability), the introduced sequential procedure can be useful. It provides insights regarding the most important groups: which sparse group combinations are stable in the sense that they are frequently selected and achieve a good performance. These criteria can be critically informative in situations where feature groups must be obtained from different data sources that are associated with further costs. (3) If the correlation strengths of features within groups are very diverse, all of the introduced methods might fail to reflect the true underlying importance of the feature groups. The size of this effect



**Fig. 11** CFEPs for the prediction of the personality trait conscientiousness.  $g_1$  describes the first principal component of the respective group, and  $g_2$  describes the second. More details about the features can be found in Appendix D and on the supplemental website [https://compstat-lmu.shinyapps.io/Personality\\_Prediction/](https://compstat-lmu.shinyapps.io/Personality_Prediction/) for Stachl et al. (2020a)

depends heavily on how well the fitted model captures the true underlying relationship between features. Especially when using random forests, we showed that all of the methods lead to misleading results. (4) Groups with many features might tend to have a higher grouped importance score than groups with fewer features. Normalizing the grouped importance score leads to an average score per feature. However, this might result in choosing groups where grouped scores are smaller than those of other groups and, hence, contain less (performance-based) information than others. When using GSI, users can extract additional feature-level information to gain more insights into the group scores. Specifically, we showed that single feature Shapley importance scores add up to GSI when no higher-order interactions between groups are present. As third main contribution we proposed the CFEP, which is another global interpretation method that allows visualizations of the combined effect of multiple features on the prediction of an ML model. By applying a sparse SPCA, we received more meaningful and interpretable results for the final CFEPs compared to its unsupervised counterpart. We also demonstrated the suitability of the method in our real data example from computational psychology. Although, we only considered a numeric feature space here, all methods are in general also applicable to mixed feature spaces. However, in the presence of categorical features, a suitable dimension reduction method for CFEP must be chosen.

Here, we have focused on knowledge-driven feature groupings. However, the introduced methods could also be applied to data-driven groups (e.g., via shared variance).

Notably, their interpretation is only meaningful if groups can be described by some underlying factor. This might be a good application for interpretable latent variables to find causal relationships between feature groups and predictions of ML models. Additionally, with regard to highly correlated feature groups that cannot be grouped naturally, a data-driven approach might be more suitable.

It is our goal that this article not only provides a helpful reference for researchers in selecting appropriate interpretation methods when features can be grouped, but also that it inspires future research in this area.

**Author Contributions** Conceptualization: QA, JH, GC, BB; Methodology: QA, JH, GC; Formal analysis and investigation: QA, JH, GC; Writing - original draft preparation: QA, JH; Writing - review and editing: GC, CS, BB; Investigation: QA, JH; Visualization: QA, JH; Validation: QA, JH, GC; Software: QA, JH; Funding acquisition: GC, CS, BB; Supervision: GC, BB.

**Funding** Open Access funding enabled and organized by Projekt DEAL. This work has been partially supported by the German Federal Ministry of Education and Research (BMBF) under Grant No. 01IS18036A, the Bavarian State Ministry of Science and the Arts in the framework of the Centre Digitisation.Bavaria (ZD.B), the Bavarian Ministry of Economic Affairs, Regional Development and Energy as part of the program “Bayerischen Verbundförderprogramms (BayVFP)—Förderlinie Digitalisierung—Förderbereich Informations- und Kommunikationstechnik” under the Grant DIK-2106-0007 // DIK0260/02, a Google research grant, the LMU-excellence initiative, and the National Science Foundation (NSF) Award SES-1758835. The authors of this work take full responsibility for its content.

**Availability of data and materials** All data are created or provided in the following public git-repository: [https://github.com/slds-lmu/grouped\\_feat\\_imp\\_and\\_effects](https://github.com/slds-lmu/grouped_feat_imp_and_effects).

**Code Availability** The implementation of the proposed methods and reproducible scripts for the experimental analysis are provided in the following public git-repository: [https://github.com/slds-lmu/grouped\\_feat\\_imp\\_and\\_effects](https://github.com/slds-lmu/grouped_feat_imp_and_effects).

## Declarations

**Conflict of interest** Not applicable.

**Consent to participate** Not applicable.

**Ethics approval** Not applicable.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## Appendix A Motivational example for grouped importance methods

In some settings, permuting single features individually might not be meaningful, for example, when categorical features are dummy-encoded. Table 6 shows for two

**Table 6** We draw 1000 samples of two independent categorical random variables  $X_1, X_2 \in \{1, 2, 3, 4\}$  where the categories 1 and 2 occur four times more frequently than 3 and 4

Method	$X_1$	$X_{2,2}, X_{2,3}, X_{2,4}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
Individually permuted	2.63	–	2.45	1.00	1.71
Group-wise permuted	2.63	2.65		–	

Consider the target  $y = 5 \cdot \mathbb{1}_{X_1 \neq 1} + 5 \cdot \mathbb{1}_{X_2 \neq 1} + \epsilon$  with  $\epsilon \sim \mathcal{N}(0, 1)$ . Both categorical features have the same influence on the target. We explicitly dummy encode  $X_2$  using  $X_2 = 1$  as the reference category to obtain 3 binary features  $X_{2,k} = \mathbb{1}_{X_2=k}, k \in \{2, 3, 4\}$ . We fit a linear model using the categorical feature  $X_1$  and the binary features  $X_{2,2}, X_{2,3}, X_{2,4}$ . Here, we want to illustrate why it makes more sense to permute the 3 binary features jointly rather than individually, since they naturally belong together. As expected, permuting the binary features  $X_{2,2}, X_{2,3}, X_{2,4}$  jointly as a group yields a comparable importance to  $X_1$ . However, permuting each binary feature individually gives different importance scores making it unclear how important  $X_2$  is compared to  $X_1$

equally important categorical features that if one feature is dummy-encoded (here:  $X_2$ ), then all resulting binary features must be permuted as a group to obtain a comparable importance score to  $X_1$ . Hence, settings like in Table 6 or as described in Sects. 1 or 1.2 point out the need of grouped importance methods.

## Appendix B Shapley importance

### B.1 Properties of the grouped Shapley importance

For single features<sup>12</sup>  $x_i \in \{1, \dots, p\}$ , which are divided into  $l$  groups, we define the marginal contribution for  $x_i$  as

$$\Delta_{\{x_i\}}(S) = v(S \cup \{x_i\}) - v(S),$$

for  $S \subset \{1, \dots, p\} \setminus \{x_i\}$ . The Shapley importance for single features  $\phi(x_i)$  can also be defined analogously to (15). One interesting question is, does the GSI for a group  $G \subset \{1, \dots, p\}$  decompose into the sum of Shapley importances of features in  $G$ ? In the following, we want to analyze the remainder

$$R = \phi(G) - \sum_{i \in G} \phi(x_i). \tag{B1}$$

Similar to the functional ANOVA decomposition (Hooker 2004), we assume, that the value function for a coalition  $S \subset \{1, \dots, p\}$  can be broken down into main and interaction effects

$$v(S) = v_0 + \sum_{x_i \in S} v(x_i) + \sum_{i \neq j} \epsilon_{ij} + \sum_{i \neq j \neq k} \epsilon_{ijk} + \dots, \tag{B2}$$

<sup>12</sup> Remember the one-to-one association of the numbers  $1, \dots, p$  and the features  $\mathbf{x}_1, \dots, \mathbf{x}_p$

where  $\epsilon_{i\dots m}$  is the effect of the interaction between the features  $x_i, \dots, x_m \in S$ . A needed requirement to apply this decomposition is that each of the functional terms has zero means, hence they need to be centralized. The considered intercept shift is stored in  $v_0$ . To receive a unique decomposition, the orthogonality between the functional terms needs to be fulfilled which is not the case in the presence of correlated features. Hooker (2007) therefore suggests the generalized functional ANOVA which replaces the orthogonality property with a hierarchical orthogonality condition and which is a weighted version of the standard functional ANOVA (Hooker 2004). However, we do not try to estimate or calculate the decomposed function terms, we only use the (valid) assumption that a function can be decomposed as in Eq. (B2) to show how GSI relates to Shapley importance for individual features. Hence, we are not directly interested in a unique solution of the decomposition.

With the assumption in Eq. (B2), it follows that the Shapley importance of a single feature  $x_1$  (without loss of generality) can be written as

$$\phi(x_1) = v(x_1) + \frac{1}{2} \left( \sum_{i \neq 1}^p \epsilon_{1i} \right) + \frac{1}{3} \left( \sum_{i \neq j \neq 1}^p \epsilon_{1ij} \right) + \dots + \frac{1}{p} \epsilon_{1\dots p}. \quad (\text{B3})$$

The value function of the feature  $x_1$  contributes to the Shapley importance with the weight 1 and all possible interaction effects with feature  $x_1$  contribute with the reciprocal length of the interaction effect. We proved this assertion in Appendix B.2. Similar to (B3), the GSI of a group  $G_1$  (w.l.o.g.) can be written as

$$\phi(G_1) = v(G_1) + \frac{1}{2} \left( \sum_{i \neq 1}^k \epsilon_{G_1 G_i} \right) + \frac{1}{3} \left( \sum_{i \neq j \neq 1}^k \epsilon_{G_1 G_i G_j} \right) + \dots + \frac{1}{k} \epsilon_{G_1 \dots G_k} \quad (\text{B4})$$

where  $\epsilon_{G_1 \dots G_k}$  is the (non-computable) interaction effect between features of groups  $G_1, \dots, G_k$ , where each group provides at least one feature. By using Eq. (B2) on  $v(G_1)$ , we get:

$$v(G_1) = \sum_{i \in G_1} v(x_i) + \sum_{i \neq j \in G_1} \epsilon_{ij} + \sum_{i \neq j \neq k \in G_1} \epsilon_{ijk} + \dots \quad (\text{B5})$$

Looking back at Eq. (B1), a lot of terms cancel out by using Eqs. (B3) and (B5). The term  $v(G_1)$ , meaning all main effects  $v(x_i), i \in G_1$ , and all interaction effects  $\epsilon_{i, \dots, k}, 1 \leq k \leq |G_1|$  between features within  $G_1$ , cancels out entirely.<sup>13</sup> Furthermore, at least all two-way interaction effects between groups  $\epsilon_{G_1 G_i}, i = 2, \dots, k$  cancel out. A combination of higher-order interaction terms between features of  $G_1$  and  $\{1, \dots, p\} \setminus G_1$  remain.<sup>14</sup> This means that the remainder  $R$  is (usually) not equal to zero in case the applied algorithm learned a higher-order interaction between features

<sup>13</sup> Note,  $v(G_1)$  cancels out, meaning that these interaction terms cannot be computed directly but are assumed to affect the ‘‘payout’’ of the value function.

<sup>14</sup> They mostly only partly cancel out, depending on the number of features within the groups  $G_1, \dots, G_k$ .



of the regarded group and other groups. The higher the remainder, the larger the higher-order interaction effect. Thus, the remainder can be used as a quantification of learned higher-order interaction effects between features of different groups.

### B.2 Proof of Properties

Assume, that the value function for a coalition  $S \subset \{x_1, \dots, x_p\}$  can be broken down into main and interaction effects:

$$v(S) = \sum_{x_i \in S} v(x_i) + \sum_{i_1 \neq i_2} \epsilon_{i_1 i_2} + \sum_{i_1 \neq i_2 \neq i_3} \epsilon_{i_1 i_2 i_3} + \dots,$$

the Shapley importance of a single feature  $x_1$  can be written as

$$\phi(x_1) = v(x_1) + \frac{1}{2} \left( \sum_{i \neq 1}^p \epsilon_{1i} \right) + \frac{1}{3} \left( \sum_{i \neq j \neq 1}^p \epsilon_{1ij} \right) + \dots + \frac{1}{p} \epsilon_{1\dots p}.$$

**Proof** Let  $N = \{x_2, \dots, x_p\}$ . The general formula for the Shapley importance is given by:

$$\phi_p(x_1) = \sum_{S \subset N \setminus \{x_1\}} \frac{(p-1-|S|)! \cdot |S|!}{p!} (v(S \cup \{x_1\}) - v(S)) \tag{B6}$$

With assumption (B2) the term  $v(S \cup \{x_1\}) - v(S)$  will reduce to:

$$v(S \cup \{x_1\}) - v(S) = v(x_1) + \sum_{i_1 \neq 1}^p \epsilon_{1i_1} + \dots + \sum_{i_1 \neq \dots \neq i_{|S|} \neq 1}^p \epsilon_{1i_1 \dots i_{|S|}} \tag{B7}$$

It is the sum of  $v(x_1)$  and all interactions with feature  $x_1$  of sizes  $2, \dots, |S| + 1$ . All other terms without feature  $x_1$  cancel out.

Equation (B6) consists of many summands of the form (B7). The term  $v(x_1)$  appears once for every subset  $S \subset N \setminus \{x_1\}$ . There are  $\binom{p-1}{|S|}$  different subsets of size  $|S|$ . Only looking at the summands with the term  $v(x_1)$ , Eq. (B6) reduces to

$$\sum_{|S|=0}^{p-1} \frac{(p-1-|S|)! \cdot |S|!}{p!} \binom{p-1}{|S|} v(x_1) = v(x_1). \tag{B8}$$

For the interaction terms, we first start counting the interaction term  $\epsilon_{12}$  of size 2, as an example. For  $|S| = 0$ , there are zero terms of  $\epsilon_{12}$ . For  $|S| = 1$ , the term  $\epsilon_{12}$  only appears once, when  $S = \{x_2\}$ . For  $|S| = 2$ , the term  $\epsilon_{12}$  appears  $p - 2$  times, once for each subset  $S = \{x_2, x_j\}$ , for  $3 \leq j \leq p$ . For  $|S| = 3$ , we have  $\binom{p-2}{2}$  times the term  $\epsilon_{12}$ , again, once for each subset  $S = \{x_2, x_j, x_k\}$ , for  $3 \leq j \neq k \leq p$ . This pattern goes on until there are  $\binom{p-2}{p-2}$  terms of  $\epsilon_{12}$  for  $|S| = p - 1$ . Now, we look at

the interaction terms  $\epsilon_{1i_1 \dots i_{k-1}}$  of size  $k$ . Following the pattern, which we just derived, there are zero terms of  $\epsilon_{1i_1 \dots i_{k-1}}$  for  $|S| \leq k - 2$  and  $\binom{p-k}{|S|-k+1}$  terms of  $\epsilon_{1i_1 \dots i_{k-1}}$  for  $k \leq |S| \leq p - 1$ . If we only look at the interaction terms  $\epsilon_{1i_1 \dots i_{k-1}}$  of size  $k$  and following the Eq. (B6), we get

$$\sum_{|S|=k-1}^{p-1} \frac{(p-1-|S|)! \cdot |S|!}{p!} \binom{p-k}{|S|-k+1} \epsilon_{1i_1 \dots i_{k-1}} = \frac{1}{k} \epsilon_{1i_1 \dots i_{k-1}},$$

which was left to show the assertion. □

## Appendix C More details on dimension reduction techniques

### C.1 Principal component analysis

PCA only considers the data matrix  $\mathbf{X}$  and does not take the target vector  $\mathbf{Y}$  into account. This procedure is thus unsupervised.

Given a centering Matrix

$$\mathbf{H} = \mathbf{I} - n^{-1} \mathbf{e} \mathbf{e}^T, \tag{C9}$$

where  $\mathbf{e}$  is an  $n$ -dimensional vector of all ones. The centered matrix is  $\mathbf{X}_C = \mathbf{H}\mathbf{X}$ . The sample covariance matrix of  $\mathbf{X}$  can be written as:

$$\mathbf{S}_X := \frac{1}{n} \mathbf{X}_C^T \mathbf{X}_C = \frac{1}{n} \mathbf{X}^T \mathbf{H} \mathbf{H} \mathbf{X} \tag{C10}$$

The goal is to maximize the total variance of projected data, which is equivalent to maximizing trace of the sample covariance matrix. Equation (C10) can also be written as  $\mathbf{S}_X = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_C^{(i)} \mathbf{x}_C^{(i)T}$ , where  $\mathbf{x}_C^{(i)}$  corresponds to the  $i$ -th row of  $\mathbf{X}_C$ . By projecting each data point by some unknown vectors  $\mathbf{v}_j, j = 1, \dots, p$ , we get the projected variance for each  $j = 1, \dots, p$ , which is:

$$\frac{1}{n} \sum_{i=1}^n \mathbf{v}_j^T \mathbf{x}_C^{(i)} \mathbf{x}_C^{(i)T} \mathbf{v}_j = \mathbf{v}_j^T \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_C^{(i)} \mathbf{x}_C^{(i)T} \right) \mathbf{v}_j = \mathbf{v}_j^T \mathbf{S}_X \mathbf{v}_j.$$

Let  $\mathbf{V} \in \mathbb{R}^{p \times p}$  be the full projection matrix. The projected total variance is  $tr(\mathbf{V}^T \mathbf{S}_X \mathbf{V})$ , and by ignoring constant terms, PCA finds a solution to the problem

$$\underset{\mathbf{V}}{\operatorname{argmax}} \operatorname{tr}(\mathbf{V}^T \mathbf{S}_X \mathbf{V}) = \underset{\mathbf{V}}{\operatorname{argmax}} \operatorname{tr}(\mathbf{V}^T \mathbf{X}^T \mathbf{H} \mathbf{H} \mathbf{X} \mathbf{V}) \tag{C11}$$

with an Eigen decomposition of the covariance matrix  $\mathbf{S}_X$ . The resulting Eigen vectors thus maximize the variation of projected data.

## C.2 Measuring statistical dependence with Hilbert Schmidt norms

In Gretton et al. (2005) a more generalized measure of dependence between variables  $X$  and  $Y$  was introduced:

Two random variables  $X$  and  $Y$  are independent if and only if any bounded continuous function of them are uncorrelated.

In more detail, this means that any pairs  $(X, Y), (X, Y^2), (X^2, Y), (\cos(X), \log(Y)), \dots$  have to be uncorrelated. The resulting independence measure is called the Hilbert-Schmidt Independence Criterion (HSIC). For the analysis of this independence measure, it is necessary to analyze functions on random variables. Therefore theory of Hilbert spaces and concepts of functional analysis are necessary for a thorough analysis, but they are not part of this paper. For an extensive discussion of Hilbert spaces, especially reproducing kernel hilbert spaces (RKHS) we refer to Hein and Bousquet (2004).

Let  $\mathcal{F}$  be a separable RKHS containing all bounded continuous functions from  $\mathcal{X}$  to  $\mathbb{R}$ . The associated kernel shall be denoted by  $\mathbf{K} \in \mathbb{R}^{n \times n}$ , with  $\mathbf{K}_{ij} = k(x_i, x_j)$ . Concurrently, let  $\mathcal{G}$  be a separable RKHS with bounded continuous functions from  $\mathcal{Y}$  to  $\mathbb{R}$  and associated kernel  $\mathbf{L} \in \mathbb{R}^{n \times n}$ , with  $\mathbf{L}_{ij} = l(y_i, y_j)$ .

We are particularly interested in the cross variance between  $f$  and  $g$ :

$$Cov(f(x), g(y)) = \mathbb{E}_{x,y}[f(x)g(y)] - \mathbb{E}_x[f(x)]\mathbb{E}_y[g(y)] \tag{C12}$$

A function, which maps one element from one hilbert space to another hilbert space is called *operator*. A theorem (see e.g. Fukumizu et al. 2004) states, that there exists a unique operator  $C_{X,Y} : \mathcal{G} \rightarrow \mathcal{F}$  with

$$\langle f, C_{x,y}(g) \rangle_{\mathcal{F}} = Cov(f(x), g(y)). \tag{C13}$$

The Hilbert-Schmidt Independence Criterion (HSIC) is defined as the squared Hilbert-Schmidt norm of the cross-covariance operator  $C$ :

$$HSIC(P_{\mathcal{X},\mathcal{Y}}, \mathcal{F}, \mathcal{G}) = \|C_{x,y}\|_{HS}^2 \tag{C14}$$

$\|C_{x,y}\|_{HS}^2 = 0$  if and only if the random variables  $\mathcal{X}$  and  $\mathcal{Y}$  are independent. For a detailed discussion and derivation of the HSIC independence measure, we refer to Gretton et al. (2005). The HSIC measure was used for feature selection in Song et al. (2007) or for supervised principal components in Barshan et al. (2011).

### C.2.1 Empirical HSIC

For a dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$  the empirical HSIC is:

$$HSIC(\mathcal{D}, F, G) = (n - 1)^{-2}tr(\mathbf{KHLH}) = (n - 1)^{-2}tr(\mathbf{HKHL}), \tag{C15}$$

where  $\mathbf{H}$  is the centering matrix from (C9). A high level of dependency between two kernels yields a high HSIC value.

### C.3 Supervised sparse principal components

In the process of finding interpretable latent variables, which also incorporate dependencies to a target variable, the Sparse Supervised Principal Components (SPCA), which was introduced in Sharifzadeh et al. (2017), is a suitable method for our application.

For sparse SPCA the kernel matrix  $K$  is defined as  $K = X V V^T X^T$  with a constraint for unit length and an  $L_1$  penalty for sparsity. By ignoring constant terms, we get the optimization problem:

$$\operatorname{argmax}_{\mathbf{V}} \operatorname{tr}(\mathbf{H}\mathbf{K}\mathbf{H}\mathbf{L}) = \operatorname{argmax}_{\mathbf{V}} \operatorname{tr}(\mathbf{H}\mathbf{X}\mathbf{V}\mathbf{V}^T\mathbf{X}^T\mathbf{H}\mathbf{L}) \quad (\text{C16})$$

$$= \operatorname{argmax}_{\mathbf{V}} \operatorname{tr}(\mathbf{V}^T\mathbf{X}^T\mathbf{H}\mathbf{L}\mathbf{H}\mathbf{X}\mathbf{V}) \quad (\text{C17})$$

$$s.t. \quad \mathbf{V}^T\mathbf{V} = \mathbf{I}, \quad |\mathbf{V}| \leq c. \quad (\text{C18})$$

Note, that without the sparsity constraint, (C17) reduces to (C11), when choosing  $\mathbf{L} = \mathbf{I}$ . Already explained in Barshan et al. (2011), PCA is a special form of their Supervised PCA, where setting  $\mathbf{L} = \mathbf{I}$  is a kernel, which only captures similarity between a point and itself. Maximizing dependency between  $\mathbf{K}$  and the identity matrix corresponds to retaining maximal diversity between observations.

Now, an arbitrary  $\mathbf{L}$  can be decomposed as  $\mathbf{L} = \Delta\Delta^T$ , since  $\mathbf{L}$ , as a kernel matrix, is positive definite and symmetric. Defining  $\Psi := \Delta^T\mathbf{H}\mathbf{X} \in \mathbb{R}^{n \times p}$ , the objective function (C17) can be rewritten as:

$$\operatorname{argmax}_{\mathbf{V}} \operatorname{tr}(\mathbf{V}^T\Psi^T\Psi\mathbf{V}) \quad s.t. \quad \mathbf{V}^T\mathbf{V} = \mathbf{I}, \quad |\mathbf{V}| \leq c. \quad (\text{C19})$$

Using the singular value decomposition (SVD), the matrix  $\Psi$  with  $\operatorname{rank}(\Psi) = m \leq n$  can be written as a product of matrices:

$$\Psi = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T \quad s.t. \quad \mathbf{U}^T\mathbf{U} = \mathbf{I}_n, \quad \mathbf{V}\mathbf{V}^T = \mathbf{I}_p, \quad \mathbf{\Lambda} = I(\lambda_1, \dots, \lambda_m, 0, \dots, 0), \quad (\text{C20})$$

where  $\mathbf{U} \in \mathbb{R}^{n \times n}$  and  $\mathbf{V} \in \mathbb{R}^{p \times p}$  are orthogonal matrices, and  $\mathbf{\Lambda} \in \mathbb{R}^{n \times p}$  is a diagonal matrix, with descending diagonal entries  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0$ . It is easy to see that the columns of  $\mathbf{V}$  are Eigen vectors of the matrix  $\Psi^T\Psi$ , since the following Eigen value decomposition holds:

$$\Psi^T\Psi = \mathbf{V}\mathbf{\Lambda}\mathbf{U}^T\mathbf{U}\mathbf{\Lambda}\mathbf{V}^T = \mathbf{V}(\mathbf{\Lambda}^2)\mathbf{V}^T. \quad (\text{C21})$$

The sparse SPCA problem (C19) now becomes a matrix decomposition problem of the matrix  $\Psi$ , when adding an  $L_1$  penalty on the matrix  $\mathbf{V}$ , since the columns of  $\mathbf{V}$ , being Eigen vectors of  $\Psi^T\Psi$ , maximize  $\operatorname{tr}(\mathbf{V}^T\Psi^T\Psi\mathbf{V})$ .

With an  $L_1$  penalty on  $\mathbf{V}$ , this problem is a *penalized matrix decomposition* problem (PMD, Witten et al. (2009)).

Recalling our original problem of finding interpretable latent variables that also depend on a target variable, the rank  $m$  matrix decomposition of  $\Psi$  may not be desirable. It can be shown (e.g. Eckart and Young 1936) that the best low rank ( $r \leq m$ ) approximation of  $\Psi$  is calculated by the first  $r$  singular values of  $\Lambda$  and the first  $r$  singular vectors of  $\mathbf{U}$  and  $\mathbf{V}$ . With  $\mathbf{u}_i$  being the  $i$ -th column of  $\mathbf{U}$  and  $\mathbf{v}_i$  being the  $i$ -th column of  $\mathbf{V}$ , the best low rank approximation can thus be written as:

$$\sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{v}_i^T = \underset{\hat{\Psi}}{\operatorname{argmin}} \|\Psi - \hat{\Psi}\|_F^2, \tag{C22}$$

subject to the squared Frobenius-norm ( $A \in \mathbb{R}^{m \times n}$ :  $\|A\|_F^2 = \sum_{i=1}^n \sum_{j=1}^m |a_{ij}|^2$ ). The following equality was demonstrated in Witten et al. (2009):

$$\frac{1}{2} \|\Psi - \mathbf{U}\Lambda\mathbf{V}^T\|_F^2 = \frac{1}{2} \|\Psi\|_F^2 - \sum_{i=1}^r \mathbf{u}_i^T \Psi \mathbf{v}_i \lambda_i + \frac{1}{2} \sum_{i=1}^r \lambda_i^2. \tag{C23}$$

The minimization problem (C22) thus becomes a maximization problem, by ignoring the constant terms. Sharifzadeh et al. (2017) added additional  $L_2$  constraints on  $\mathbf{u}_i$  and  $\mathbf{v}_i$ , an  $L_1$  constraint on  $v_i$  for sparsity and an orthogonality constraint for  $u_i$ :

$$\underset{\mathbf{u}_i, \mathbf{v}_i}{\operatorname{argmax}} \mathbf{u}_i^T \Psi \mathbf{v}_i \text{ s.t. } \|\mathbf{u}_i\|_2 \leq 1, \|\mathbf{v}_i\|_2 \leq 1, \|\mathbf{v}_i\|_1 \leq c, \mathbf{u}_i \perp \mathbf{u}_1, \dots, \mathbf{u}_{i-1} \tag{C24}$$

The  $L_2$  constraints do not force unit length to avoid non convex optimization problems. Witten et al. (2009) discuss how to solve many penalized matrix decomposition problems of this kind. Without the orthogonality constraint, they call this particular problem PMD( $\cdot, L_1$ ). The solution to this problem is discussed in detail in Sharifzadeh et al. (2017). A software implementation is available with the R-package PMA by Witten and Tibshirani (2020), which we will use for our demonstrations. Problem (C24) does not yield orthogonal sparse vectors  $\mathbf{v}_i$ , Witten et al. (2009) state that these vectors are unlikely to be very correlated, since the vectors  $\mathbf{v}_i$  are associated with orthogonal vectors  $\mathbf{u}_i, i = 1, \dots, r$ .

### C.3.1 Choice of the Kernel

For sparse SPCA the kernel  $\mathbf{K}$  has been predefined as. The choice of the kernel  $\mathbf{L}$ , however, has a decisive impact on how the dependencies are modeled. Song et al. (2012) discuss the kernel choice for different situations. For binary classification, one may simply choose

$$l(y_i, y_j) = y_i y_j, \text{ where } y_i, y_j \in \{\pm 1\}, \tag{C25}$$

or a weighted version, giving different weights on positive and negative labels. For multiclass classification a possible kernel is

$$l(y_i, y_j) = c_y \delta_{y_i, y_j}, \text{ where } c_y > 0. \tag{C26}$$

For regression one can also use a linear kernel  $l(y_i, y_j) = y_i, y_j$ , but then only simple linear correlations between features and the target variable can be detected. A more universal choice is the radial basis function (RBF) kernel:

$$l(y_i, y_j) = \exp\left(-\frac{\|y_i - y_j\|^2}{2\sigma^2}\right). \quad (\text{C27})$$

The choice of the bandwidth  $2\sigma^2$  is extremely important. For example, if  $2\sigma^2 \rightarrow 0$ , the matrix  $\mathbf{L}$  becomes the identity matrix. Or if  $2\sigma^2 \rightarrow \infty$ , all entries of  $\mathbf{L}$  are 1. In both cases, all relevant information of the dependency between features and the target variable is lost. Besides the bandwidth  $2\sigma$ , the kernel matrix  $L$  depends only on the pairwise distances  $\|y_i - y_j\|^2$ . A reasonable, and heuristically well performing (Pfister et al. 2017) choice is  $2\sigma^2 = \text{median}(\|y_i - y_j\|^2 : i > j)$ . However, it might also be possible and advantageous to use other kernels that are selected to be particularly efficient in detecting certain kinds of dependencies.

### C.3.2 Choice of $c$

Witten et al. (2009) explained how PMD can be used to impute missing data. The main idea is simply to exclude missing entries from the maximization problem (C24) and impute missing values by the low rank approximation matrix  $\mathbf{U}\mathbf{\Lambda}\mathbf{V}^\top$ . This procedure can also be used for finding optimal values for  $c$  by a cross-validation approach. The test data consists of leaving out some entries of the matrix  $\Psi$  (not entire rows or columns, but individual elements of the matrix), yielding a matrix with missing entries  $\tilde{\Psi}$ . For candidate values  $c_i, i = 1, \dots, k$ , calculate the  $\text{PMD}(\cdot, L_1)$  and record the mean squared error over the missing elements of  $\tilde{\Psi}$  and the estimate  $\mathbf{U}\mathbf{\Lambda}\mathbf{V}^\top$ . The true values of the missing values of  $\tilde{\Psi}$  are available in the original data  $\Psi$ . The optimal value  $c^*$  corresponds to the best candidate value  $c_j$ , which minimizes the mean squared error.

However, such a cross-validation approach for the search for  $c$  is not always necessary. If the method is used as a descriptive method to better understand the underlying structure of the data, a small value of  $c$  can be chosen to achieve a desired sparsity.

## Appendix D Feature description for smartphone sensor data

See Table 7.

**Table 7** Description of features used for CFEPs in Sect. 7

Feature	Description
daily_mean_num_unique_Weather_weekend	Mean number of different weather apps used each day on weekends
daily_mean_num_Weather	Mean number of weather apps used each day
daily_mean_num_unique_Weather_week	Mean number of different weather apps used each day on weekdays
daily_mean_num_unique_Weather	Mean number of different weather apps used each day
daily_mean_num_unique_apps	Mean number of different apps used each day
daily_mean_num_unique_apps_week	Mean number of different apps used each day on weekdays
daily_mean_num_unique_apps_weekend	Mean number of different apps used each day on weekends
daily_mean_sum_events_night	Number of all events during the night averaged for each day
daily_mean_dur_all	Duration of all events averaged for each day
daily_sd_sum_intereventall	Sd of the sum of all inter-event time intervals for each day
daily_mean_num_uniq_song	Mean number of different songs listened to each day
daily_mean_num_song	Mean number of songs listened to each day
daily_mean_duration_music	Mean duration of music apps used each day

## References

- Allaire J, Gandrud C, Russell K, et al (2017) networkD3: D3 JavaScript network graphs from R. <https://CRAN.R-project.org/package=networkD3>, R package version 0.4
- Alon U, Barkai N, Notterman DA et al (1999) Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proc Natl Acad Sci* 96(12):6745–6750
- Amoukou SI, Brunel NJB, Salaün T (2021) The shapley value of coalition of variables provides better explanations. [arXiv:2103.13342](https://arxiv.org/abs/2103.13342)
- Apley DW, Zhu J (2019) Visualizing the effects of predictor variables in black box supervised learning models. [arXiv:1612.08468](https://arxiv.org/abs/1612.08468)
- Bair E, Hastie T, Paul D et al (2006) Prediction by supervised principal components. *J Am Stat Assoc* 101(473):119–137
- Barshan E, Ghodsi A, Azimifar Z et al (2011) Supervised principal component analysis: visualization, classification and regression on subspaces and submanifolds. *Pattern Recogn* 44(7):1357–1371
- Berk R, Sherman L, Barnes G et al (2009) Forecasting murder within a population of probationers and parolees: a high stakes application of statistical learning. *J R Stat Soc A Stat Soc* 172(1):191–211
- Breiman L (2001) Random forests. *Mach Learn* 45(1):5–32
- Brenning A (2021) Transforming feature space to interpret machine learning models. [arXiv:2104.04295](https://arxiv.org/abs/2104.04295)
- Caputo B, Sim K, Furesjö F, et al (2002) Appearance-based object recognition using SVMs: Which kernel should I use. In: Proceedings of the NIPS workshop on statistical methods for computational experiments in visual processing and computer vision, Red Hook, NY, USA
- Casalicchio G, Molnar C, Bischl B (2019) Visualizing the feature importance for black box models. Springer International Publishing. *Machine Learning and Knowledge Discovery in Databases*, pp 655–670
- Chakraborty D, Pal NR (2008) Selecting useful groups of features in a connectionist framework. *IEEE Trans Neural Netw* 19(3):381–396

- Cohen SB, Ruppín E, Dror G (2005) Feature selection based on the Shapley value. In: Kaelbling LP, Saffiotti A (eds) IJCAI-05, Proceedings of the nineteenth international joint conference on artificial intelligence, Edinburgh, Scotland, UK, July 30–August 5, 2005. Professional Book Center, pp 665–670
- Covert I, Lundberg SM, Lee SI (2020) Understanding global feature contributions with additive importance measures. *Adv Neural Inf Process Syst* 33:17212–17223
- de Mijolla D, Frye C, Kunesch M, et al (2020) Human-interpretable model explainability on high-dimensional data. *CoRR* [arXiv:2010.07384](https://arxiv.org/abs/2010.07384)
- Eckart C, Young G (1936) The approximation of one matrix by another of lower rank. *Psychometrika* 1(3):211–218
- Fisher A, Rudin C, Dominici F (2019) All models are wrong, but many are useful: learning a variable’s importance by studying an entire class of prediction models simultaneously. *J Mach Learn Res* 20(177):1–81
- Friedman JH (2001) Greedy function approximation: a gradient boosting machine. *Ann Stat*, 1189–1232
- Friedman J, Hastie T, Tibshirani R (2010) A note on the group lasso and a sparse group lasso. [arXiv:1001.0736](https://arxiv.org/abs/1001.0736)
- Fuchs K, Scheipl F, Greven S (2015) Penalized scalar-on-functions regression with interaction term. *Comput Stat Data Anal* 81:38–51
- Fukumizu K, Bach FR, Jordan MI (2004) Dimensionality reduction for supervised learning with reproducing kernel Hilbert spaces. *J Mach Learn Res* 5:73–99
- Goldberg LR (1990) An alternative “description of personality”: the big-five factor structure. *J Person Soc Psychol* 59:1216–1229
- Goldstein A, Kapelner A, Bleich J et al (2013) Peeking inside the black box: visualizing statistical learning with plots of individual conditional expectation. *J Comput Gr Stat* 24:44–65
- Gregorova M, Kalousis A, Marchand-Maillet S (2018) Structured nonlinear variable selection. In: Globerson A, Silva R (eds) Proceedings of the thirty-fourth conference on uncertainty in artificial intelligence, UAI 2018, Monterey, California, USA, August 6–10, 2018. AUAI Press, pp 23–32
- Gregorutti B, Michel P, Saint-Pierre P (2015) Grouped variable importance with random forests and application to multiple functional data analysis. *Comput Stat Data Anal* 90:15–35
- Gretton A, Bousquet O, Smola A, et al (2005) Measuring statistical dependence with Hilbert-Schmidt norms. In: International conference on algorithmic learning theory. Springer, pp 63–77
- Guyon I, Weston J, Barnhill S et al (2002) Gene selection for cancer classification using support vector machines. *Mach Learn* 46(1–3):389–422
- Harari GM, Gosling SD, Wang R et al (2015) Capturing situational information with smartphones and mobile sensing methods. *Eur J Pers* 29(5):509–511
- Harari GM, Lane ND, Wang R et al (2016) Using smartphones to collect behavioral data in psychological science: opportunities, practical considerations, and challenges. *Perspect Psychol Sci* 11(6):838–854
- Harari GM, Müller SR, Aung MS et al (2017) Smartphone sensing methods for studying behavior in everyday life. *Curr Opin Behav Sci* 18:83–90
- Harari GM, Müller SR, Stachl C et al (2019) Sensing sociability: individual differences in young adults’ conversation, calling, texting, and app use behaviors in daily life. *J Person Soc Psychol* 119:204
- He Z, Yu W (2010) Stable feature selection for biomarker discovery. *Comput Biol Chem* 34:215–225
- Hein M, Bousquet O (2004) Kernels, Associated structures and generalizations, Max Planck Institute for Biological Cybernetics
- Shipp MA, Ross KN, Tamayo P et al (2002) Diffuse large b-cell lymphoma outcome prediction by gene-expression profiling and supervised machine learning. *Nat Med* 8(1):68–74
- Hooker G (2004) Discovering additive structure in black box functions. In: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, pp 575–580
- Hooker G (2007) Generalized functional anova diagnostics for high-dimensional functions of dependent variables. *J Comput Graph Stat* 16(3):709–732
- Hooker G, Mentch L (2019) Please stop permuting features: an explanation and alternatives. [arXiv:1905.03151](https://arxiv.org/abs/1905.03151)
- Jackson JJ, Wood D, Bogg T et al (2010) What do conscientious people do? Development and validation of the behavioral indicators of conscientiousness (bic). *J Res Pers* 44(4):501–511
- Jaeger J, Sengupta R, Ruzzo W (2003) Improved gene selection for classification of microarrays. *Pac Symp Biocomput Pac Symp Biocomput* 8:53–64
- Jolliffe IT (1986) Principal component analysis. Springer, New York
- Kolenik T, Gams M (2021) Intelligent cognitive assistants for attitude and behavior change support in mental health: state-of-the-art technical review. *Electronics* 10(11):1250



- Lei J, G'Sell M, Rinaldo A et al (2018) Distribution-free predictive inference for regression. *J Am Stat Assoc* 113(523):1094–1111
- Lipton ZC (2018) The mythos of model interpretability: in machine learning, the concept of interpretability is both important and slippery. *Queue* 16(3):31–57
- Lozano AC, Abe N, Liu Y et al (2009) Grouped graphical granger modeling for gene expression regulatory networks discovery. *Bioinformatics* 25(12):i110–i118
- Lundberg SM, Erion GG, Lee S (2018) Consistent individualized feature attribution for tree ensembles. *CoRR arXiv:1802.03888*
- Lundberg SM, Lee SI (2017) A unified approach to interpreting model predictions. In: Proceedings of the 31st international conference on neural information processing systems. Curran Associates Inc., Red Hook, NY, USA, NIPS'17, pp 4768–4777
- Meier L, Van De Geer S, Bühlmann P (2008) The group lasso for logistic regression. *J R Stat Soc Ser B (Stat Methodol)* 70(1):53–71
- Meinshausen N, Bühlmann P (2010) Stability selection. *J R Stat Soc Ser B (Stat Methodol)* 72(4):417–473
- Miller G (2012) The smartphone psychology manifesto. *Perspect Psychol Sci* 7(3):221–237
- Molnar C (2019) Interpretable machine learning. <https://christophm.github.io/interpretable-ml-book/>
- Molnar C, König G, Bischl B, et al (2020a) Model-agnostic feature importance and effects with dependent features—a conditional subgroup approach. *arXiv:2006.04628*
- Molnar C, König G, Herbringer J, et al (2020b) General pitfalls of model-agnostic interpretation methods for machine learning models. *arXiv preprint arXiv:2007.04131*
- Nicodemus K, Malley J, Strobl C, et al (2010) The behaviour of random forest permutation-based variable importance measures under predictor correlation. *BMC Bioinform* 11–110
- Onnela JP, Rauch SL (2016) Harnessing smartphone-based digital phenotyping to enhance behavioral and mental health. *Neuropsychopharmacology* 41(7):1691–1696
- Ozer DJ, Benet-Martínez V (2006) Personality and the prediction of consequential outcomes. *Annu Rev Psychol* 57:401–421
- Park MY, Hastie T, Tibshirani R (2006) Averaged gene expressions for regression. *Biostatistics* 8(2):212–227
- Pfister N, Bühlmann P, Schölkopf B et al (2017) Kernel-based tests for joint independence. *J R Stat Soc Ser B (Stat Methodol)* 80(1):5–31
- Rachuri KK, Musolesi M, Mascolo C, et al (2010) Emotionsense: a mobile phones based adaptive platform for experimental social psychology research. In: *UbiComp'10—Proceedings of the 2010 ACM conference on ubiquitous computing*
- Raento M, Oulasvirta A, Eagle N (2009) Smartphones: an emerging tool for social scientists. *Social Methods Res* 37(3):426–454
- Rapaport F, Barillot E, Vert JP (2008) Classification of Arraycgh data using fused SVM. *Bioinformatics* 24(13):i375–i382
- Saab S, Lattie EG, Schueller SM et al (2016) The relationship between mobile phone location sensor data and depressive symptom severity. *PeerJ* 4:e2537
- Schoedel R, Au Q, Völkel ST et al (2018) Digital footprints of sensation seeking. *Zeitschrift für Psychologie* 226(4):232–245
- Schoedel R, Pargent F, Au Q et al (2020) To challenge the morning lark and the night owl: using smartphone sensing data to investigate day-night behaviour patterns. *Eur J Personal* 34:733–752
- Scholbeck CA, Molnar C, Heumann C et al (2020) Sampling, intervention, prediction, aggregation: a generalized framework for model-agnostic interpretations. In: Cellier P, Driessens K (eds) *Machine learning and knowledge discovery in databases*. Springer, Cham, pp 205–216
- Schuerk T, Kaltefleiter LJ, Au JQ et al (2019) Enter the wild: autistic traits and their relationship to mentalizing and social interaction in everyday life. *J Autism Dev Disorders* 49:4193–4208
- Seedorff N, Brown G (2021) totalvis: a principal components approach to visualizing total effects in black box models. *SN Comput Sci* 2(3):1–12
- Servia-Rodríguez S, Rachuri KK, Mascolo C, et al (2017) Mobile sensing at the service of mental well-being: A large-scale longitudinal study. In: 26th international world wide web conference, WWW 2017. International World Wide Web Conferences Steering Committee, pp 103–112
- Shapley LS (1953) A value for n-person games. *Contrib Theory Games* 2(28):307–317
- Sharifzadeh S, Ghodsi A, Clemmensen LH et al (2017) Sparse supervised principal component analysis (sspca) for dimension reduction and variable selection. *Eng Appl Artif Intell* 65:168–177

- Song L, Smola A, Gretton A et al (2012) Feature selection via dependence maximization. *J Mach Learn Res* 13:1393–1434
- Song L, Smola A, Gretton A, et al (2007) Supervised feature selection via dependence estimation. In: *Proceedings of the 24th international conference on Machine learning*, pp 823–830
- Stachl C, Hilbert S, Au JQ et al (2017) Personality traits predict smartphone usage. *Eur J Pers* 31(6):701–722
- Stachl C, Au Q, Schoedel R et al (2020a) Predicting personality from patterns of behavior collected with smartphones. *Proc Natl Acad Sci* 117:17680–17687
- Stachl C, Pargent F, Hilbert S et al (2020b) Personality research and assessment in the era of machine learning. *Eur J Personal* 34:613–631
- Strobl C, Boulesteix AL, Kneib T et al (2008) Conditional variable importance for random forests. *BMC Bioinform* 9:307
- Thomé S (2018) Mobile phone use and mental health; A review of the research that takes a psychological perspective on exposure. *Int J Environ Res Public Health* 15(12):2692
- Tibshirani R (1996) Regression shrinkage and selection via the lasso. *J Roy Stat Soc Ser B (Methodol)* 58(1):267–288
- Toloşi L, Lengauer T (2011) Classification with correlated features: unreliability of feature ranking and solutions. *Bioinformatics* 27(14):1986–1994
- Tripathi S, Hemachandra N, Trivedi P (2020) Interpretable feature subset selection: a Shapley value based approach. In: *Proceedings of 2020 IEEE international conference on big data, special session on explainable artificial intelligence in safety critical systems*
- Valentin S, Harkotte M, Popov T (2020) Interpreting neural decoding models using grouped model reliance. *PLOS Comput Biol* 16(1):e1007148
- Venables B, Ripley B (2002) *Modern applied statistics with S*
- Watson DS, Wright MN (2019) Testing conditional independence in supervised learning algorithms. [arXiv:1901.09917](https://arxiv.org/abs/1901.09917)
- Williamson BD, Gilbert PB, Simon NR, et al (2020) A unified approach for inference on algorithm-agnostic variable importance. [arXiv:2004.03683](https://arxiv.org/abs/2004.03683)
- Williamson B, Feng J (2020) Efficient nonparametric statistical inference on population feature importance using Shapley values. In: *International conference on machine learning*, PMLR, pp 10282–10291
- Witten D, Tibshirani R (2020) PMA: penalized multivariate analysis. *R Package Vers* 1(2):1
- Witten DM, Tibshirani R, Hastie T (2009) A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics* 10(3):515–534
- Wold S, Albano C, Dunn WJ et al (1984) *Multivariate data analysis in chemistry*. Springer, Dordrecht, pp 17–95
- Yarkoni T, Westfall J (2017) Choosing prediction over explanation in psychology: lessons from machine learning. *Perspect Psychol Sci* 12(6):1100–1122
- Yuan M, Lin Y (2006) Model selection and estimation in regression with grouped variables. *J R Stat Soc Ser B (Stat Methodol)* 68(1):49–67

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.