

Nonlinear Feature Aggregation: Two Algorithms driven by Theory

Anonymous authors

Paper under double-blind review

Abstract

Many real-world machine learning applications are characterized by a huge number of features, leading to computational and memory issues, as well as the risk of overfitting. Ideally, only relevant and non-redundant features should be considered to preserve the complete information of the original data and limit the dimensionality. Dimensionality reduction and feature selection are common preprocessing techniques addressing the challenge of efficiently dealing with high-dimensional data. Dimensionality reduction methods control the number of features in the dataset while preserving its structure and minimizing information loss. Feature selection aims to identify the most relevant features for a task, discarding the less informative ones. Previous works have proposed approaches that aggregate features depending on their correlation without discarding any of them and preserving their interpretability through aggregation with the mean. A limitation of methods based on correlation is the assumption of linearity in the relationship between features and target. In this paper, we relax such an assumption in two ways. First, we propose a bias-variance analysis for general models with additive Gaussian noise, leading to a dimensionality reduction algorithm (NonLinCFA) which aggregates non-linear transformations of features with a generic aggregation function. Then, we extend the approach assuming that a generalized (non-)linear model regulates the relationship between features and target. A deviance analysis leads to a second dimensionality reduction algorithm (GenLinCFA), applicable to a larger class of regression problems and classification settings. Finally, we test the algorithms on synthetic and real-world datasets, performing regression and classification tasks, showing competitive performances.

1 Introduction

Dimensionality reduction is an essential technique in *machine learning* (ML), employed to limit the number of features or dimensions of datasets. It has been successfully applied in a large variety of fields where high-dimensional data needs to be analyzed, classified, visualized, or interpreted. For instance, in computer vision and in natural language processing, data is often represented as high-dimensional vectors; in bioinformatics, high-throughput biological data like DNA sequences are often represented as high-dimensional data; in earth sciences, meteorological variables have high-dimensional measurements in different locations. In these contexts, datasets are often characterized by a large number of highly correlated features. Projecting them into a lower dimensional space is crucial to simplify data representation and enhance the performance of models, mitigating overfitting, and limiting the computational complexity. However, dimensionality reduction may lead to loss of information and interpretability.

Another method to reduce the dimension of the feature space is *feature selection*, which aims to identify the most relevant features from a dataset. The importance of feature selection lies in its ability to improve model performance, reduce computational resources, and simplify the model's interpretation. Its main desirable property compared to dimensionality reduction is interpretability, since the reduced features are simply a subset of the original ones. However, these techniques usually discard features that may be exploited to reduce the variance.

In Bonetti et al. (2023), the authors propose a dimensionality reduction approach that *aggregates* subsets of features with their mean, if this is convenient in terms of bias-variance tradeoff. In this way, the algorithm preserves the interpretability of the transformed features, while exploiting the information of each feature. Assuming linearity in the underlying data generation process and applying linear regression, the asymptotic analysis on which the algorithm is based on suggests aggregating two features with their mean if their correlation is *sufficiently large*. Compared to other dimensionality reduction techniques, the aggregation based on the mean preserves the interpretability of the newly generated feature, since the mean function is a transformation that a domain expert can understand without any additional explanation by ML experts (see Kovalerchuk et al. (2021) for the definition of interpretability in these terms).

Contributions In this work, after introducing the notation and formulation of the problem and reviewing the main dimensionality reduction methods (Section 2), we design principled algorithms that generalize this approach in two ways (Section 3). First, we relax the assumption of linearity of the relationship between the features and the target. In Section 4, we propose a dimensionality reduction technique that still considers linear regression as a supervised learning technique, but we allow the inputs to be *generic non-linear transformations* of the original features and the aggregation function to be a *generic non-linear function* (instead of the mean as in Bonetti et al. (2023)). This way, we allow the designer to ponder the suitable balance between the interpretability of the result (e.g., preferring simple feature and aggregation function) and the incorporation of complex non-linear relationships. Second, we analyse *generalized non-linear models*, assuming the expected value of the target to be a generic transformation of the features mapped through the link function (Section 5). We extend the proposed algorithm in this context, which is particularly useful when the Gaussianity assumption of the noise model is unrealistic or in the case of classification problems. Finally, we apply the two algorithms to classification and regression datasets, showing competitive results w.r.t. state-of-the-art methods (Section 6). We conclude the paper in Section 7 summarising its contributions and possible future developments. The paper is accompanied by an appendix. In particular, Appendix A and Appendix C respectively contain proofs and technical results related to Section 4 and Section 5. Appendix B shows an additional bi-variate analysis for generalized linear models, which is the starting point of the analysis that leads to the algorithm presented in the main paper. Finally, Appendix D presents additional experiments and gives more details on the experiments discussed in the main paper.

2 Preliminaries

Notation Given N samples and D features, let $\mathbf{X} \in \mathbb{R}^{N \times D}$ and $\mathbf{y} \in \mathbb{R}^N$ (resp. $\mathbf{y} \in \{0, 1\}^N$ for classification) be the feature matrix and target vector, where $P_{X,Y}$ denotes the joint probability distribution. The j -th row of the matrix \mathbf{X} is denoted with \mathbf{x}_j , while its i -th element is denoted with x_i , and it is called *feature*. Similarly, y_j is the j -th element of the target vector \mathbf{y} . σ_a^2 , $cov(a, b)$, $\rho_{a,b}$ and $\hat{\sigma}_a^2$, $\hat{cov}(a, b)$, $\hat{\rho}_{a,b}$ denote the variance of a random variable a , its covariance, and correlation with the random variable b and their estimators, respectively. Finally, $\mathbb{E}_a[f(a)]$ and $var_a(f(a))$ are the expected value and the variance of a function $f(\cdot)$ of the random variable a w.r.t. its distribution.

Data Generation Processes We consider three scenarios for describing the data generation processes.

- *Non-linear model*: we consider a general non-linear relationship f between the features \mathbf{x} and the target y , with additive Gaussian noise:

$$y = f(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2). \quad (1)$$

- *Generalized linear model*: we assume the distribution of the target to be part of the canonical exponential family:

$$Y|X \sim \exp\left(\frac{y\boldsymbol{\theta} - b(\boldsymbol{\theta})}{\phi}\right) + c(y, \phi), \quad (2)$$

where $\boldsymbol{\theta} = \mathbf{X}\mathbf{w}$ is the parameter of the family, ϕ , $b(\cdot)$ and $c(\cdot)$ are respectively a known scale parameter, a function of the parameter $\boldsymbol{\theta}$ characterizing the assumed distribution and a normalizing

function independent from θ . In this setting, the expected value of the target is $\mathbb{E}_\theta[y] = b'(\theta)$. Moreover, in generalized linear models, the expected value of the target given the features is a nonlinear transformation of the linear combination of their values:

$$\mu(\mathbf{x}) = \mathbb{E}[y|\mathbf{x}] = g^{-1}(\mathbf{x}^\top \mathbf{w}) \iff g(\mu(\mathbf{x})) = \mathbf{x}^\top \mathbf{w}, \quad (3)$$

where $g(\cdot)$ is the *canonical link function*. A complete analysis of generalized linear models can be found in (McCullagh, 2019).

- *Generalized non-linear model*: to further generalize the approach, we consider a generic non-linear transformation $f(\cdot)$ inside the link function ($\theta = f(\mathbf{x})$), such that the parameter θ is not constrained to be a linear transformation of the features:

$$\mu(\mathbf{x}) = \mathbb{E}[y|\mathbf{x}] = g^{-1}(f(\mathbf{x})) \iff g(\mu(\mathbf{x})) = f(\mathbf{x}). \quad (4)$$

Dimensionality Reduction via Aggregation A dimensionality reduction algorithm maps the original D -dimensional feature matrix \mathbf{X} into a reduced dataset $\mathbf{U} = \psi(\mathbf{X}) \in \mathbb{R}^{N \times d}$ with $d < D$, being ψ a transformation function designed to maximize a specific performance measure. In this paper we consider non-linear transformations of the original features $\phi_i(X) : \mathbb{R}^D \rightarrow \mathbb{R}$ and a generic aggregation function $h(\cdot) : \mathbb{R}^2 \rightarrow \mathbb{R}$, and we aggregate couples of them through the aggregation function $h(\phi_i(X), \phi_j(X))$.

The choice to consider both non-linear transformation of the original features $\phi_i(\cdot)$ as inputs and a generic user-defined aggregation function $h(\cdot)$ is made to produce theoretical results that hold in a general setting. Some considerations follows:

- In line with the work of (Bonetti et al., 2023), the main applicative interest is to consider the D features $\{x_1, \dots, x_D\}$ as inputs and the mean as aggregation function. This way, the original (continuous) features are aggregated via a simple transformation, preserving the interpretability of the reduced features. However, in some applications, a non-linear relationship between some features and the target may be known, and the proposed algorithms can handle this situation as well. In the experimental section we will provide an example of this with synthetic data, where a quadratic relationship between the features and the target holds, by considering quadratic inputs $\phi_i(X) = x_i^2$ and the mean as aggregation function h .
- The choice of considering firstly non-linear transformations of the features and then an aggregation function may be replaced by a unique function that aggregates the original features non-linearly combining them (e.g., when a quadratic relationship holds it is possible to directly consider the sum of squares as unique transformation function ψ). However, this choice is made to underline the two distinct meanings of these functions. Firstly, the non-linear transformations $\phi_i(\cdot)$ of the inputs are user-defined modifications of the original features that allow a linear regression model to handle non-linear relationships and are not aimed to reduce the dimension. Then, the aggregation function $h(\cdot)$ is a characteristic of the proposed algorithms that is responsible to perform the dimensionality reduction via aggregation. Indeed, the proposed algorithms allow to decide the aggregation function that is more meaningful for a specific application, or to test multiple aggregation functions to select the best performing one.

In this context, we assume ϵ to be a noise signal, independent of X , and we denote with \hat{w}_i and \hat{y} the estimated coefficients and the predicted target. Finally, to simplify the computations, we assume each expected value to be zero: $\mathbb{E}[\phi_i(X)] = \mathbb{E}[Y] = \mathbb{E}[h(\phi(X))] = \mathbb{E}[f(X)] = 0$.

Performance Indexes The Mean Squared Error (*MSE*), that is usually adopted as loss measure in regression problems that assume Gaussianity, can be decomposed into three terms (bias-variance decomposition (Hastie et al., 2009)):

$$\begin{aligned} \underbrace{\mathbb{E}_{x,y,\mathcal{T}}[(\mathcal{M}_{\mathcal{T}}(x) - y)^2]}_{\text{MSE}} &= \underbrace{\mathbb{E}_{x,\mathcal{T}}[(\mathcal{M}_{\mathcal{T}}(x) - \bar{\mathcal{M}}(x))^2]}_{\text{variance}} \\ &\quad + \underbrace{\mathbb{E}_x[(\bar{\mathcal{M}}(x) - \bar{y}(x))^2]}_{\text{bias}} + \underbrace{\mathbb{E}_{x,y}[(\bar{y}(x) - y)^2]}_{\text{noise}}, \end{aligned} \quad (5)$$

with x, y features and target of a test sample, \mathcal{T} training set, $\mathcal{M}_{\mathcal{T}}(\cdot)$ ML model trained with \mathcal{T} , $\bar{\mathcal{M}}(\cdot)$ its expected value w.r.t. \mathcal{T} and \bar{y} expected value of the test target y w.r.t. the input features x . We will consider the *MSE* to guarantee that the advantage in terms of reduction of variance (due to the dimensionality reduction) is larger than the disadvantage in terms of increase of bias.

The *deviance* is a measure of goodness of fit, usually considered in generalized linear models to extend the *MSE* to non-Gaussian settings (McCullagh, 2019). Recalling that in our setting $\theta = f(\mathbf{X})$ (or $\theta = \mathbf{X}\mathbf{w}$ assuming linearity), the deviance measures how the likelihood of the estimated model deviates from the best one:

$$D^*(\theta, \hat{\theta}) := \frac{D(\theta, \hat{\theta})}{\phi} = -2 \left[\ell(\hat{\theta}) - \ell(\theta) \right] = \frac{2}{\phi} [y(\theta - \hat{\theta}) - (b(\theta) - b(\hat{\theta}))], \quad (6)$$

where D is the deviance, D^* is the scaled deviance, $\ell(\theta), \ell(\hat{\theta})$ are the log-likelihood of the model and of the estimated one, ϕ and $b(\cdot)$ are the known scale parameter and specific function of the distribution assumed. The last equation holds only if the distribution of the target belongs to the *canonical* exponential family, which will be assumed throughout the paper. When dealing with generalized (non-)linear models we will therefore consider the expected value of the scaled deviance to compare two different models. Performing dimensionality reduction, the deviance decreases when the reduction of variance due to the aggregation is convenient w.r.t. the loss of information.

2.1 Dimensionality Reduction Methods

This section contains a literature survey on dimensionality reduction. More extensive surveys can be found in (Van Der Maaten et al., 2009; Sorzano et al., 2014; Ayesha et al., 2020) and applicative results in (Zebari et al., 2020; Espadoto et al., 2021).

Linear Dimensionality Reduction Principal Components Analysis (PCA) (Pearson, 1901; Hotelling, 1933) is a popular linear dimensionality reduction method that embeds high dimensional data into a linear subspace, trying to preserve the variance of the original dataset. This method performs linear projections of possibly all the original features with different coefficients: interpretability can be difficult and it may suffer from the curse of dimensionality. Several linear dimensionality reduction approaches have been proposed to overcome the issues of PCA or to perform particular projections. Among these, svPCA (Ulfarsson & Solo, 2011) forces most of the weights of the projection to be zero, improving the interpretability but discarding the contribution of many features. Independent Component Analysis (ICA) (Hyvarinen, 1999) is an information-theoretic approach that looks for independent components designed for multichannel data. Locality Preserving Projections (LPP) (He & Niyogi, 2003) is an unsupervised linear method that tries to preserve the local structure of the original data. Linear Discriminant Analysis (LDA) (Fisher, 1936) is a supervised technique that finds a linear combination of features that identifies the projection that separates the target classes. A broader overview of linear dimensionality reduction techniques can be found in (Cunningham & Ghahramani, 2015) and applicative results of linear methods in (Reddy et al., 2020).

Non-linear Dimensionality Reduction The limits of linear projections have been overcome by resorting to non-linear methods. Kernel PCA (Shawe-Taylor & Cristianini, 2004) is a variation of PCA that allows non-linearity by combining PCA with a kernel. Sammon Mapping (Sammon, 1969) is an unsupervised algorithm that maps the data trying to preserve pairwise distance among data globally, similarly to Multidimensional scaling (Kruskal & Wish, 1978). Isomap (Tenenbaum et al., 2000) is an extension of MDS aimed to preserve the geodesic distance among features. Locally Linear Embedding (LLE) (Roweis & Saul, 2000) is a method aimed at preserving the distance among data, focusing on local similarity. More recently, many applications have embedded the dimensionality reduction phase in the learning process of a neural network, typically applying convolutional neural networks or autoencoders (Hinton & Salakhutdinov, 2006; Zebari et al., 2020). Among supervised methods (Chao et al., 2019), Supervised PCA (Bair et al., 2006; Barshan et al., 2011) projects the data onto a subspace where the features are uncorrelated, simultaneously maximizing the dependency between the reduced dataset and the target. NMF-based algorithms (Jing et al., 2012; Lu et al., 2017) focus on the non-negativity property of features, which is not a general property of applicative problems. Finally, manifold-based methods perform supervised non-linear projections, usually adjusting an unsupervised approach to take into consideration the target (Neighborhood Components

Analysis (NCA) (Goldberger et al., 2004), supervised Isomap (Ribeiro et al., 2008; Zhang et al., 2018) and supervised LLE (Zhang, 2009) are examples of these approaches).

3 Proposed Algorithms: a Methodological Perspective

Non-Linear Correlated Features Aggregation Starting from the dimensionality reduction algorithm proposed in Bonetti et al. (2023), we introduce a similar approach relaxing the linearity assumptions, named *Non-Linear Correlated Features Aggregation (NonLinCFA)*. Let the relationship between the D features x_i and the target y be non-linear ($y = f(\mathbf{x})$), with additive Gaussian noise (Equation 1). We iteratively compare, in terms of Mean Squared Error (MSE), two linear regression models. Given two non-linear transformations of the original features $\phi_1(\mathbf{x}), \phi_2(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}$, that allow to account for non-linearities in the linear regression, the first model considers them as separate inputs, while the second aggregates them with a generic aggregation function $h(\cdot) : \mathbb{R}^2 \rightarrow \mathbb{R}$:

$$\begin{cases} \hat{y} = \hat{w}_1 \phi_1(\mathbf{x}) + \hat{w}_2 \phi_2(\mathbf{x}) \\ \hat{y} = \hat{w} h(\phi_1(\mathbf{x}), \phi_2(\mathbf{x})) = \hat{w} h(\phi(\mathbf{x})). \end{cases} \quad (7)$$

Intuitively, the second model has similar or better performances if the correlation between each input $\phi_1(\mathbf{x}), \phi_2(\mathbf{x})$ and the function $f(\mathbf{x})$ that generates the target ($\rho_{\phi_1, f}, \rho_{\phi_2, f}$) is not much larger than the one between their aggregation $h(\phi_1(\mathbf{x}), \phi_2(\mathbf{x}))$ and $f(\mathbf{x})$ ($\rho_{h(\phi_1, \phi_2), f}$). In this case, it is convenient to reduce the dimensionality of the inputs by aggregating them through the function $h(\cdot)$.

Remark 1 (About linear features ϕ). *Considering the case $\phi_1(x) = x_i, \phi_2(x) = x_j, h(\phi_1, \phi_2) = \frac{\phi_1 + \phi_2}{2}$, the algorithm aggregates the two variables x_i, x_j with their means if it is convenient in terms of MSE. In this case, the performance is preserved, and the aggregation keeps the reduced feature interpretable. Generally, the algorithm allows any zero-mean non-linear transformation of the input features and aggregation function, to balance the complexity and interpretability of the dimensionality reduction to be adapted to the specific problem.*

Remark 2 (About considering two features at a time). *Considering two inputs at a time can be convenient for huge-dimensional input spaces. The proposed algorithm iteratively compares couples of inputs, identifies a couple to aggregate and considers that aggregation as a single feature for the subsequent iterations.*

Generalized-Linear Correlated Features Aggregation The second algorithm we propose, *Generalized-Linear Correlated Features Aggregation (GenLinCFA)*, can be applied to problems where the target does not follow Gaussian distributions, including some classification problems. We assume the target to follow a distribution in the canonical exponential family, and we consider generalized linear models with canonical link function. Firstly, we compare the expected deviance of two models in the bivariate setting. Assuming that the expected value of the target is a linear combination of the two features, transformed by the canonical link ($\mathbb{E}[y|\mathbf{x}] = g^{-1}(w_1 x_1 + w_2 x_2)$), the first model considers the two features x_1, x_2 separately, while the second model considers their aggregation with the mean $\frac{x_1 + x_2}{2}$:

$$\begin{cases} \hat{y} = g^{-1}(\hat{w}_1 x_1 + \hat{w}_2 x_2) \\ \hat{y} = g^{-1}(\hat{w} \cdot \frac{x_1 + x_2}{2}), \end{cases} \quad (8)$$

where $\hat{w}_1, \hat{w}_2, \hat{w}$ are the coefficients estimated from data by the two models.

Then, we generalize this approach assuming a non-linear relationship between D features and the expected value of the target ($\mathbb{E}[y|\mathbf{x}] = g^{-1}(f(\mathbf{x}))$). We analyse again the performances of two models, in terms of expected deviance:

$$\begin{cases} \hat{y} = g^{-1}(\hat{w}_1 \phi_1(x) + \hat{w}_2 \phi_2(x)) \\ \hat{y} = g^{-1}(\hat{w} \cdot h(\phi_1(x), \phi_2(x))). \end{cases} \quad (9)$$

As discussed for the first algorithm, we consider two nonlinear transformations $\phi_1(\cdot), \phi_2(\cdot)$ of the features as inputs and we compare a model that considers them separately with a model that aggregates them through a

nonlinear function $h(\cdot)$. Again, the algorithm allows choosing different non-linear transformations of features and aggregation depending on the problem, and it iteratively compares couples of inputs.

4 Non-Linear Correlated Features Aggregation

This section describes the first proposed algorithm, *NonLinCFA*, designed to aggregate features in regression problems. We consider a problem with D features and a non-linear relationship between the features and the target (Equation 1).

Given a training dataset with N samples, $\mathbf{X} \in \mathbb{R}^{N \times D}$, we compare the *MSE* of the two models described in Equation 7. As discussed in Section 3, we, therefore, compare a bivariate with a univariate linear regression, where the two nonlinear functions of features $\phi_1(x), \phi_2(x)$ are aggregated through a function $h(\cdot)$. The variance decreases due to the reduction of the dimension of the hypothesis space and the bias increases due to the aggregation. Therefore, we analyse these two quantities, that together with the irreducible error compose the *MSE* (Equation 5), to guarantee that it does not increase significantly after the aggregation.

4.1 Theoretical Analysis

We firstly show the decrease of variance due to the aggregation in the asymptotic case. The finite sample result and the proofs can be found in Appendix A.

Theorem 1. *Let the relationship between the features and the target be nonlinear with additive Gaussian noise (Equation 1). Let also each estimator to converge in probability to the quantity that it estimates. In the asymptotic case, the decrease of variance between a bivariate linear regression and the univariate case where the two features are aggregated (Equation 7) is:*

$$\Delta_{var}^{n \rightarrow \infty} = \frac{\sigma^2}{(n-1)}. \quad (10)$$

Remark 3. *The asymptotic result of Equation 10 follows the intuition that there is more variability in the prediction of a bivariate regression than in the univariate case if there is a small number of samples or a high variability of the noise.*

As a second result, we show the asymptotic increase of bias due to aggregation. A finite-sample result and related proofs are reported in Appendix A.

Theorem 2. *Let the relationship between features and target be nonlinear with additive Gaussian noise (Equation 1) and each estimator converge in probability. The asymptotic increase of bias between bivariate and univariate linear regression, with the two features aggregated with a function $h(\cdot)$ (Equation 7) is:*

$$\begin{aligned} \Delta_{bias}^{n \rightarrow \infty} = & -\frac{\text{cov}(f, h(\phi_1, \phi_2))^2}{\sigma_{h(\phi_1, \phi_2)}^2} \\ & + \frac{\sigma_{\phi_1}^2 \text{cov}(\phi_2, f)^2 + \sigma_{\phi_2}^2 \text{cov}(\phi_1, f)^2 - 2\text{cov}(\phi_1, f)\text{cov}(\phi_2, f)\text{cov}(\phi_1, \phi_2)}{\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}(\phi_1, \phi_2)^2}. \end{aligned} \quad (11)$$

Remark 4. *Intuitively, the asymptotic result of Equation 11 suggests that it is convenient to aggregate the two features $\phi_1(\mathbf{x}), \phi_2(\mathbf{x})$ if there is large covariance between their aggregation and the target $f(\mathbf{x})$, or the covariance between the two features and the target is small, considering them singularly.*

We can now introduce the main theoretical result of this section: the asymptotic guarantee that ensures the *MSE* to not worsen after the aggregation.

Theorem 3. *Let the relationship between features and target be nonlinear with additive Gaussian noise (Equation 1) and each estimator converge in probability. The asymptotic MSE of a bivariate linear regression is not greater than the univariate case with the two inputs $\phi_1(\mathbf{x}), \phi_2(\mathbf{x})$ aggregated with a function $h(\cdot)$ (Equation 7) if and only if:*

$$\frac{\sigma^2}{\sigma_f^2(n-1)} \geq \frac{\rho_{\phi_1, f}^2 + \rho_{\phi_2, f}^2 - 2\rho_{\phi_1, f}\rho_{\phi_2, f}\rho_{\phi_1, \phi_2}}{1 - \rho_{\phi_1, \phi_2}^2} - \rho_{f, h(\phi_1, \phi_2)}^2 = R_{f, \phi_1 \phi_2}^2 - R_{f, h(\phi_1, \phi_2)}^2. \quad (12)$$

Proof. The result follows imposing $\Delta_{var}^{n \rightarrow \infty} \geq \Delta_{bias}^{n \rightarrow \infty}$ from Equation 10 and 11. \square

Remark 5. *Intuitively, the aggregation is convenient if:*

- *there is much noise in the model or a small number of samples (first term);*
- *real model and aggregated feature share a lot of information (third term);*
- *the two features do not share much information with the real model (small second term). Indeed, $R_{f, \phi_1 \phi_2}^2$ is the coefficient of multiple correlation, indicating how well the target can be linearly predicted with a set of features (Keith, 2019).*

Remark 6. *In the bivariate linear case, where $\phi_1 = x_1$, $\phi_2 = x_2$, $h(\phi_1, \phi_2) = \frac{x_1 + x_2}{2}$, $f(x) = w_1 x_1 + w_2 x_2$, Equation 7 becomes $\rho_{x_1, x_2} \geq 1 - \frac{2\sigma^2}{(n-1)(w_1 - w_2)^2}$. This is in line with the result found in Bonetti et al. (2023) with linearity assumptions.*

Remark 7. *Assuming unitary variances $\sigma_{\phi_1}^2 = \sigma_{\phi_2}^2 = 1$ and uncorrelated combinations of features $\rho_{\phi_1, \phi_2} = 0$, the right hand side of Equation 12 becomes $\rho_{\phi_1, f}^2 + \rho_{\phi_2, f}^2 - \rho_{f, h(\phi_1, \phi_2)}^2$. This follows the intuition that the correlation between the real model and the aggregation should not be worse than the correlation between the real model and each individual transformation of features $\phi_i(\mathbf{x})$.*

4.2 Proposed algorithm: NonLinCFA

Starting from the result of Equation 12, this section introduces the algorithm proposed in this paper to perform dimensionality reduction in regression settings. The algorithm focuses on identifying if it is possible to add one feature in an aggregation, iteratively comparing a bivariate with a univariate approach. Algorithm 1 shows the pseudo-code of the proposed algorithm *Non-Linear Correlated Features Aggregation* (NonLinCFA). The main function of the algorithm (*NonLinCFA*) generates d partitions of the D inputs ϕ_1, \dots, ϕ_D : at each iteration, it calls an auxiliary function (*compute_threshold*) to evaluate the performance in terms of coefficient of determination (*R2score*) of two models.

Firstly the bivariate model is composed aggregating inputs already assigned to the current partition $h(\phi_{\mathcal{P}})$ and a feature not already assigned to any partition ϕ_j . Then, the univariate linear regression considers the aggregation of the already selected inputs and the one under analysis ($h(\phi_{\mathcal{P}}, \phi_j)$) as a single feature. Their difference is an estimate of the right hand side of Equation 12: if it is *sufficiently small*, the algorithm adds the input ϕ_j to the current cluster. The hyperparameter ϵ regulates the propensity of the algorithm to aggregate features: when is large the algorithm is prone to aggregate, while small values are more conservative. This hyperparameter substitutes the left hand side term of Equation 12, which is difficult to be estimated in practice since it depends on the variance of the real model. Finally, the algorithm computes the aggregations of each element of the identified partition \mathcal{P} , returning the set of aggregated features $\{\bar{h}_{\phi}^1, \dots, \bar{h}_{\phi}^d\}$.

As discussed in the previous sections, a particularly meaningful version of the algorithm can be obtained considering as inputs the D features $\{x_1, \dots, x_D\}$ and the mean as aggregation. This way, the algorithm identifies groups of features to aggregate with their mean, preserving the interpretability.

5 Generalized-Linear Correlated Features Aggregation

This section describes a second algorithm, *GenLinCFA*, that relaxes the Gaussianity assumption. We consider D features and we assume the model in the canonical exponential family (Equation 2). The expected value of the target is a general function of the features $f(\mathbf{x})$, transformed by the linking function $g(\cdot)$ (Equation 4). In a first analysis we assume the function $f(\cdot)$ to be linear, modeling the expected value of the target as in Equation 3, and we compare the bivariate case with separate features and the univariate case with their aggregation (Equation 8). The bivariate analysis and some additional results are shown in Appendix B. As a second more general approach, given a training dataset with N samples $\mathbf{X} \in \mathbb{R}^{N \times D}$, we compare the *expected deviance* of the two models described in Equation 9. We therefore compare a bivariate with a univariate regression, where two nonlinear functions of features $\phi_1(x), \phi_2(x)$ are aggregated through a function $h(\cdot)$. The main difference is the additional nonlinear transformation of the predicted linear model through the linking function $g(\cdot)$ (Equation 9), which characterizes the generalized non-linear model under analysis.

Algorithm 1 NonLinCFA: Non-Linear Correlated Features Aggregation

Input: D combinations of features $\{\phi_1(x), \dots, \phi_D(x)\}$; target y ; N samples, aggregation function $h(\cdot)$, tolerance ϵ

Output: reduced features $\{\bar{h}_\phi^1, \dots, \bar{h}_\phi^d\}$, with $d \leq D$

```

function COMPUTE_THRESHOLD( $h, \phi_P, \phi_j, y, \epsilon$ )                                ▷ Threshold from Equation 12
     $R1 \leftarrow \text{R2score}(h(\phi_P), \phi_j, y)$ 
     $R2 \leftarrow \text{R2score}(h(\phi_P, \phi_j), y)$ 
    return  $R2 - R1 - \epsilon$ 
end function

function NONLINCFA( $Input$ )                                                        ▷ Main function
     $\mathcal{P} \leftarrow \{\}$                                                             ▷ Partition of the features
     $\mathcal{V} \leftarrow \{\}$                                                             ▷ Set of already considered features
    for each  $i \in \{1, \dots, D\}$  do
        if  $i \notin \mathcal{V}$  then
             $\mathcal{P} \leftarrow \{i\}$ 
             $\mathcal{V} \leftarrow \mathcal{V} \cup \{i\}$ 
            for each  $j \in \{i+1, \dots, D\}$  do
                 $threshold \leftarrow \text{COMPUTE\_THRESHOLD}(h, \phi_P, \phi_j, y, \epsilon)$ 
                if  $threshold \leq 0$  then                                          ▷ Aggregate the features
                     $\mathcal{P} \leftarrow \mathcal{P} \cup \{j\}$ 
                     $\mathcal{V} \leftarrow \mathcal{V} \cup \{j\}$ 
                end if
            end for
             $\mathcal{P} \leftarrow \mathcal{P} \cup \{\mathcal{P}\}$ 
        end if
    end for
     $d \leftarrow |\mathcal{P}|$ 
    for each  $k \in \{1, \dots, d\}$  do
         $\bar{h}_\phi^k = h(\phi_{\mathcal{P}_k})$ 
    end for
    return  $\{\bar{h}_\phi^1, \dots, \bar{h}_\phi^d\}$ 
end function

```

5.1 Theoretical analysis

Starting from Equation 6, given two estimators $\hat{\theta}, \bar{\theta}$ of the parameter θ , the increase of expected deviance between the two models is:

$$\mathbb{E}_{\mathbf{x}, y, \mathcal{T}}[D^*(\theta, \hat{\theta}) - D^*(\theta, \bar{\theta})] = \frac{2}{\phi} \mathbb{E}_{\mathbf{x}, y, \mathcal{T}}[y(\bar{\theta} - \hat{\theta}) - (b(\bar{\theta}) - b(\hat{\theta}))]. \quad (13)$$

Considering the two models under analysis (Equation 9), recalling that we defined $\theta = f(\mathbf{x})$ and that the two estimators of θ in the two cases are $\hat{\theta} = \hat{w}h(\phi_1, \phi_2)$ and $\bar{\theta} = \hat{w}_1\phi_1 - \hat{w}_2\phi_2$, the expected increase of deviance becomes:

$$\begin{aligned} & \frac{2}{\phi} \left\{ \mathbb{E}_{\mathbf{x}, y} [y \cdot (\mathbb{E}_{\mathcal{T}}[\hat{w}_1]\phi_1 + \mathbb{E}_{\mathcal{T}}[\hat{w}_2]\phi_2 - \mathbb{E}_{\mathcal{T}}[\hat{w}]h(\phi_1, \phi_2))] \right. \\ & \left. - \mathbb{E}_{\mathbf{x}, y, \mathcal{T}} [b(\hat{w}_1\phi_1 + \hat{w}_2\phi_2) - b(\hat{w}h(\phi_1, \phi_2))] \right\}. \end{aligned} \quad (14)$$

The following theorem provides the main theoretical result of this setting: a second-order approximated upper bound of the quantity under analysis, whose derivation can be found in Appendix C.

Theorem 4. *Let M, m be the largest and smallest absolute value of the following expected values of coefficients: $\mathbb{E}_{\mathcal{T}}[\hat{w}_1]$, $\mathbb{E}_{\mathcal{T}}[\hat{w}_2]$, $\mathbb{E}_{\mathcal{T}}[\hat{w}]$, $\mathbb{E}_{\mathcal{T}}[\hat{w}_1^2]$, $\mathbb{E}_{\mathcal{T}}[\hat{w}_2^2]$, $\mathbb{E}_{\mathcal{T}}[\hat{w}_1\hat{w}_2]$, $\mathbb{E}_{\mathcal{T}}[\hat{w}^2]$. Let the real model belong to the*

canonical exponential family and $\Delta(D^*)$ be the expected increase of deviance due to the aggregation of the two transformed features $\phi_1(x), \phi_2(x)$ through an aggregation function $h(\cdot)$:

$$\Delta(D^*) \leq \frac{2}{\phi} \left\{ \left[M \cdot (|\text{cov}(\phi_1, y)| + |\text{cov}(\phi_2, y)|) - m \cdot |\text{cov}(h(\phi_1, \phi_2), y)| \right] - \frac{1}{2} b''(0) \cdot \left(m \cdot \sigma_{\phi_1 + \phi_2}^2 - M \cdot \sigma_{h(\phi_1, \phi_2)}^2 \right) \right\}. \quad (15)$$

Remark 8. Theorem 4 follows the intuition that it is convenient to aggregate two inputs when the variance of their sum is large or the variance of their aggregation $h(\phi_1, \phi_2)$ is small. Moreover, it is convenient to aggregate when the absolute value of the covariance between each feature and the target is small or the absolute value between the aggregated feature and the target is large.

Remark 9. Assuming a Gaussian distribution of the target given the features, the asymptotic increase of the expected deviance is equal to the asymptotic increase of MSE due to the aggregation. The proof of this considerations and additional technical results can be found in Appendix C.

5.2 Proposed algorithm: GenLinCFA

Following the theoretical analysis performed, the second algorithm proposed in this paper is *GenLinCFA*: similarly to *NonLinCFA*, it iteratively compares tuples of inputs to decide if it is convenient, in terms of expected deviance, to substitute them with their aggregation through a user-defined function $h(\cdot)$.

The proposed algorithm is a variation of Algorithm 1. Its peculiarity is the different way to compute the threshold that suggests that it is convenient to aggregate two inputs. Algorithm 2 shows the pseudo-code of the function *compute_threshold*, which is the only difference w.r.t. Algorithm 1. Indeed, starting from the result of Equation 15, the algorithm aggregates two features if the upper bound of the increase of expected deviance due to the aggregation is smaller than 0. The variance and covariance of the quantities that appear in the upper bound are estimated from data, while the constant $\frac{m}{M}$ is replaced by an hyperparameter ϵ , which regulates the propensity of the algorithm to aggregate. Large values give to the algorithm more propensity to aggregate, while small values require more information provided by the aggregation or more noise in the original features to perform the aggregation.

As for *NonLinCFA*, a specific version of the algorithm that preserves interpretability is to consider the inputs equal to the D features $\{x_1, \dots, x_D\}$ and the mean as aggregation function $h(x) = \frac{1}{|x|} \sum_{x_i \in x} x_i$.

Algorithm 2 GenLinCFA: Generalized-Linear Correlated Features Aggregation

Input: D combinations of features $\{\phi_1(x), \dots, \phi_D(x)\}$; target y ; N samples, aggregation function $h(\cdot)$, parameter ϵ

Output: reduced features $\{\bar{h}_\phi^1, \dots, \bar{h}_\phi^d\}$, with $d \leq D$

function COMPUTE_THRESHOLD($h, \phi_P, \phi_j, y, \epsilon$) ▷ Threshold from Equation 15

$L \leftarrow |\text{cov}(\phi_P, y)| + |\text{cov}(\phi_j, y)| + \frac{1}{2} b''(0) \hat{\sigma}_{h(\phi)}^2$

$R \leftarrow |\text{cov}(h(\phi), y)| + \frac{1}{2} b''(0) \hat{\sigma}_{\phi_P + \phi_j}^2$

return $L - \epsilon R$

end function

function GENLINCFA(*Input*) ▷ Main function

Equal to *NonLinCFA* function in Algorithm 1

The only difference is the renewed COMPUTE_THRESHOLD function

return Reduced features as described in *Algorithm 1*

end function

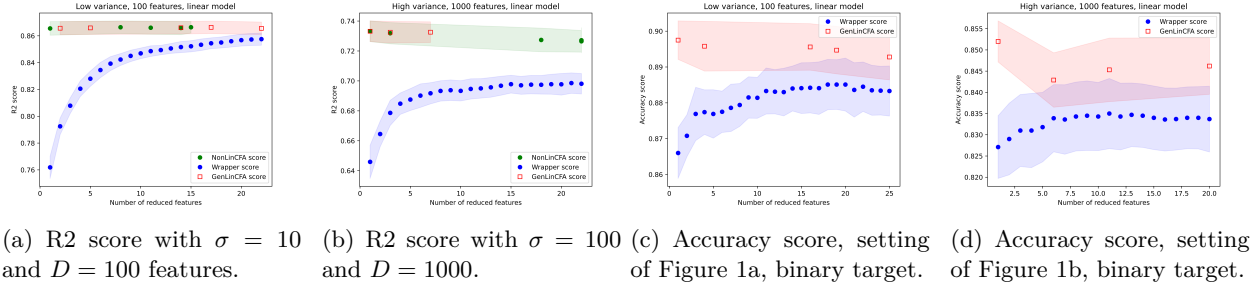


Figure 1: Application of NonLinCFA and GenLinCFA in regression and classification. Confidence intervals show test performances with different numbers of reduced features considering different hyperparameters.

6 Experiments

This section describes the experiments performed on synthetic and real-world datasets to empirically validate the proposed algorithms. Additional details on datasets, methodologies and results can be found in Appendix D. Code and datasets can be found at the following anonymous link <https://www.dropbox.com/s/eoqyrs3o0ymh4o4/nonlinearFeatureAggregation.zip?dl=0> and they will be made publicly available on github upon acceptance.

6.1 Synthetic experiments

In this subsection, we validate the proposed algorithms with synthetic experiments. We design two regression problems with $n = 3000$ samples (randomly considering 2000 samples for training and 1000 samples for testing) and we repeat the experiments ten times to produce confidence intervals. In the first regression setting, $D = 100$ linearly correlated features are considered, designing a target that linearly depends on each of them, with additive Gaussian noise with standard deviation $\sigma = 10$. In the second linear regression setting, more features are considered ($D = 1000$) and a more variable noise is added to the target ($\sigma = 100$). To validate GenLinCFA also in classification settings, we repeat the two experiments applying the sign function to the target, which transforms it into a binary variable. More details can be found in Appendix D.

Figure 1 shows the results of the application of the methods in terms of R^2 score on the test set. In both settings, NonLinCFA and GenLinCFA have been applied, considering the mean as an aggregation function. Then, linear regression was performed on the reduced features to evaluate the performances. As a baseline, a wrapper forward feature selection has been applied, with a number of features up to the maximum number of features extracted by the two methods. The experiments were repeated 10 times to produce confidence intervals. From Figure 1a, it is possible to see that the two algorithms extract a comparable number of features, and the performance is already satisfactory with a small number of reduced features, while many features are needed for the wrapper method to reach similar results. Figure 1b reports the results in a more noisy environment with more features. With the same hyperparameters, NonLinCFA and GenLinCFA are respectively less and more prone to aggregate. As an additional comparison, LinCFA algorithm has been applied, since in this setting we are assuming the linearity of the model. The two proposed algorithms have similar performances w.r.t. LinCFA, but with the advantage to be more prone to further reduce the dimension.

To test the GenLinCFA algorithm also in a classification setting, the same two experiments are repeated, applying the sign function to the target and evaluating performances in terms of test accuracy (compared again with a wrapper forward feature selection). From Figure 1c-1d, it is possible to conclude that, again, the performances are similar or better w.r.t. the wrapper baseline and that, considering the same values of the hyperparameter, the algorithm is more prone to aggregate in a more complex and noisy context. Appendix D reports detailed results and confidence intervals of the four experiments performed as well as additional synthetic experiments showing that the two proposed algorithms can deal with non-linear transformations of the input features and with nonlinear aggregations (in particular, a quadratic relationship

between the features and the target is analysed, both considering the original features and the sum of squares as aggregation function or the squared features and the mean as aggregation function).

6.2 Real World Experiments

To conclude the empirical analysis of the proposed algorithms, this section reports some experiments on real datasets. This analysis has been conducted evaluating test performances in terms of the coefficient of determination of the linear regression (or accuracy of the logistic regression for classification tasks), considering the reduced features as inputs. The two dimensionality reduction methods introduced in this paper have been performed with different hyperparameters, considering the original features as inputs ($\phi_i(\mathbf{x}) = x_i$) and performing the mean as aggregation function, which preserves the interpretability of the reduced features. The results are compared with different state-of-the-art dimensionality reduction methods with different characteristics (discussed in Section 2.1): linear unsupervised (PCA), linear supervised (LDA, LinCFA), nonlinear unsupervised (kernel PCA, Isomap, LLE) and nonlinear supervised (supervised PCA, NCA). For all the approaches, the first 66% of the samples of each dataset has been considered to for training and validation purposes and the remaining 33% to test the results. Confidence intervals on the test performances have been obtained bootstrapping the training and validation set with five different seeds.

Table 1 and 2 respectively report a subset of test performances of three datasets from Kaggle and the UCI ML repository and four datasets extracted by the authors from meteorological data. In particular, The *Finance* dataset has been retrieved from Kaggle¹, while the *Bankruptcy*² and *Parkinson*³ datasets have been retrieved from the UCI ML repository. The first climate dataset (*Climate*, *Climate(Class.)*) considers temperature, precipitation (retrieved from (Didan, 2015; Cornes et al., 2018)) and the state of vegetation (retrieved from (Zellner, 2022)) of neighbouring basins as features to predict the state of vegetation of a sub-basin of the Po River. The second dataset only focuses on temperature and precipitation features. For each climate dataset, two versions are considered: the original regression problem and a binary classification task obtained considering the values below and above the average respectively as class 0 and 1, representing conditions of water scarcity and abundance. These dataset are an example of the main applicative interest of the authors, which is to aggregate features representing measurements of the same variables at different locations through their mean, which remains interpretable (e.g., the mean of temperature measurements over a sub-region has a clear meaning for a climatologist) and significantly reduces the dimension and the autocorrelation among the reduced features. In both tables are reported the number of reduced dimensions and the performance test score, considering the best validation hyperparameter for NonLinCFA and GenLinCFA as well as the performance of LinCFA and of the best performing algorithm among the aforementioned state-of-the-art methods considered. Some NA values associated with *NonLinCFA* and *LinCFA* are due to the fact that the two algorithms cannot be applied on classification tasks, which are considered to further investigate the *GenLinCFA* algorithm. The complete results related to all the baseline algorithms and to the different choices of hyperparameters can be found in Appendix D.

From the results it is possible to conclude that the proposed algorithms have competitive performances, with the advantage of preserving the interpretability of the reduced features, that are aggregated with the mean. In some cases (*Finance*, *Climate*, *Climate(Class.)*) the proposed algorithms outperform the compared methods, while in other situations they perform similarly (*Bankruptcy*, *Parkinson*). In the second meteorological case (*Climate II*, *Climate II(class.)*) the Kernel PCA algorithm has better performances, showing that in some settings it is necessary to balance between the loss of information and the interpretability of the reduced features.

7 Conclusions

In this paper, we have deepened the study of dimensionality reduction to account for non-linear effects, focusing on preserving both information and interpretability. The non-linearity has been accounted for in both the deterministic mapping function and the noise model, considering the exponential family of

¹<https://www.kaggle.com/datasets/dgawlik/nyse>

²<https://archive.ics.uci.edu/ml/datasets/Polish+companies+bankruptcy+data>

³<https://archive.ics.uci.edu/dataset/470/parkinson+s+disease+classification>

Table 1: Experiments on real world datasets from Kaggle and UCI ML repository. Total number of samples n splitted into train (66%) and test (33%) sets.

Quantity	Finance	Bankruptcy (class.)	Parkinson (class.)
# samples n	1299	1084	384
# features D	75	65	753
Reduced Dimension			
NonLinCFA	7.4 ± 0.4	NA	NA
GenLinCFA	8.0 ± 1.4	27.6 ± 5.6	23.4 ± 1.1
LinCFA	11.4 ± 0.7	NA	NA
Best baseline	36.0 ± 10.8	12.0 ± 9.5	28.2 ± 15.3
Test performance	R² score	Accuracy score	Accuracy score
NonLinCFA	0.8136 ± 0.0036	NA	NA
GenLinCFA	0.8119 ± 0.0010	0.7503 ± 0.0012	0.8016 ± 0.0069
LinCFA	0.8010 ± 0.0128	NA	NA
Best baseline	0.7764 ± 0.0118	0.7637 ± 0.0079	0.7952 ± 0.0181

Table 2: Experiments on climate datasets. The total number of samples n has been splitted into train (66% of data) and test (33% of data) sets.

	Climate	Climate (Class.)	Climate II	Climate II (class.)
# samples n	981	981	867	867
# features D	1991	1991	2408	2408
Reduced Dimension				
NonLinCFA	16.0 ± 0.8	NA	7.0 ± 0.9	NA
GenLinCFA	13.8 ± 3.0	17.5 ± 3.3	7.2 ± 1.1	26.0 ± 6.1
LinCFA	38.2 ± 1.6	NA	222.0 ± 2.7	NA
Best baseline	41.8 ± 2.5	37.0 ± 6.6	21.8 ± 9.5	11.1 ± 2.3
Test performance	R² score	Accuracy score	R² score	Accuracy score
NonLinCFA	0.9395 ± 0.0125	NA	0.2949 ± 0.0156	NA
GenLinCFA	0.9275 ± 0.0004	0.9107 ± 0.0022	0.2841 ± 0.0051	0.7127 ± 0.0159
LinCFA	0.9007 ± 0.0310	NA	-1.2861 ± 0.2322	NA
Best baseline	0.8454 ± 0.0049	0.8827 ± 0.0098	0.3889 ± 0.0199	0.7640 ± 0.0062

distributions. The resulting algorithms aggregate, in the most general case, non-linear aggregation of non-linear features. Theoretical results have been provided to investigate the performance of the aggregation either in terms of MSE or increase of deviance. The experimental validation illustrates that our algorithms outperform the proposed baselines both in synthetically generated environments and in a real-world domain, or they have competitive results, with the advantage of letting the user define the most suitable aggregation function that, in most cases, has been selected as the mean to preserve the interpretability of the reduced features. Future works will include the consideration of additional indexes of statistical dependence, other than the correlation and covariance, to perform the aggregation (e.g., mutual information). Additionally, the proposed algorithms will be further applied to investigate their impact on the detection of the state of vegetation of European river basins.

References

- Shaeela Ayesha, Muhammad Kashif Hanif, and Ramzan Talib. Overview and comparative study of dimensionality reduction techniques for high dimensional data. *Information Fusion*, 59:44–58, 2020.
- Eric Bair, Trevor Hastie, Debashis Paul, and Robert Tibshirani. Prediction by supervised principal components. *Journal of the American Statistical Association*, 101(473):119–137, 2006. doi: 10.1198/016214505000000628.

- Elnaz Barshan, Ali Ghodsi, Zohreh Azimifar, and Mansoor Jahromi. Supervised principal component analysis: Visualization, classification and regression on subspaces and submanifolds. *Pattern Recognition*, 44: 1357–1371, 07 2011. doi: 10.1016/j.patcog.2010.12.015.
- Paolo Bonetti, Alberto Maria Metelli, and Marcello Restelli. Interpretable linear dimensionality reduction based on bias-variance analysis, 2023.
- Guoqing Chao, Yuan Luo, and Weiping Ding. Recent advances in supervised dimension reduction: A survey. *Machine Learning and Knowledge Extraction*, 1(1):341–358, 2019. ISSN 2504-4990. doi: 10.3390/make1010020.
- Richard C. Cornes, Gerard van der Schrier, Else J. M. van den Besselaar, and Philip D. Jones. An ensemble version of the e-obs temperature and precipitation data sets. *Journal of Geophysical Research: Atmospheres*, 123(17):9391–9409, 2018. doi: <https://doi.org/10.1029/2017JD028200>.
- John P. Cunningham and Zoubin Ghahramani. Linear dimensionality reduction: Survey, insights, and generalizations. *Journal of Machine Learning Research*, 16(89):2859–2900, 2015.
- K. Didan. Myd13q1 modis/aqua vegetation indices 16-day l3 global 250m sin grid v006 [data set]. *NASA EOSDIS Land Processes DAAC.*, 2015. doi: 10.1109/72.761722.
- Mateus Espadoto, Rafael Messias Martins, Andreas Kerren, Nina Sumiko Tomita Hirata, and Alexandru Cristian Telea. Toward a quantitative survey of dimension reduction techniques. *IEEE Transactions on Visualization and Computer Graphics*, 27:2153–2173, 2021.
- Ronald A Fisher. The use of multiple measurements in taxonomic problems. *Annals of eugenics*, 7(2): 179–188, 1936.
- Jacob Goldberger, Sam T. Roweis, Geoffrey E. Hinton, and Ruslan Salakhutdinov. Neighbourhood components analysis. pp. 513–520, 2004.
- T. Hastie, R. Tibshirani, and J.H. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer series in statistics. 2009. ISBN 9780387848846.
- Xiaofei He and Partha Niyogi. Locality preserving projections. *Advances in neural information processing systems*, 16, 2003.
- G. E. Hinton and R. R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006. doi: 10.1126/science.1127647.
- Harold Hotelling. Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology*, 24:498–520, 1933.
- A. Hyvarinen. Fast and robust fixed-point algorithms for independent component analysis. *IEEE Transactions on Neural Networks*, 10(3):626–634, 1999. doi: 10.1109/72.761722.
- Liping Jing, Chao Zhang, and Michael K. Ng. Snmfca: Supervised nmf-based image classification and annotation. *IEEE Transactions on Image Processing*, 21(11):4508–4521, 2012. doi: 10.1109/TIP.2012.2206040.
- Timothy Z Keith. *Multiple regression and beyond: An introduction to multiple regression and structural equation modeling*. Routledge, 2019.
- Boris Kovalerchuk, Muhammad Aurangzeb Ahmad, and Ankur Teredesai. Survey of explainable machine learning with visual and granular methods beyond quasi-explanations. *Interpretable artificial intelligence: A perspective of granular computing*, pp. 217–267, 2021.
- Joseph B Kruskal and Myron Wish. *Multidimensional scaling*, volume 11. Sage, 1978.

- Yuwu Lu, Zhihui Lai, Yong Xu, Xuelong Li, David Zhang, and Chun Yuan. Nonnegative discriminant matrix factorization. *IEEE Transactions on Circuits and Systems for Video Technology*, 27(7):1392–1405, 2017. doi: 10.1109/TCSVT.2016.2539779.
- Peter McCullagh. *Generalized linear models*. Routledge, 2019.
- Karl F.R.S. Pearson. Liii. on lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11):559–572, 1901. doi: 10.1080/14786440109462720.
- G. Thippa Reddy, M. Praveen Kumar Reddy, Kuruva Lakshmana, Rajesh Kaluri, Dharmendra Singh Rajput, Gautam Srivastava, and Thar Baker. Analysis of dimensionality reduction techniques on big data. *IEEE Access*, 8:54776–54788, 2020. doi: 10.1109/ACCESS.2020.2980942.
- Bernardete Ribeiro, Armando Vieira, and João Carvalho das Neves. Supervised isomap with dissimilarity measures in embedding learning. In *Progress in Pattern Recognition, Image Analysis and Applications*. Springer Berlin Heidelberg, 2008. ISBN 978-3-540-85920-8.
- Sam T. Roweis and Lawrence K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000. doi: 10.1126/science.290.5500.2323.
- J.W. Sammon. A nonlinear mapping for data structure analysis. *IEEE Transactions on Computers*, C-18(5):401–409, 1969. doi: 10.1109/T-C.1969.222678.
- John Shawe-Taylor and Nello Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004. doi: 10.1017/CBO9780511809682.
- C. O. S. Sorzano, J. Vargas, and A. Pascual Montano. A survey of dimensionality reduction techniques, 2014.
- Joshua Tenenbaum, Vin Silva, and John Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290:2319–2323, 01 2000.
- Magnus Ulfarsson and Victor Solo. Vector l0 sparse variable pca. *Signal Processing, IEEE Transactions on*, 59:1949 – 1958, 05 2011. doi: 10.1109/TSP.2011.2112653.
- Laurens Van Der Maaten, Eric Postma, and Jaap Van den Herik. Dimensionality reduction: a comparative review. *J Mach Learn Res*, 10:66–71, 2009.
- Rizgar Ramadhan Zebari, Adnan Mohsin Abdulazeez, Diyar Qader Zeebaree, Dilovan Asaad Zebari, and Jwan Najeeb Saeed. A comprehensive review of dimensionality reduction techniques for feature selection and feature extraction. 2020.
- Peter Zellner. Vegetation health index - 231 m 8 days (version 1.0) [data set]. *Eurac Research*, 2022. doi: <https://doi.org/10.48784/161b3496-534a-11ec-b78a-02000a08f41d>.
- Shi-qing Zhang. Enhanced supervised locally linear embedding. *Pattern Recognition Letters*, 30:1208–1218, 10 2009. doi: 10.1016/j.patrec.2009.05.011.
- Yan Zhang, Zhao Zhang, Jie Qin, Li Zhang, Bing Li, and Fanzhang Li. Semi-supervised local multi-manifold isomap by linear embedding for feature extraction. *Pattern Recognition*, 76, 2018. ISSN 0031-3203.

A Non-Linear Correlated Features Aggregation: additional proofs and results

This section contains additional results and proofs related to Section 4 of the main paper.

Firstly, we introduce the finite-sample increase of variance due to the aggregation, of which we reported the asymptotic version in Theorem 1 of the main paper.

Theorem 5. *Let the relationship between the features and the target be nonlinear with additive Gaussian noise (Equation 1 of the main paper). The decrease of variance between a bivariate linear regression and the univariate case where the two features are aggregated (Equation 7 of the main paper) is:*

$$\Delta_{var} = \frac{\sigma^2}{(n-1)} \left[\frac{\sigma_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 + \sigma_{\phi_2}^2 \hat{\sigma}_{\phi_1}^2 - 2cov(\phi_1, \phi_2)c\hat{ov}(\phi_1, \phi_2)}{(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - c\hat{ov}(\phi_1, \phi_2)^2)} - \frac{\sigma_{h(\phi_1, \phi_2)}^2}{\hat{\sigma}_{h(\phi_1, \phi_2)}^2} \right]. \quad (16)$$

Proof of Theorem 5

Recalling the result:

$$\Delta_{var} = \frac{\sigma^2}{(n-1)} \left[\frac{\sigma_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 + \sigma_{\phi_2}^2 \hat{\sigma}_{\phi_1}^2 - 2cov(\phi_1, \phi_2)c\hat{ov}(\phi_1, \phi_2)}{(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - c\hat{ov}(\phi_1, \phi_2)^2)} - \frac{\sigma_{h(\phi(x))}^2}{\hat{\sigma}_{h(\phi(x))}^2} \right],$$

we will compute the variance for the bivariate case and for the univariate case, and then compute their difference. To do so, we need to start by computing the variance of the estimators.

Lemma 1. *In the one dimensional case $\hat{y} = \hat{w}h(\phi(x))$:*

$$var_D(\hat{w}|\mathbf{X}) = \frac{\sigma^2}{(n-1)\hat{\sigma}_{h(\phi(x))}^2}.$$

In the two dimensional case $\hat{y} = \hat{w}_1\phi_1(x) + \hat{w}_2\phi_2(x)$:

$$var_D(\hat{w}|\mathbf{X}) = \frac{\sigma^2}{(n-1)(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - c\hat{ov}(\phi_1, \phi_2)^2)} \begin{bmatrix} \hat{\sigma}_{\phi_2}^2 & -c\hat{ov}(\phi_1, \phi_2) \\ -c\hat{ov}(\phi_1, \phi_2) & \hat{\sigma}_{\phi_1}^2 \end{bmatrix}.$$

Proof. For the one dimensional result, denoting with $\mathbf{h}(\phi)$ the N dimensional vector of aggregated samples $[h(\phi_1(x^1), \phi_2(x^1)), \dots, h(\phi_1(x^N), \phi_2(x^N))]$:

$$\begin{aligned} var_D(\hat{w}|\mathbf{X}) &= var_D((\mathbf{h}(\phi)^\top \mathbf{h}(\phi))^{-1} \mathbf{h}(\phi)^\top y | \mathbf{X}) = (\mathbf{h}(\phi)^\top \mathbf{h}(\phi))^{-1} \sigma^2 \\ &= [h(\phi(x^1)) \quad \dots \quad h(\phi(x^n))] \begin{bmatrix} h(\phi(x^1)) \\ \dots \\ h(\phi(x^n)) \end{bmatrix} \sigma^2 \\ &= \frac{\sigma^2}{\sum_{i=1}^N h(\phi(x^i))^2} = \frac{\sigma^2}{(n-1)\hat{\sigma}_{h(\phi(x))}^2}. \end{aligned}$$

In the first equality it is exploited the closed formula to estimate the linear regression coefficients in linear regression, the second equality exploits the property of the variance to extract the constant matrices and the definition of variance of the target y . Finally, the third and fourth equalities are simple algebraic computations.

In the two dimensional setting, denoting with $\Phi = \begin{bmatrix} \phi_1(x^1) & \phi_2(x^1) \\ \dots & \dots \\ \phi_1(x^n) & \phi_2(x^n) \end{bmatrix}$ the $N \times 2$ matrix of samples:

$$var_D(\hat{w}|\mathbf{X}) = (\Phi^\top \Phi)^{-1} \sigma^2 = \left(\begin{bmatrix} \phi_1(x^1) & \dots & \phi_1(x^n) \\ \phi_2(x^1) & \dots & \phi_2(x^n) \end{bmatrix} \begin{bmatrix} \phi_1(x^1) & \phi_2(x^1) \\ \dots & \dots \\ \phi_1(x^n) & \phi_2(x^n) \end{bmatrix} \right)^{-1} \sigma^2$$

$$\begin{aligned}
&= \left(\begin{bmatrix} \phi_1(x^1)^2 + \dots + \phi_1(x^n)^2 & \phi_1(x^1)\phi_2(x^1) + \dots + \phi_1(x^n)\phi_2(x^n) \\ \phi_1(x^1)\phi_2(x^1) + \dots + \phi_1(x^n)\phi_2(x^n) & \phi_2(x^1)^2 + \dots + \phi_2(x^n)^2 \end{bmatrix} \right)^{-1} \sigma^2 \\
&= \frac{\sigma^2}{(n-1)(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - \hat{cov}(\phi_1, \phi_2)^2)} \begin{bmatrix} \hat{\sigma}_{\phi_2}^2 & -\hat{cov}(\phi_1, \phi_2) \\ -\hat{cov}(\phi_1, \phi_2) & \hat{\sigma}_{\phi_1}^2 \end{bmatrix}.
\end{aligned}$$

The first equivalence follows again from the closed form solution of linear regression and from the variance of the target, while the others follow from algebraic computations. \square

We are now ready to derive the expression of variance of the model in the two cases.

Theorem 6. *In the one dimensional case $\hat{y} = \hat{w}g(\phi(x))$:*

$$\mathbb{E}[(h_D(x) - \bar{h}(x))^2 | \mathbf{X}] = \frac{\sigma^2}{(n-1)} \cdot \frac{\sigma_{h(\phi(x))}^2}{\hat{\sigma}_{h(\phi(x))}^2}.$$

In the two dimensional case $\hat{y} = \hat{w}_1\phi_1(x) + \hat{w}_2\phi_2(x)$:

$$\mathbb{E}[(h_D(x) - \bar{h}(x))^2 | \mathbf{X}] = \frac{\sigma^2}{(n-1)} \cdot \frac{\sigma_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 + \sigma_{\phi_2}^2 \hat{\sigma}_{\phi_1}^2 - 2cov(\phi_1, \phi_2)\hat{cov}(\phi_1, \phi_2)}{(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - \hat{cov}(\phi_1, \phi_2)^2)}.$$

Proof. One dimensional, exploiting the independence between training, test samples and the definition of model variance and the assumption of expected values equal to zero of $h(\phi(x))$:

$$\begin{aligned}
\mathbb{E}_{X,Y,D}[(\mathcal{M}_D(x) - \bar{\mathcal{M}}(x))^2] &= \mathbb{E}_{X,Y,D}[(\hat{w}h(\phi(x)) - \mathbb{E}_D[\hat{w}h(\phi(x))])^2] \\
&= \mathbb{E}_{X,Y,D}[(h(\phi(x))(\hat{w} - \mathbb{E}_D[\hat{w}]))^2] \\
&= \mathbb{E}_{X,Y}[h(\phi(x))^2] \mathbb{E}_D[(\hat{w} - \mathbb{E}_D[\hat{w}])^2] \\
&= var_X(h(\phi(x))) \cdot var_D(\hat{w}).
\end{aligned}$$

Conditioning on the features training set:

$$\mathbb{E}_{X,Y,D}[(\mathcal{M}_D(x) - \bar{\mathcal{M}}(x))^2 | \mathbf{X}] = \sigma_{h(\phi(x))}^2 \cdot \frac{\sigma^2}{(n-1)\hat{\sigma}_{h(\phi(x))}^2}.$$

Two dimensional, exploiting again the independence between train and test set and the assumption of expected values equal to zero of $\phi_1(x)$ and $\phi_2(x)$:

$$\begin{aligned}
\mathbb{E}_{X,Y,D}[(\mathcal{M}_D(x) - \bar{\mathcal{M}}(x))^2] &= \mathbb{E}_{X,Y,D}[(\hat{w}_1\phi_1 + \hat{w}_2\phi_2 - \mathbb{E}_D[\hat{w}_1\phi_1 + \hat{w}_2\phi_2])^2] \\
&= \mathbb{E}_{X,Y,D}[(\phi_1(\hat{w}_1 - \mathbb{E}_D[\hat{w}_1]) + \phi_2(\hat{w}_2 - \mathbb{E}_D[\hat{w}_2]))^2] \\
&= \mathbb{E}_{X,Y,D}[(\phi_1(\hat{w}_1 - \mathbb{E}_D[\hat{w}_1]))^2] + \mathbb{E}_{X,Y,D}[(\phi_2(\hat{w}_2 - \mathbb{E}_D[\hat{w}_2]))^2] \\
&\quad + 2\mathbb{E}_{X,Y,D}[\phi_1\phi_2(\hat{w}_1 - \mathbb{E}_D[\hat{w}_1])(\hat{w}_2 - \mathbb{E}_D[\hat{w}_2])] \\
&= var_X(\phi_1)var_D(\hat{w}_1) + var_X(\phi_2)var_D(\hat{w}_2) + 2cov_X(\phi_1, \phi_2)cov_D(\hat{w}_1, \hat{w}_2).
\end{aligned}$$

Conditioning on the features training set:

$$\begin{aligned}
&\mathbb{E}_{X,Y,D}[(\mathcal{M}_D(x) - \bar{\mathcal{M}}(x))^2 | \mathbf{X}] \\
&= var_X(\phi_1 | \mathbf{X})var_D(\hat{w}_1 | \mathbf{X}) + var_X(\phi_2 | \mathbf{X})var_D(\hat{w}_2 | \mathbf{X}) \\
&\quad + 2cov_X(\phi_1, \phi_2 | \mathbf{X})cov_D(\hat{w}_1, \hat{w}_2 | \mathbf{X}) \\
&= \frac{\sigma^2}{(n-1)} \cdot \frac{\sigma_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 + \sigma_{\phi_2}^2 \hat{\sigma}_{\phi_1}^2 - 2cov(\phi_1, \phi_2)\hat{cov}(\phi_1, \phi_2)}{(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - \hat{cov}(\phi_1, \phi_2)^2)}.
\end{aligned}$$

\square

In conclusion, Theorem 5 trivially follows by computing the difference between the two results of the theorem, i.e., between the variance of the two and the one dimensional settings.

Moreover Theorem 1, in the main paper, follows from this result substituting each estimator with the quantity that it estimates, since it is assume to be convergent in probability to its value.

Proof of Theorem 2

Recalling the asymptotic result of the theorem:

$$\Delta_{bias}^{n \rightarrow \infty} = -\frac{cov(f, h(\phi(x)))^2}{\sigma_{h(\phi(x))}^2} + \frac{\sigma_{\phi_1}^2 cov(\phi_2, f)^2 + \sigma_{\phi_2}^2 cov(\phi_1, f)^2 - 2cov(\phi_1, f)cov(\phi_2, f)cov(\phi_1, \phi_2)}{\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - cov(\phi_1, \phi_2)^2},$$

we firstly need to derive the expected value of the estimators.

Lemma 2. *The expected value of the coefficient of the one dimensional regression is:*

$$\mathbb{E}_D[\hat{w}|\mathbf{X}] = \frac{cov(h(\phi(x)), f(x))}{\hat{\sigma}_{h(\phi(x))}^2}.$$

For the two dimensional regression:

$$\mathbb{E}_D[\hat{w}|\mathbf{X}] = \frac{1}{(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - cov(\phi_1, \phi_2)^2)} \begin{bmatrix} \hat{\sigma}_{\phi_2}^2 cov(\phi_1, f) - cov(\phi_1, \phi_2)cov(\phi_2, f) \\ \hat{\sigma}_{\phi_1}^2 cov(\phi_2, f) - cov(\phi_1, \phi_2)cov(\phi_1, f) \end{bmatrix}.$$

Proof. In the one-dimensional case, exploiting the closed form solution of the estimate of the linear regression coefficient:

$$\begin{aligned} \mathbb{E}_D[\hat{w}|\mathbf{X}] &= (h(\Phi)^T h(\Phi))^{-1} h(\Phi)^T f(\mathbf{X}) \\ &= \frac{1}{(n-1)\hat{\sigma}_{h(\phi(x))}^2} \begin{bmatrix} h(\phi(x^1)) & \dots & h(\phi(x^n)) \end{bmatrix} \begin{bmatrix} f(x^1) \\ \dots \\ f(x^n) \end{bmatrix} \\ &= \frac{cov(h(\phi(x)), f(x))}{\hat{\sigma}_{h(\phi(x))}^2}. \end{aligned}$$

In the two dimensional case, substituting the expression of the estimated coefficients:

$$\begin{aligned} \mathbb{E}_D(\hat{w}|\mathbf{X}) &= (\Phi^T \Phi)^{-1} \Phi^T f(\mathbf{X}) \\ &= \left(\begin{bmatrix} \phi_1(x^1) & \dots & \phi_1(x^n) \\ \phi_2(x^1) & \dots & \phi_2(x^n) \end{bmatrix} \begin{bmatrix} \phi_1(x^1) & \phi_2(x^1) \\ \dots & \dots \\ \phi_1(x^n) & \phi_2(x^n) \end{bmatrix} \right)^{-1} \Phi^T f(\mathbf{X}) \\ &= \frac{1}{(n-1)(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - cov(\phi_1, \phi_2)^2)} \\ &\quad \times \begin{bmatrix} \hat{\sigma}_{\phi_2}^2 & -cov(\phi_1, \phi_2) \\ -cov(\phi_1, \phi_2) & \hat{\sigma}_{\phi_1}^2 \end{bmatrix} \begin{bmatrix} \phi_1(x^1) & \dots & \phi_1(x^n) \\ \phi_2(x^1) & \dots & \phi_2(x^n) \end{bmatrix} \begin{bmatrix} f(x^1) \\ \dots \\ f(x^n) \end{bmatrix} \\ &= \frac{1}{(n-1)(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - cov(\phi_1, \phi_2)^2)} \\ &\quad \times \begin{bmatrix} \hat{\sigma}_{\phi_2}^2 \sum_i f(x^i) \phi_1(x^i) - cov(\phi_1, \phi_2) \sum_i f(x^i) \phi_2(x^i) \\ -cov(\phi_1, \phi_2) \sum_i f(x^i) \phi_1(x^i) + \hat{\sigma}_{\phi_1}^2 \sum_i f(x^i) \phi_2(x^i) \end{bmatrix} \end{aligned}$$

$$= \frac{1}{(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - \hat{cov}(\phi_1, \phi_2)^2)} \left[\hat{\sigma}_{\phi_2}^2 \hat{cov}(\phi_1, f) - \hat{cov}(\phi_1, \phi_2) \hat{cov}(\phi_2, f) \right].$$

□

The next step of the proof is to derive the expression of (squared) bias of the two linear regression settings.

Theorem 7. *Let $h = h(\phi(x))$, $f = f(x)$.*

In the one dimensional linear regression $\hat{y} = \hat{w}h$ the squared bias is:

$$\mathbb{E}[(\bar{\mathcal{M}}(x) - \bar{y})^2 | \mathbf{X}] = \sigma_f^2 + \frac{\sigma_h^2 \hat{cov}(f, h)^2 - 2cov(f, h) \hat{cov}(f, h) \hat{\sigma}_h^2}{\hat{\sigma}_h^4}.$$

In the two dimensional case $\hat{y} = \hat{w}_1 \phi_1 + \hat{w}_2 \phi_2$ the squared bias is:

$$\begin{aligned} \mathbb{E}[(\bar{\mathcal{M}}(x) - \bar{y})^2 | \mathbf{X}] &= \sigma_{\phi_1}^2 \mathbb{E}_D[\hat{w}_1 | \mathbf{X}]^2 + \sigma_{\phi_2}^2 \mathbb{E}_D[\hat{w}_2 | \mathbf{X}]^2 \\ &\quad + 2cov(\phi_1, \phi_2) \mathbb{E}_D[\hat{w}_1 | \mathbf{X}] \mathbb{E}_D[\hat{w}_2 | \mathbf{X}] + \sigma_f^2 \\ &\quad - 2cov(\phi_1, f) \mathbb{E}_D[\hat{w}_1 | \mathbf{X}] - 2cov(\phi_2, f) \mathbb{E}_D[\hat{w}_2 | \mathbf{X}]. \end{aligned}$$

Proof. In the one dimensional case, exploiting the independence between the train and test set:

$$\begin{aligned} \mathbb{E}_{X,Y,D}[(\bar{\mathcal{M}}(x) - \bar{y})^2 | \mathbf{X}] &= \mathbb{E}_{X,Y,D}[(\mathbb{E}_D[\hat{w}h(\phi(x))] - f(x))^2 | \mathbf{X}] \\ &= \mathbb{E}_X[(h(\phi(x)) \mathbb{E}_D[\hat{w} | \mathbf{X}] - f(x))^2] \\ &= \mathbb{E}_X[h(\phi(x))^2 \mathbb{E}_D[\hat{w} | \mathbf{X}]^2 + \mathbb{E}_X[f(x)^2] - 2\mathbb{E}_X[f(x)h(\phi(x))] \mathbb{E}_D[\hat{w} | \mathbf{X}]] \\ &= \frac{\sigma_h^2}{\hat{\sigma}_h^4} \hat{cov}(f, h)^2 + \sigma_f^2 - 2cov(f, h) \frac{\hat{cov}(f, h)}{\hat{\sigma}_h^2}, \end{aligned}$$

where the last equation follows substituting the expected value of the coefficient with its expression derived in the previous lemma and exploiting the assumption of null expected value of the aggregation function h .

In the two dimensional setting, exploiting again the independence between train and test set and the null assumption of the expected value of the two functions of features ϕ_1, ϕ_2 :

$$\begin{aligned} \mathbb{E}_{X,Y,D}[(\bar{\mathcal{M}}(x) - \bar{y})^2 | \mathbf{X}] &= \mathbb{E}_{X,Y}[(\phi_1(x) \mathbb{E}_D[\hat{w}_1] + \phi_2(x) \mathbb{E}_D[\hat{w}_2] - f(x))^2 | \mathbf{X}] \\ &= \sigma_{\phi_1}^2 \mathbb{E}_D[\hat{w}_1 | \mathbf{X}]^2 + \sigma_{\phi_2}^2 \mathbb{E}_D[\hat{w}_2 | \mathbf{X}]^2 + 2cov(\phi_1, \phi_2) \mathbb{E}_D[\hat{w}_1 | \mathbf{X}] \mathbb{E}_D[\hat{w}_2 | \mathbf{X}] \\ &\quad + \sigma_f^2 - 2\mathbb{E}[f(x)(\phi_1 \mathbb{E}_D[\hat{w}_1] + \phi_2 \mathbb{E}_D[\hat{w}_2]) | \mathbf{X}] \\ &= \sigma_{\phi_1}^2 \mathbb{E}_D[\hat{w}_1 | \mathbf{X}]^2 + \sigma_{\phi_2}^2 \mathbb{E}_D[\hat{w}_2 | \mathbf{X}]^2 + 2cov(\phi_1, \phi_2) \mathbb{E}_D[\hat{w}_1 | \mathbf{X}] \mathbb{E}_D[\hat{w}_2 | \mathbf{X}] \\ &\quad + \sigma_f^2 - 2cov(\phi_1, f) \mathbb{E}_D[\hat{w}_1 | \mathbf{X}] - 2cov(\phi_2, f) \mathbb{E}_D[\hat{w}_2 | \mathbf{X}]. \end{aligned}$$

□

Remark 10. *In the two dimensional result of Theorem 7 the expected values of the coefficients, found in Lemma 2, are not explicitly inserted, to improve readability. The full expression of squared bias, in the bivariate case, would be:*

$$\begin{aligned} &\mathbb{E}[(\bar{\mathcal{M}}(x) - \bar{y})^2 | \mathbf{X}] \\ &= \sigma_{\phi_1}^2 \left(\frac{\hat{\sigma}_{\phi_2}^2 \hat{cov}(\phi_1, f) - \hat{cov}(\phi_1, \phi_2) \hat{cov}(\phi_2, f)}{(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - \hat{cov}(\phi_1, \phi_2)^2)} \right)^2 \\ &\quad + \sigma_{\phi_2}^2 \left(\frac{\hat{\sigma}_{\phi_1}^2 \hat{cov}(\phi_2, f) - \hat{cov}(\phi_1, \phi_2) \hat{cov}(\phi_1, f)}{(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - \hat{cov}(\phi_1, \phi_2)^2)} \right)^2 \\ &\quad + 2cov(\phi_1, \phi_2) \end{aligned}$$

$$\begin{aligned}
& \times \left(\frac{(\hat{\sigma}_{\phi_2}^2 \hat{c}ov(\phi_1, f) - \hat{c}ov(\phi_1, \phi_2) \hat{c}ov(\phi_2, f))(\hat{\sigma}_{\phi_1}^2 \hat{c}ov(\phi_2, f) - \hat{c}ov(\phi_1, \phi_2) \hat{c}ov(\phi_1, f))}{(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - \hat{c}ov(\phi_1, \phi_2)^2)^2} \right) \\
& + \sigma_f^2 \\
& - 2\hat{c}ov(\phi_1, f) \left(\frac{\hat{\sigma}_{\phi_2}^2 \hat{c}ov(\phi_1, f) - \hat{c}ov(\phi_1, \phi_2) \hat{c}ov(\phi_2, f)}{(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - \hat{c}ov(\phi_1, \phi_2)^2)} \right) \\
& - 2\hat{c}ov(\phi_2, f) \left(\frac{\hat{\sigma}_{\phi_1}^2 \hat{c}ov(\phi_2, f) - \hat{c}ov(\phi_1, \phi_2) \hat{c}ov(\phi_1, f)}{(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - \hat{c}ov(\phi_1, \phi_2)^2)} \right).
\end{aligned}$$

That, after algebraic computations, is equal to:

$$\begin{aligned}
& \mathbb{E}[(\bar{\mathcal{M}}(x) - \bar{y})^2 | \mathbf{X}] \\
& = \sigma_f^2 + \frac{1}{(\hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 - \hat{c}ov(\phi_1, \phi_2)^2)^2} \times \left\{ \begin{aligned}
& \sigma_{\phi_1}^2 \hat{\sigma}_{\phi_2}^4 \hat{c}ov(\phi_1, f)^2 + \sigma_{\phi_1}^2 \hat{c}ov(\phi_1, \phi_2)^2 \hat{c}ov(\phi_2, f)^2 \\
& - 2\sigma_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 \hat{c}ov(\phi_1, f) \hat{c}ov(\phi_2, f) \hat{c}ov(\phi_1, \phi_2) \\
& + \sigma_{\phi_2}^2 \hat{\sigma}_{\phi_1}^4 \hat{c}ov(\phi_2, f)^2 + \sigma_{\phi_2}^2 \hat{c}ov(\phi_1, \phi_2)^2 \hat{c}ov(\phi_1, f)^2 \\
& - 2\sigma_{\phi_2}^2 \hat{\sigma}_{\phi_1}^2 \hat{c}ov(\phi_1, f) \hat{c}ov(\phi_2, f) \hat{c}ov(\phi_1, \phi_2) \\
& + 2\hat{c}ov(\phi_1, \phi_2) \hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 \hat{c}ov(\phi_1, f) \hat{c}ov(\phi_2, f) \\
& + 2\hat{c}ov(\phi_1, \phi_2) \hat{c}ov(\phi_1, f) \hat{c}ov(\phi_2, f) \hat{c}ov(\phi_1, \phi_2)^2 \\
& - 2\hat{c}ov(\phi_1, \phi_2) \hat{\sigma}_{\phi_2}^2 \hat{c}ov(\phi_1, f)^2 \hat{c}ov(\phi_1, \phi_2) \\
& - 2\hat{c}ov(\phi_1, \phi_2) \hat{\sigma}_{\phi_1}^2 \hat{c}ov(\phi_2, f)^2 \hat{c}ov(\phi_1, \phi_2) \\
& - 2\hat{c}ov(\phi_1, f) \hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^4 \hat{c}ov(\phi_1, f) \\
& + 2\hat{c}ov(\phi_1, f) \hat{\sigma}_{\phi_2}^2 \hat{c}ov(\phi_1, f) \hat{c}ov(\phi_1, \phi_2)^2 \\
& + 2\hat{c}ov(\phi_1, f) \hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 \hat{c}ov(\phi_2, f) \hat{c}ov(\phi_1, \phi_2) \\
& - 2\hat{c}ov(\phi_1, f) \hat{c}ov(\phi_2, f) \hat{c}ov(\phi_1, \phi_2)^3 \\
& - 2\hat{c}ov(\phi_2, f) \hat{\sigma}_{\phi_2}^2 \hat{\sigma}_{\phi_1}^4 \hat{c}ov(\phi_2, f) \\
& + 2\hat{c}ov(\phi_2, f) \hat{\sigma}_{\phi_1}^2 \hat{c}ov(\phi_2, f) \hat{c}ov(\phi_1, \phi_2)^2 \\
& + 2\hat{c}ov(\phi_2, f) \hat{\sigma}_{\phi_1}^2 \hat{\sigma}_{\phi_2}^2 \hat{c}ov(\phi_1, f) \hat{c}ov(\phi_1, \phi_2) \\
& - 2\hat{c}ov(\phi_2, f) \hat{c}ov(\phi_1, f) \hat{c}ov(\phi_1, \phi_2)^3
\end{aligned} \right\}.
\end{aligned}$$

Lemma 3. In the asymptotic case, considering each estimator convergent in probability to the quantity that it estimates, the (squared) bias for the univariate and bivariate linear regression under analysis is respectively:

$$\begin{aligned}
bias_{1D}^{n \rightarrow \infty} &= \sigma_f^2 - \frac{cov(f, h)^2}{\sigma_h^2}, \\
bias_{2D}^{n \rightarrow \infty} &= \sigma_f^2 + \frac{2cov(\phi_1, f)cov(\phi_2, f)cov(\phi_1, \phi_2) - \sigma_{\phi_1}^2 cov(\phi_2, f)^2 - \sigma_{\phi_2}^2 cov(\phi_1, f)^2}{\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - cov(\phi_1, \phi_2)^2}.
\end{aligned}$$

Proof. To prove the two results it is enough to start from the results of Theorem 7 and Remark 10 and substitute each estimator with the quantity that it estimates, to which it converges in probability. \square

Theorem 2, that is reported in the main paper, is finally proved from the two asymptotic quantities derived in Lemma 3, subtracting the univariate to the bivariate result.

B Generalized-Linear Correlated Features Aggregation: bivariate analysis

In this section we show a first bivariate result related to generalized linear models. As in the general case, we assume the conditional distribution of the target given the feature to belong to the canonical exponential family:

$$f_{\theta}(y) = \exp\left(\frac{y\theta - b(\theta)}{\phi}\right) + c(y, \phi).$$

Moreover, we assume that the expected value of the target is a linear combination of the inputs, subsequently transformed with the canonical link function:

$$\mathbb{E}[y|x] = g^{-1}(w_1x_1 + w_2x_2).$$

In this setting we compare a bivariate model with the univariate one that substitutes the two features with their aggregation through a zero-mean function $h(\cdot)$:

$$\begin{cases} \hat{y} = g^{-1}(\hat{w}_1x_1 + \hat{w}_2x_2) \\ \hat{y} = g^{-1}(\hat{w} \cdot \frac{x_1+x_2}{2}). \end{cases}$$

As a first result, we prove the following lemma that justifies the adoption of the expected deviance as a goodness of fit measure.

Lemma 4. *Assuming a bivariate generalized linear model and considering a Gaussian distribution of the target given the features $y \sim \mathcal{N}(\mu, \sigma^2) \implies f_{\theta}(y) = \exp\{\frac{y\theta - \theta^2/2}{\phi^2} + c(y, \phi)\}$, $\theta = \mu$, $\phi = \sigma$ the difference of performance in terms of MSE between the bivariate linear regression and the univariate one that aggregates the two input features with their mean is equivalent to the difference of deviance between the two models.*

Proof. Recalling the definition of deviance:

$$D^*(\theta, \hat{\theta}) := \frac{D(\theta, \hat{\theta})}{\phi} = -2 \left[\ell(\hat{\theta}) - \ell(\theta) \right] = \frac{2}{\phi} [y(\theta - \hat{\theta}) - (b(\theta) - b(\hat{\theta}))],$$

in the Gaussian setting the parameter $\theta = \mu = w_1x_1 + w_2x_2$ represents the mean of the distribution of the target y . Defining $\hat{\theta} = \hat{\mu}_1$ the mean of the target estimated with the univariate regression and $\bar{\theta} = \hat{\mu}_2$ the one of the bivariate case, the increase of expected deviance between the two models is:

$$\mathbb{E}_{x,y,\mathcal{T}}[D^*(\theta, \hat{\theta}) - D^*(\theta, \bar{\theta})] = \mathbb{E}_{x,y,\mathcal{T}}[D^*(\mu, \hat{\mu}_1) - D^*(\mu, \hat{\mu}_2)].$$

Moreover, assuming Gaussianity we have $b(\theta) = \frac{\theta^2}{2}$ and the link function $g(\cdot)$ is the identity function (indeed, in the bivariate linear regression we just compare the prediction $\hat{y} = \hat{w}_1x_1 + \hat{w}_2x_2$ with the univariate $\hat{y} = \hat{w} \frac{x_1+x_2}{2}$).

Therefore, the expected increase of deviance becomes:

$$\begin{aligned} \mathbb{E}_{x,y,\mathcal{T}}[D^*(\mu, \hat{\mu}_1) - D^*(\mu, \hat{\mu}_2)] &= \frac{2}{\phi} \mathbb{E}_{x,y,\mathcal{T}}[y(\hat{\mu}_2 - \hat{\mu}_1) - (b(\hat{\mu}_2) - b(\hat{\mu}_1))] \\ &= \frac{2}{\sigma^2} \mathbb{E}_{x,y,\mathcal{T}} \left[y(\hat{w}_1x_1 + \hat{w}_2x_2 - \hat{w} \frac{x_1+x_2}{2}) - \frac{(\hat{w}_1x_1 + \hat{w}_2x_2)^2}{2} - \frac{(\hat{w}\bar{x})^2}{2} \right] \\ &= \frac{2}{\sigma^2} \left\{ \mathbb{E}_{x,y} \left[y \cdot \left(w_1x_1 + w_2x_2 - \frac{2(w_1\hat{\sigma}_{x_1}^2 + w_2\hat{\sigma}_{x_2}^2 + (w_1+w_2)c\hat{ov}(x_1, x_2))}{\hat{\sigma}_{x_1}^2 + \hat{\sigma}_{x_2}^2 + 2c\hat{ov}(x_1, x_2)} \cdot \frac{x_1+x_2}{2} \right) \right] \right. \\ &\quad \left. - \mathbb{E}_{x,y} \mathbb{E}_{\mathcal{T}} \left[\frac{(\hat{w}_1x_1 + \hat{w}_2x_2)^2}{2} - \frac{(\hat{w}\bar{x})^2}{2} \right] \right\}. \end{aligned}$$

The last equation follows from the independence of train and test set and the expression of the expected value of the regression coefficients (Equation A2,A3 of Bonetti et al. (2023)).

Asymptotically, the first expected value is:

$$\begin{aligned}
& \mathbb{E}_{x,y} \left[y \cdot \left(w_1 x_1 + w_2 x_2 - \frac{2(w_1 \sigma_{x_1}^2 + w_2 \sigma_{x_2}^2 + (w_1 + w_2) \text{cov}(x_1, x_2))}{\sigma_{x_1}^2 + \sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2)} \cdot \frac{x_1 + x_2}{2} \right) \right] \\
&= \mathbb{E}_{x,y} \left[y \cdot \left(w_1 x_1 + w_2 x_2 - \frac{(w_1 \sigma_{x_1}^2 + w_2 \sigma_{x_2}^2 + (w_1 + w_2) \text{cov}(x_1, x_2))}{\sigma_{x_1}^2 + \sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2)} \cdot (x_1 + x_2) \right) \right] \\
&= w_1 \mathbb{E}_{x,y} [x_1 y] + w_2 \mathbb{E}_{x,y} [x_2 y] \\
&\quad - \frac{(w_1 \sigma_{x_1}^2 + w_2 \sigma_{x_2}^2 + (w_1 + w_2) \text{cov}(x_1, x_2))}{\sigma_{x_1}^2 + \sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2)} \cdot (\mathbb{E}_{x,y} [x_1 y] + \mathbb{E}_{x,y} [x_2 y]).
\end{aligned}$$

Recalling the definition of covariance and the zero-mean assumption of the expected values:

$$\begin{aligned}
\mathbb{E}_{x,y} [x_1 y] &= \text{cov}(x_1, y) - \mathbb{E}_x [x_1] \mathbb{E}_y [y] \\
&= \text{cov}(x_1, w_1 x_1 + w_2 x_2 + \varepsilon) - 0 = w_1 \sigma_{x_1}^2 + w_2 \text{cov}(x_1, x_2).
\end{aligned}$$

A similar argumentation holds for $\mathbb{E}_{x,y} [x_2 y]$, therefore the expected value under analysis becomes:

$$\begin{aligned}
& w_1^2 \sigma_{x_1}^2 + w_2^2 \sigma_{x_2}^2 + 2w_1 w_2 \text{cov}(x_1, x_2) \\
& \quad - \frac{(w_1 \sigma_{x_1}^2 + w_2 \sigma_{x_2}^2 + (w_1 + w_2) \text{cov}(x_1, x_2))^2}{\sigma_{x_1}^2 + \sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2)}.
\end{aligned}$$

The second expected value that appears in the equation representing the increase of deviance is asymptotically equal to:

$$\begin{aligned}
& \mathbb{E}_{x,y} \mathbb{E}_{\mathcal{T}} \left[\frac{(\hat{w}_1 x_1 + \hat{w}_2 x_2)^2}{2} - \frac{(\hat{w} \bar{x})^2}{2} \right] = \frac{1}{2} \mathbb{E}_{x,y} \mathbb{E}_{\mathcal{T}} [\hat{w}_1^2 x_1^2 + \hat{w}_2^2 x_2^2 + 2w_1 w_2 x_1 x_2 - \hat{w}^2 \bar{x}^2] \\
&= \mathbb{E}_{\mathcal{T}} [\hat{w}_1^2] \mathbb{E}_x [x_1^2] + \mathbb{E}_{\mathcal{T}} [\hat{w}_2^2] \mathbb{E}_x [x_2^2] + 2 \mathbb{E}_{\mathcal{T}} [\hat{w}_1 \hat{w}_2] \mathbb{E}_x [x_1 x_2] - \mathbb{E}_{\mathcal{T}} [\hat{w}^2] \mathbb{E}_x [\bar{x}^2] \\
&= \sigma_{x_1}^2 \left(w_1^2 + \frac{\sigma^2 \sigma_{x_2}^2}{(n-1)(\sigma_{x_1}^2 \sigma_{x_2}^2 - \text{cov}^2(x_1, x_2))} \right) \\
&+ \sigma_{x_2}^2 \left(w_2^2 + \frac{\sigma^2 \sigma_{x_1}^2}{(n-1)(\sigma_{x_1}^2 \sigma_{x_2}^2 - \text{cov}^2(x_2, x_2))} \right) \\
&+ 2 \text{cov}(x_1, x_2) \left(w_1 w_2 - \frac{\sigma^2 \text{cov}(x_1, x_2)}{(n-1)(\sigma_{x_1}^2 \sigma_{x_2}^2 - \text{cov}^2(x_1, x_2))} \right) \\
&- \frac{1}{4} \left[(\sigma_{x_1}^2 + \sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2)) \right. \\
&\quad \times \left(\frac{\sigma^2}{(n-1) \cdot \frac{1}{4} (\sigma_{x_1}^2 + \sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2))} \right. \\
&\quad \left. \left. + \left(\frac{2(w_1 \sigma_{x_1}^2 + w_2 \sigma_{x_2}^2 + (w_1 + w_2) \text{cov}(x_1, x_2))}{\sigma_{x_1}^2 + \sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2)} \right)^2 \right) \right] \\
&= w_1^2 \sigma_{x_1}^2 + w_2^2 \sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2) w_1 w_2 + \frac{2\sigma^2 (\sigma_{x_1}^2 \sigma_{x_2}^2 - \text{cov}^2(x_1, x_2))}{(n-1)(\sigma_{x_1}^2 \sigma_{x_2}^2 - \text{cov}^2(x_1, x_2))}
\end{aligned}$$

$$\begin{aligned}
& - \left[\frac{\sigma^2}{n-1} + \frac{(w_1\sigma_{x_1}^2 + w_2\sigma_{x_2}^2 + (w_1 + w_2) \cdot \text{cov}(x_1, x_2))^2}{\sigma_{x_1}^2 + \sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2)} \right] \\
& = w_1^2\sigma_{x_1}^2 + w_2^2\sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2) w_1 w_2 + \frac{\sigma^2}{n-1} \\
& - \frac{(w_1\sigma_{x_1}^2 + w_2\sigma_{x_2}^2 + (w_1 + w_2) \text{cov}(x_1, x_2))^2}{\sigma_{x_1}^2 + \sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2)},
\end{aligned}$$

where the last two equalities follow again from the expression of expected value and variance of the regression coefficients.

Combining the two asymptotic terms, the expected increase of deviance becomes:

$$\begin{aligned}
\mathbb{E}_{x,y,\mathcal{T}}[D^*(\mu, \hat{\mu}_1) - D^*(\mu, \hat{\mu}_2)] &= \frac{2}{\sigma^2} \left(w_1^2\sigma_{x_1}^2 + w_2^2\sigma_{x_2}^2 + 2w_1w_2 \text{cov}(x_1, x_2) \right. \\
& \quad \left. - \frac{(w_1\sigma_{x_1}^2 + w_2\sigma_{x_2}^2 + (w_1 + w_2) \text{cov}(x_1, x_2))^2}{\sigma_{x_1}^2 + \sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2)} \right) \\
& - \frac{1}{\sigma^2} \left(w_1^2\sigma_{x_1}^2 + w_2^2\sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2) w_1 w_2 + \frac{\sigma^2}{n-1} \right. \\
& \quad \left. - \frac{(w_1\sigma_{x_1}^2 + w_2\sigma_{x_2}^2 + (w_1 + w_2) \text{cov}(x_1, x_2))^2}{\sigma_{x_1}^2 + \sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2)} \right) \\
& = \frac{1}{\sigma^2} \left(w_1^2\sigma_{x_1}^2 + w_2^2\sigma_{x_2}^2 + 2w_1w_2 \text{cov}(x_1, x_2) \right. \\
& \quad \left. - \frac{(w_1\sigma_{x_1}^2 + w_2\sigma_{x_2}^2 + (w_1 + w_2) \text{cov}(x_1, x_2))^2}{\sigma_{x_1}^2 + \sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2)} - \frac{\sigma^2}{n-1} \right) \\
& = \frac{1}{\sigma^2(\sigma_{x_1}^2 + \sigma_{x_2}^2 + 2 \text{cov}(x_1, x_2))} \left(w_1^2\sigma_{x_1}^2\sigma_{x_2}^2 + w_2^2\sigma_{x_1}^2\sigma_{x_2}^2 \right. \\
& \quad \left. - w_1^2 \text{cov}^2(x_1, x_2) - w_2^2 \text{cov}^2(x_1, x_2) - 2w_1w_2\sigma_{x_1}^2\sigma_{x_2}^2 + 2w_1w_2 \text{cov}^2(x_1, x_2) \right) \\
& - \frac{1}{\sigma^2} \cdot \frac{\sigma^2}{n-1}.
\end{aligned}$$

Recalling that the asymptotic increase of bias in the linear bivariate setting is:

$$\frac{\sigma_{x_1}^2 \sigma_{x_2}^2 (1 - \rho_{x_1, x_2}^2) (w_1 - w_2)^2}{\sigma_{x_1 + x_2}^2},$$

and the asymptotic decrease of variance is $\frac{\sigma^2}{n-1}$ (respectively Equation 9 and 13 of Bonetti et al. (2023)), the first is equal to the first term found comparing the deviance and the second one is equal to the second one, proving the equivalence. \square

After having verified that, in linear contexts, the analysis of deviance is tight w.r.t. the mean squared error, we are now ready to introduce the main result of this section.

Theorem 8. *Let M be the maximum between the following differences: $\mathbb{E}_{\mathcal{T}}[\hat{w}_1] - \frac{1}{2}\mathbb{E}_{\mathcal{T}}[\hat{w}]$, $\mathbb{E}_{\mathcal{T}}[\hat{w}_2] - \frac{1}{2}\mathbb{E}_{\mathcal{T}}[\hat{w}]$, $\mathbb{E}_{\mathcal{T}}[\hat{w}_1^2] - \frac{1}{4}\mathbb{E}_{\mathcal{T}}[\hat{w}^2]$, $\mathbb{E}_{\mathcal{T}}[\hat{w}_2^2] - \frac{1}{4}\mathbb{E}_{\mathcal{T}}[\hat{w}^2]$, $\mathbb{E}_{\mathcal{T}}[\hat{w}_1\hat{w}_2] - \frac{1}{4}\mathbb{E}_{\mathcal{T}}[\hat{w}^2]$, and let m be the minimum of the same quantities. Moreover, let the real model belong to the canonical exponential family.*

Defining $\Delta(D^*)$ as the expected increase of deviance due to the aggregation of the two features x_1, x_2 with their mean in the bivariate setting, it is equal to:

$$\Delta(D^*) \leq \frac{2}{\phi} \left\{ M \cdot (|\text{cov}(x_1, y)| + |\text{cov}(x_2, y)|) - \frac{1}{2} m \cdot b''(0) \cdot (\sigma_{x_1+x_2}^2) \right\}. \quad (17)$$

Proof. In the bivariate setting, the expected increase of deviance under analysis is equal to:

$$\begin{aligned} & \frac{2}{\phi} \{ E_{x,y} [y \cdot (\mathbb{E}_{\mathcal{T}} [\hat{w}_1] x_1 + \mathbb{E}_{\mathcal{T}} [\hat{w}_2] x_2 - \mathbb{E}_{\mathcal{T}} [\hat{w}] \bar{x})] \\ & - (\mathbb{E}_{x,y} \mathbb{E}_{\mathcal{T}} [b(\hat{w}_1 x_1 + \hat{w}_2 x_2) - b(\hat{w} \bar{x})]) \}. \end{aligned}$$

First expected value The first expected value of the previous equation can be rewritten as:

$$\begin{aligned} & E_{x,y} [y \cdot (\mathbb{E}_{\mathcal{T}} [\hat{w}_1] x_1 + \mathbb{E}_{\mathcal{T}} [\hat{w}_2] x_2 - \mathbb{E}_{\mathcal{T}} [\hat{w}] \bar{x})] \\ &= E_{x,y} [y x_1] \mathbb{E}_{\mathcal{T}} [\hat{w}_1] + E_{x,y} [y x_2] \mathbb{E}_{\mathcal{T}} [\hat{w}_2] - E_{x,y} \left[y \frac{x_1 + x_2}{2} \right] \mathbb{E}_{\mathcal{T}} [\hat{w}] \\ &= \text{cov}(x_1, y) \mathbb{E}_{\mathcal{T}} [\hat{w}_1] + \text{cov}(x_2, y) \mathbb{E}_{\mathcal{T}} [\hat{w}_2] - \frac{1}{2} \text{cov}(x_1, y) \mathbb{E}_{\mathcal{T}} [\hat{w}] - \frac{1}{2} \text{cov}(x_2, y) \mathbb{E}_{\mathcal{T}} [\hat{w}] \\ &= \text{cov}(x_1, y) \left(\mathbb{E}_{\mathcal{T}} [\hat{w}_1] - \frac{1}{2} \mathbb{E}_{\mathcal{T}} [\hat{w}] \right) + \text{cov}(x_2, y) \left(\mathbb{E}_{\mathcal{T}} [\hat{w}_2] - \frac{1}{2} \mathbb{E}_{\mathcal{T}} [\hat{w}] \right). \end{aligned}$$

Second expected value The second expected value that appears in the expected increase of deviance under analysis, $\mathbb{E}_{x,y} \mathbb{E}_{\mathcal{T}} [b(\hat{w}_1 x_1 + \hat{w}_2 x_2) - b(\hat{w} \bar{x})]$, depends on the function $b(\cdot)$, that is a specific parameter of the distribution. In order to derive a result that holds for any distribution belonging to the canonical exponential family, we use a second order Taylor approximation of the function, centered in $\theta_0 = 0$.

$$\begin{aligned} & \begin{cases} b(\hat{w}_1 x_1 + \hat{w}_2 x_2) \simeq b(0) + b'(0) \cdot (\hat{w}_1 x_1 + \hat{w}_2 x_2) + \frac{1}{2} b''(0) \cdot (\hat{w}_1 x_1 + \hat{w}_2 x_2)^2 \\ b(\hat{w} \bar{x}) \simeq b(0) + b'(0) \cdot (\hat{w} \bar{x}) + \frac{1}{2} b''(0) \cdot (\hat{w} \bar{x})^2. \end{cases} \end{aligned}$$

Since the first-order terms vanish because of the zero-mean assumption, the approximated expected value is therefore:

$$\begin{aligned} & \mathbb{E}_{x,y,\mathcal{T}} [b(\hat{w}_1 x_1 + \hat{w}_2 x_2) - b(\hat{w} \bar{x})] \simeq \frac{1}{2} b''(0) \cdot \mathbb{E}_{x,y,\mathcal{T}} \left[((\hat{w}_1 x_1 + \hat{w}_2 x_2)^2 - (\hat{w} \bar{x})^2) \right] \\ &= \frac{1}{2} b''(0) \cdot \left[\sigma_{x_1}^2 \cdot \mathbb{E}_{\mathcal{T}} [\hat{w}_1^2 - \frac{\hat{w}^2}{4}] + \sigma_{x_2}^2 \cdot \mathbb{E}_{\mathcal{T}} [\hat{w}_2^2 - \frac{\hat{w}^2}{4}] \right. \\ & \quad \left. + 2 \text{cov}(x_1, x_2) \cdot \mathbb{E}_{\mathcal{T}} [\hat{w}_1 \hat{w}_2 - \frac{\hat{w}^2}{4}] \right], \\ & \iff \frac{1}{2} b''(0) \cdot [\text{var}(\hat{w}_1 x_1 + \hat{w}_2 x_2) - \text{var}(\hat{w} \bar{x})]. \end{aligned}$$

Expected increase of deviance We are now ready to combine the two expressions:

$$\begin{aligned} \Delta(D^*) &= \frac{2}{\phi} \left\{ E_{x,y} [y \cdot (\mathbb{E}_{\mathcal{T}} [\hat{w}_1] x_1 + \mathbb{E}_{\mathcal{T}} [\hat{w}_2] x_2 - \mathbb{E}_{\mathcal{T}} [\hat{w}] \bar{x})] \right. \\ & \quad \left. - (\mathbb{E}_{x,y} \mathbb{E}_{\mathcal{T}} [b(\hat{w}_1 x_1 + \hat{w}_2 x_2) - b(\hat{w} \bar{x})]) \right\} \\ &= \frac{2}{\phi} \left\{ \text{cov}(x_1, y) \left(\mathbb{E}_{\mathcal{T}} [\hat{w}_1] - \frac{1}{2} \mathbb{E}_{\mathcal{T}} [\hat{w}] \right) + \text{cov}(x_2, y) \left(\mathbb{E}_{\mathcal{T}} [\hat{w}_2] - \frac{1}{2} \mathbb{E}_{\mathcal{T}} [\hat{w}] \right) \right. \\ & \quad \left. - \frac{1}{2} b''(0) \left[\sigma_{x_1}^2 \left(\mathbb{E}_{\mathcal{T}} [\hat{w}_1^2] - \frac{1}{4} \mathbb{E}_{\mathcal{T}} [\hat{w}^2] \right) + \sigma_{x_2}^2 \left(\mathbb{E}_{\mathcal{T}} [\hat{w}_2^2] - \frac{1}{4} \mathbb{E}_{\mathcal{T}} [\hat{w}^2] \right) \right. \right. \end{aligned}$$

$$+ 2 \text{cov}(x_1, x_2) \left(\mathbb{E}_{\mathcal{T}}[\hat{w}_1 \hat{w}_2] - \frac{1}{4} \mathbb{E}_{\mathcal{T}}[\hat{w}^2] \right) \Bigg\}.$$

Finally, substituting with: M the maximum difference of the expected value of the coefficients that appear in the equation $(\mathbb{E}_{\mathcal{T}}[\hat{w}_1] - \frac{1}{2}\mathbb{E}_{\mathcal{T}}[\hat{w}], \mathbb{E}_{\mathcal{T}}[\hat{w}_2] - \frac{1}{2}\mathbb{E}_{\mathcal{T}}[\hat{w}], \mathbb{E}_{\mathcal{T}}[\hat{w}_1^2] - \frac{1}{4}\mathbb{E}_{\mathcal{T}}[\hat{w}^2], \mathbb{E}_{\mathcal{T}}[\hat{w}_2^2] - \frac{1}{4}\mathbb{E}_{\mathcal{T}}[\hat{w}^2], \mathbb{E}_{\mathcal{T}}[\hat{w}_1 \hat{w}_2] - \frac{1}{4}\mathbb{E}_{\mathcal{T}}[\hat{w}^2])$, and with m the minimum of the same quantities, the result follows. \square

We conclude this bivariate analysis with a justification of the choice to center in 0 the Taylor expansion in the proof above and providing an intuitive interpretation of the result of the theorem.

Remark 11. *Considering the center in zero ($\theta_0 = 0$) for the second-order Taylor expansion of the function $b(\cdot)$ in the analysis of deviance, in the linear asymptotic case, the second order Taylor expansion is an exact approximation of the term.*

Proof. In the linear case, recalling that $b(\theta) = \frac{\theta^2}{2}$, the expected value that contains the function $b(\cdot)$ in the proof, is:

$$\begin{aligned} \mathbb{E}[b(\hat{w}_1 x_1 + \hat{w}_2 x_2) - b(\hat{w} \bar{x})] &= \frac{1}{2} \mathbb{E}[(\hat{w}_1 x_1 + \hat{w}_2 x_2)^2 - (\hat{w} \bar{x})^2] \\ &= \frac{1}{2} \mathbb{E} \left[\hat{w}_1^2 x_1^2 + \hat{w}_2^2 x_2^2 + 2\hat{w}_1 \hat{w}_2 x_1 x_2 - \frac{\hat{w}^2 x_1^2}{4} - \frac{\hat{w}^2 x_2^2}{4} - \frac{\hat{w}^2 x_1 x_2}{2} \right] \\ &= \frac{1}{2} \left[\sigma_{x_1}^2 \cdot \mathbb{E}_{\mathcal{T}} \left[\hat{w}_1^2 - \frac{\hat{w}^2}{4} \right] + \sigma_{x_2}^2 \cdot \mathbb{E}_{\mathcal{T}} \left[\hat{w}_2^2 - \frac{\hat{w}^2}{4} \right] \right. \\ &\quad \left. + 2 \text{cov}(x_1, x_2) \cdot \mathbb{E}_{\mathcal{T}} \left[\hat{w}_1 \hat{w}_2 - \frac{\hat{w}^2}{4} \right] \right]. \end{aligned}$$

\square

Moreover, $b''(0) = 1$. Therefore, this quantity is equal to the general expression of the second order Taylor expansion centered in 0.

Remark 12. *The result of the theorem suggests that it is convenient to aggregate two variables when the variance of the sum between the two variables is large or the absolute value of the covariance between each feature and the target is small. Indeed, this implies respectively that there is much noise or that each feature shares a lot of information with the target individually.*

C Generalized-Linear Correlated Features Aggregation: additional proofs and results

Proof of Theorem 4

Recalling the expression of the expected increase of deviance:

$$\begin{aligned} &\frac{2}{\phi} \left\{ \mathbb{E}_{x,y} [y \cdot (\mathbb{E}_{\mathcal{T}}[\hat{w}_1] \phi_1 + \mathbb{E}_{\mathcal{T}}[\hat{w}_2] \phi_2 - \mathbb{E}_{\mathcal{T}}[\hat{w}] h(\phi_1, \phi_2))] \right. \\ &\quad \left. - \mathbb{E}_{x,y,\mathcal{T}} [b(\hat{w}_1 \phi_1 + \hat{w}_2 \phi_2) - b(\hat{w} h(\phi_1, \phi_2))] \right\}, \end{aligned} \tag{18}$$

we analyse the two expected values separately.

Exploiting the zero-mean assumption, the first expected value is equal to:

$$\begin{aligned} &E_{x,y} [y \cdot (\mathbb{E}_{\mathcal{T}}[\hat{w}_1] \phi_1 + \mathbb{E}_{\mathcal{T}}[\hat{w}_2] \phi_2 - \mathbb{E}_{\mathcal{T}}[\hat{w}] h(\phi_1, \phi_2))] \\ &= E_{x,y} [y \phi_1] \mathbb{E}_{\mathcal{T}}[\hat{w}_1] + E_{x,y} [y \phi_2] \mathbb{E}_{\mathcal{T}}[\hat{w}_2] - E_{x,y} [y h(\phi_1, \phi_2)] \mathbb{E}_{\mathcal{T}}[\hat{w}] \end{aligned}$$

$$= \text{cov}(\phi_1, y) \mathbb{E}_{\mathcal{T}} [\hat{w}_1] + \text{cov}(\phi_2, y) \mathbb{E}_{\mathcal{T}} [\hat{w}_2] - \text{cov}(h(\phi_1, \phi_2), y) \mathbb{E}_{\mathcal{T}} [\hat{w}].$$

The second expectation, on the other hand, is:

$$\mathbb{E}_{x,y,\mathcal{T}} [b(\hat{w}_1\phi_1 + \hat{w}_2\phi_2) - b(\hat{w}h(\phi_1, \phi_2))].$$

In order to obtain an expression that holds for any distribution in the canonical exponential family, similarly to the bivariante case of the previous section, we perform a second order Taylor approximation of the function $b(\cdot)$:

$$\begin{cases} b(\hat{w}_1\phi_1 + \hat{w}_2\phi_2) \simeq b(0) + b'(0) \cdot (\hat{w}_1\phi_1 + \hat{w}_2\phi_2) + \frac{1}{2}b''(0) \cdot (\hat{w}_1\phi_1 + \hat{w}_2\phi_2)^2 \\ b(\hat{w}h(\phi_1, \phi_2)) \simeq b(0) + b'(0) \cdot (\hat{w}h(\phi_1, \phi_2)) + \frac{1}{2}b''(0) \cdot (\hat{w}h(\phi_1, \phi_2))^2. \end{cases}$$

The expected value under analysis therefore becomes:

$$\begin{aligned} & \mathbb{E}_{x,y,\mathcal{T}} [b(\hat{w}_1\phi_1 + \hat{w}_2\phi_2) - b(\hat{w}h(\phi_1, \phi_2))] \\ & \simeq \frac{1}{2}b''(0) \cdot \mathbb{E}_{x,y,\mathcal{T}} [(\hat{w}_1\phi_1 + \hat{w}_2\phi_2)^2 - (\hat{w}h(\phi_1, \phi_2))^2] \\ & = \frac{1}{2}b''(0) \cdot [\sigma_{\phi_1}^2 \mathbb{E}_{\mathcal{T}} [\hat{w}_1^2] + \sigma_{\phi_2}^2 \mathbb{E}_{\mathcal{T}} [\hat{w}_2^2] \\ & \quad + 2 \text{cov}(\phi_1, \phi_2) \mathbb{E}_{\mathcal{T}} [\hat{w}_1\hat{w}_2] - \sigma_{h(\phi_1, \phi_2)}^2 \mathbb{E}_{\mathcal{T}} [\hat{w}^2]] \\ & \iff \frac{1}{2}b''(0) \cdot [\text{var}(\hat{w}_1\phi_1 + \hat{w}_2\phi_2) - \text{var}(\hat{w}h(\phi_1, \phi_2))]. \end{aligned}$$

Merging the two expected values, the expected increase of deviance finally becomes:

$$\begin{aligned} \Delta(D^*) &= \frac{2}{\phi} \left\{ \text{cov}(\phi_1, y) \mathbb{E}_{\mathcal{T}} [\hat{w}_1] + \text{cov}(\phi_2, y) \mathbb{E}_{\mathcal{T}} [\hat{w}_2] - \text{cov}(h(\phi_1, \phi_2), y) \mathbb{E}_{\mathcal{T}} [\hat{w}] \right. \\ & \quad \left. - \frac{1}{2}b''(0) \cdot [\sigma_{\phi_1}^2 \mathbb{E}_{\mathcal{T}} [\hat{w}_1^2] + \sigma_{\phi_2}^2 \mathbb{E}_{\mathcal{T}} [\hat{w}_2^2] \right. \\ & \quad \left. + 2 \text{cov}(\phi_1, \phi_2) \mathbb{E}_{\mathcal{T}} [\hat{w}_1\hat{w}_2] - \sigma_{h(\phi_1, \phi_2)}^2 \mathbb{E}_{\mathcal{T}} [\hat{w}^2]] \right\} \\ & \leq \frac{2}{\phi} \left\{ \left[M \cdot (|\text{cov}(\phi_1, y)| + |\text{cov}(\phi_2, y)|) - m \cdot |\text{cov}(h(\phi), y)| \right] \right. \\ & \quad \left. - \frac{1}{2}b''(0) \cdot \left(m \cdot \sigma_{x_1+x_2}^2 - M \cdot \sigma_{h(\phi)}^2 \right) \right\}, \end{aligned}$$

where M, m are respectively the maximum and minimum absolute values of the expected values of the coefficients and their squared values. This concludes the proof, since the second expression of the theorem follows by rearranging the terms.

C.1 Additional Gaussian considerations

In this section we prove the statement of Remark 9, assuming that the generalized linear model under analysis follows a Gaussian distribution.

Lemma 5. *Let a generalized linear model have Gaussian distribution of the target given the features ($Y|X \sim \mathbb{N}(\mu, \sigma^2)$) and compare the bivariate linear regression having as inputs two zero-mean transformations of the features $\phi_1(x), \phi_2(x)$ with a univariate linear regression having as input the aggregation of ϕ_1, ϕ_2 through a function $h(\cdot)$ (Equation 7 of the main paper). The increase of expected MSE due to the aggregation is equal to the increase of the expected deviance.*

Proof. Recalling that, with the Gaussianity assumption, the link function $g(\cdot)$ is the identity and the function $b(\cdot) = \frac{\theta^2}{2}$, the expression of the expected increase of deviance (Equation 18) becomes:

$$\frac{2}{\sigma^2} \left\{ \mathbb{E}_{x,y} [y \cdot (\mathbb{E}_{\mathcal{T}} [\hat{w}_1]\phi_1 + \mathbb{E}_{\mathcal{T}} [\hat{w}_2]\phi_2 - \mathbb{E}_{\mathcal{T}} [\hat{w}]h(\phi_1, \phi_2))] \right\}$$

$$\begin{aligned}
& -\frac{1}{2}\mathbb{E}_{x,y,\mathcal{T}} \left[(\hat{w}_1\phi_1 + \hat{w}_2\phi_2)^2 - (\hat{w}h(\phi_1, \phi_2))^2 \right] \Big\} \\
& = \left\{ \mathbb{E}_{x,y} \left[y \cdot \left(\frac{\phi_1(\sigma_{\phi_2}^2 \text{cov}(\phi_1, f) - \text{cov}(\phi_1, \phi_2)\text{cov}(\phi_2, f)) + \phi_2(\sigma_{\phi_1}^2 \text{cov}(\phi_2, f) - \text{cov}(\phi_1, \phi_2)\text{cov}(\phi_1, f))}{\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2)} \right. \right. \right. \\
& \quad \left. \left. - \frac{\text{cov}(h(\phi_1, \phi_2), f)}{\sigma_h^2(\phi_1, \phi_2)} \cdot h(\phi_1, \phi_2) \right) \right] \\
& \quad \left. - \frac{1}{2}\mathbb{E}_{x,y,\mathcal{T}} \left[\hat{w}_1^2\phi_1^2 + \hat{w}_2^2\phi_2^2 + 2\hat{w}_1\hat{w}_2\phi_1\phi_2 - \hat{w}^2\hat{h}^2(\phi_1, \phi_2) \right] \right\},
\end{aligned}$$

where the equation holds substituting the expression of the expected value of the linear regression coefficients.

First expected value.

Exploiting the independence between train and test data and the zero mean assumption, the first expected value in the equation above becomes:

$$\begin{aligned}
& \frac{\text{cov}(\phi_1, y) \cdot [\sigma_{\phi_2}^2 \text{cov}(\phi_1, f) - \text{cov}(\phi_1, \phi_2) \text{cov}(\phi_2, f)] + \text{cov}(\phi_2, y) \cdot [\sigma_{\phi_1}^2 \text{cov}(\phi_2, f) - \text{cov}(\phi_1, \phi_2) \text{cov}(\phi_1, f)]}{\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2)} \\
& - \frac{\text{cov}(h(\phi_1, \phi_2), f)}{\sigma_{h(\phi_1, \phi_2)}^2} \cdot \text{cov}(h(\phi_1, \phi_2), y) \\
& \frac{\sigma_{\phi_2}^2 \text{cov}^2(\phi_1, f) - \text{cov}(\phi_1, \phi_2) \text{cov}(\phi_1, f) \text{cov}(\phi_2, f) + \sigma_{\phi_1}^2 \text{cov}^2(\phi_2, f) - \text{cov}(\phi_1, \phi_2) \text{cov}(\phi_1, f) \text{cov}(\phi_2, f)}{\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2)} \\
& - \frac{\text{cov}^2(h(\phi_1, \phi_2), f)}{\sigma_{h(\phi_1, \phi_2)}^2} \\
& = \frac{\sigma_{\phi_2}^2 \text{cov}^2(\phi_1, f) + \sigma_{\phi_1}^2 \text{cov}^2(\phi_2, f) - 2 \text{cov}(\phi_1, \phi_2) \text{cov}(\phi_1, f) \text{cov}(\phi_2, f)}{\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2)} - \frac{\text{cov}^2(h(\phi_1, \phi_2), f)}{\sigma_{h(\phi_1, \phi_2)}^2} \\
& = \Delta_{\text{bias}}^{n \rightarrow \infty},
\end{aligned}$$

where the last expression is exactly the increase of (squared) asymptotic bias found in Equation 11.

Second expected value

Substituting the expected values of the estimates of the regression coefficients and their variances, the second expected value of the expression of the increase of deviance becomes:

$$\begin{aligned}
& \mathbb{E}_{\mathcal{T}} [\hat{w}_1^2] \sigma_{\phi_1}^2 + \mathbb{E}_{\mathcal{T}} [\hat{w}_2^2] \sigma_{\phi_2}^2 + 2\mathbb{E}_{\mathcal{T}} [\hat{w}_1 \hat{w}_2] \text{cov}(\phi_1, \phi_2) - \mathbb{E}_{\mathcal{T}} [\hat{w}^2] \sigma_{h(\phi_1, \phi_2)}^2 \\
& = \left(\text{var}(\hat{w}_1) + \mathbb{E}_{\mathcal{T}} [\hat{w}_1^2] \right) \sigma_{\phi_1}^2 + \left(\text{var}(\hat{w}_2) + \mathbb{E}_{\mathcal{T}} [\hat{w}_2^2] \right) \sigma_{\phi_2}^2 \\
& + 2(\text{cov}(\hat{w}_1, \hat{w}_2) + \mathbb{E}_{\mathcal{T}} [\hat{w}_1] \mathbb{E}_{\mathcal{T}} [\hat{w}_2]) \text{cov}(\phi_1, \phi_2) - (\text{var}(\hat{w}) + \mathbb{E}_{\mathcal{T}} [\hat{w}^2]) \sigma_{h(\phi_1, \phi_2)}^2 \\
& = \left[\frac{\sigma^2 \cdot \sigma_{\phi_2}^2}{(n-1)(\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2))} + \left(\frac{\sigma_{\phi_2}^2 \text{cov}(\phi_1, f) - \text{cov}(\phi_1, \phi_2) \text{cov}(\phi_2, f)}{\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2)} \right)^2 \right] \sigma_{\phi_1}^2 \\
& + \left[\frac{\sigma^2 \cdot \sigma_{\phi_1}^2}{(n-1)(\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2))} + \left(\frac{\sigma_{\phi_1}^2 \text{cov}(\phi_2, f) - \text{cov}(\phi_1, \phi_2) \text{cov}(\phi_1, f)}{\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2)} \right)^2 \right] \sigma_{\phi_2}^2 \\
& + 2 \text{cov}(\phi_1, \phi_2) \left[\frac{-\sigma^2 \cdot \text{cov}(\phi_1, \phi_2)}{(n-1)(\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2))} \right. \\
& \left. + \left(\frac{(\sigma_{\phi_2}^2 \text{cov}(\phi_1, f) - \text{cov}(\phi_1, \phi_2) \text{cov}(\phi_2, f)) \cdot (\sigma_{\phi_1}^2 \text{cov}(\phi_2, f) - \text{cov}(\phi_1, \phi_2) \text{cov}(\phi_1, f))}{(\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2))^2} \right) \right]
\end{aligned}$$

$$- \left(\frac{\sigma^2}{n-1} \cdot \frac{1}{\sigma_{h(\phi_1, \phi_2)}^2} + \left(\frac{\text{cov}(h(\phi_1, \phi_2), f)}{\sigma_{h(\phi_1, \phi_2)}^2} \right)^2 \right) \sigma_{h(\phi_1, \phi_2)}^2.$$

Considering the terms with $n-1$ in the denominator, they are equal to the asymptotic decrease of variance found in Equation 10:

$$\begin{aligned} \frac{\sigma^2}{n-1} \cdot \left[\frac{\sigma_{\phi_2}^2 \sigma_{\phi_1}^2 - \sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - 2 \text{cov}^2(\phi_1, \phi_2)}{\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2)} - \frac{\sigma_{h(\phi_1, \phi_2)}^2}{\sigma_{h(\phi_1, \phi_2)}^2} \right] \\ = \frac{\sigma^2}{n-1} \cdot (2-1) = \frac{\sigma^2}{n-1} = \Delta_{\text{var}}^{n \rightarrow \infty}. \end{aligned}$$

The remaining terms are:

$$\begin{aligned} \sigma_{\phi_1}^2 \left(\frac{\sigma_{\phi_2}^2 \text{cov}(\phi_1, f) - \text{cov}(\phi_1, \phi_2) \text{cov}(\phi_2, f)}{\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2)} \right)^2 + \sigma_{\phi_2}^2 \left(\frac{\sigma_{\phi_1}^2 \text{cov}(\phi_2, f) - \text{cov}(\phi_1, \phi_2) \text{cov}(\phi_1, f)}{\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2)} \right)^2 \\ + 2 \text{cov}(\phi_1, \phi_2) \left(\frac{\left(\sigma_{\phi_2}^2 \text{cov}(\phi_1, f) - \text{cov}(\phi_1, \phi_2) \text{cov}(\phi_2, f) \right) \cdot \left(\sigma_{\phi_1}^2 \text{cov}(\phi_2, f) - \text{cov}(\phi_1, \phi_2) \text{cov}(\phi_1, f) \right)}{\left(\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2) \right)^2} \right) \\ - \frac{\text{cov}^2(h(\phi_1, \phi_2), f)}{\sigma_{h(\phi_1, \phi_2)}^2}. \end{aligned}$$

The first three terms of this expression have the same denominator, therefore we can focus on their numerators:

$$\begin{aligned} \sigma_{\phi_1}^2 (\sigma_{\phi_2}^4 \text{cov}^2(\phi_1, f) + \text{cov}^2(\phi_1, \phi_2) \text{cov}^2(\phi_2, f) - 2 \sigma_{\phi_2}^2 \text{cov}(\phi_1, f) \text{cov}(\phi_2, f) \text{cov}(\phi_1, \phi_2)) \\ + \sigma_{\phi_2}^2 (\sigma_{\phi_1}^4 \text{cov}^2(\phi_2, f) + \text{cov}^2(\phi_1, \phi_2) \text{cov}^2(\phi_1, f) - 2 \sigma_{\phi_1}^2 \text{cov}(\phi_1, f) \text{cov}(\phi_2, f) \text{cov}(\phi_1, \phi_2)) \\ + 2 \text{cov}(\phi_1, \phi_2) (\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 \text{cov}(\phi_1, f) \text{cov}(\phi_2, f) - \sigma_{\phi_2}^2 \text{cov}^2(\phi_1, f) \text{cov}(\phi_1, \phi_2) \\ + \text{cov}^2(\phi_1, \phi_2) \text{cov}(\phi_1, f) \text{cov}(\phi_2, f) - \sigma_{\phi_1}^2 \text{cov}^2(\phi_2, f) \text{cov}(\phi_1, \phi_2)) \\ = \sigma_{\phi_1}^2 \sigma_{\phi_2}^4 \text{cov}^2(\phi_1, f) + \sigma_{\phi_1}^4 \sigma_{\phi_2}^2 \text{cov}^2(\phi_2, f) - \sigma_{\phi_1}^2 \text{cov}^2(\phi_2, f) \text{cov}^2(\phi_1, \phi_2) \\ - \sigma_{\phi_2}^2 \text{cov}^2(\phi_1, f) \text{cov}^2(\phi_1, \phi_2) - 2 \sigma_{\phi_1}^2 \sigma_{\phi_2}^2 \text{cov}(\phi_1, f) \text{cov}(\phi_2, f) \text{cov}(\phi_1, \phi_2) \\ + 2 \text{cov}^3(\phi_1, \phi_2) \text{cov}(\phi_1, f) \text{cov}(\phi_2, f) \\ = \sigma_{\phi_1}^2 \text{cov}^2(\phi_2, f) \cdot [\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2)] + \sigma_{\phi_2}^2 \text{cov}^2(\phi_1, f) \cdot [\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2)] \\ - 2 \text{cov}(\phi_1, f) \text{cov}(\phi_2, f) \text{cov}(\phi_1, \phi_2) \cdot [\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2)] \\ = (\sigma_{\phi_1}^2 \text{cov}^2(\phi_2, f) + \sigma_{\phi_2}^2 \text{cov}^2(\phi_1, f) - 2 \text{cov}(\phi_1, f) \text{cov}(\phi_2, f) \text{cov}(\phi_1, \phi_2)) \cdot (\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2)). \end{aligned}$$

The full term is therefore:

$$\begin{aligned} \frac{(\sigma_{\phi_1}^2 \text{cov}^2(\phi_2, f) + \sigma_{\phi_2}^2 \text{cov}^2(\phi_1, f) - 2 \text{cov}(\phi_1, f) \text{cov}(\phi_2, f) \text{cov}(\phi_1, \phi_2))}{(\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2))^2} \\ \times (\sigma_{\phi_1}^2 \sigma_{\phi_2}^2 - \text{cov}^2(\phi_1, \phi_2)) - \frac{\text{cov}^2(h(\phi_1, \phi_2), f)}{\sigma_{h(\phi_1, \phi_2)}^2} = \Delta_{\text{bias}}^{n \rightarrow \infty}. \end{aligned}$$

Summing up, the expected deviance is therefore equal to:

$$\frac{2}{\phi} \left(\Delta_{\text{bias}}^{n \rightarrow \infty} - \frac{1}{2} (\Delta_{\text{var}}^{n \rightarrow \infty} + \Delta_{\text{bias}}^{n \rightarrow \infty}) \right) = \frac{1}{\phi} (\Delta_{\text{bias}}^{n \rightarrow \infty} - \Delta_{\text{var}}^{n \rightarrow \infty}),$$

that is exactly the increase of MSE in terms of increase of bias and reduction of variance due to the aggregation found in the analysis of the NonLinCFA algorithm. \square

A second result assuming Gaussianity justifies the choice of centering the second order Taylor expansion in $\theta_0 = 0$ to approximate the function $b(\cdot)$ for any generalized linear model.

Lemma 6. *Considering the center in zero ($\theta_0 = 0$), in the Gaussian asymptotic case, the second order Taylor expansion leads to an exact approximation of the function $b(\cdot)$.*

Proof. The function $b(\cdot)$ appears in the increase of deviance analysis performed in the previous subsection in the following equation:

$$\mathbb{E}_{x,y,\tau} [b(\hat{w}_1\phi_1 + \hat{w}_2\phi_2) - b(\hat{w}h(\phi_1, \phi_2))].$$

In the linear case, recalling that $b(\theta) = \frac{\theta^2}{2}$, the expected value under analysis is equal to:

$$\begin{aligned} & \frac{1}{2} \cdot \mathbb{E}_{x,y,D} \left[(\hat{w}_1\phi_1 + \hat{w}_2\phi_2)^2 - (\hat{w}h(\phi_1, \phi_2))^2 \right] \\ &= \frac{1}{2} \cdot \left[\sigma_{\phi_1}^2 \mathbb{E}_D[\hat{w}_1^2] + \sigma_{\phi_2}^2 \mathbb{E}_D[\hat{w}_2^2] + 2 \text{cov}(\phi_1, \phi_2) \mathbb{E}_D[\hat{w}_1\hat{w}_2] - \sigma_{h(\phi_1, \phi_2)}^2 \mathbb{E}_D[\hat{w}^2] \right]. \end{aligned}$$

Moreover, $b''(0) = 1$. Therefore, this quantity is equal to the general expression of the second order Taylor expansion centered in 0 (Equation C). \square

D Experiments

D.1 Synthetic experiments

This subsection provides more details on the synthetic experiments introduced in the main paper. The first synthetic problem that we designed is composed of $D = 100$ features and standard deviation of the noise $\sigma = 10$. The first independent variable x_1 follows a uniform distribution in the interval $[0, 1]$. Any other feature x_i , $i \in \{2, \dots, 100\}$, is a linear combination between a randomly chosen previous features x_j , $j < i$ and a random variable that follows a uniform distribution in the interval $[0, 1]$ (specifically $x_i = 0.7x_j + 0.3u$, $u \sim \mathcal{U}([0, 1])$). The target variable y is finally a linear combination between the D features x_1, \dots, x_{100} , randomly sampling the coefficients from a uniform distribution in $[0, 1]$. Moreover, a Gaussian noise with standard deviation σ is added to the target. The same experiment is repeated in a more complex setting, considering $D = 1000$ features and $\sigma = 100$ as standard deviation of the additive noise.

To test the GenLinCFA algorithm in classification settings, the same two experiment described in this section have been repeated applying the sign function to the target, in order to transform the two problems into classification tasks.

In Table 3 and Table 4 it is possible to find the detailed results of the experiments performed respectively with continuous and binary target. As discussed in the main paper, also LinCFA algorithm has been applied to compare the regression results. Since it has no hyperparameters to tune, it is not reported in the tables. In the setting with $D = 100$ features it leads to $d = 39 \pm 1.52$ features and an $R2score = 0.8659 \pm 0.0049$. In the setting with $D = 1000$ features, it produces $d = 193.2 \pm 2.37$ features with an $R2score = 0.7070 \pm 0.0074$.

The two regression and the two classification synthetic experiments that have been described in this section have been repeated considering a nonlinear relationship between the features and the target. In particular, the datasets have been generated exactly in the same way, with the only difference to define the target

Table 3: 95% confidence intervals, linear synthetic setting: different hyperparameters, continuous target.

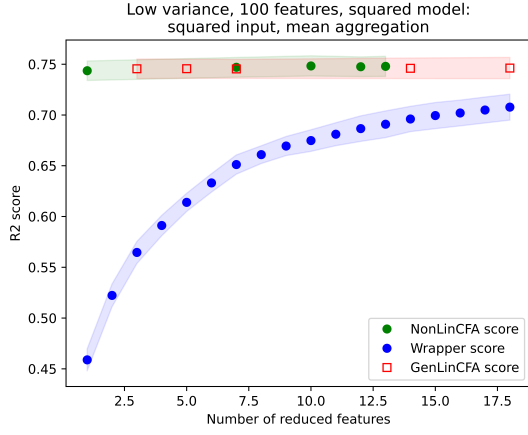
Noise variance $\sigma = 10$, number of features $D = 100$					
NonLinCFA					
Hyperparameter	$\epsilon = 0.01$	$\epsilon = 0.001$	$\epsilon = 0.0001$	$\epsilon = 0.00001$	$\epsilon = 0.000001$
Number reduced features d	1.0 ± 0.0	8.0 ± 0.88	11.4 ± 1.45	14.4 ± 1.28	14.7 ± 1.21
R^2 test score	0.8655 ± 0.0051	0.8664 ± 0.0048	0.8661 ± 0.0049	0.8659 ± 0.0050	0.8664 ± 0.0043
GenLinCFA					
Hyperparameter	$\epsilon = 0.76$	$\epsilon = 0.77$	$\epsilon = 0.78$	$\epsilon = 0.79$	$\epsilon = 0.80$
Number reduced features d	21.6 ± 1.47	16.6 ± 0.63	13.6 ± 1.31	5.4 ± 1.15	2.0 ± 0.39
R^2 test score	0.8656 ± 0.0047	0.8663 ± 0.0046	0.8660 ± 0.0046	0.8659 ± 0.0049	0.8657 ± 0.0049
Noise variance $\sigma = 100$, number of features $D = 1000$					
NonLinCFA					
Hyperparameter	$\epsilon = 0.01$	$\epsilon = 0.001$	$\epsilon = 0.0001$	$\epsilon = 0.00001$	$\epsilon = 0.000001$
Number reduced features d	1.0 ± 0.0	3.1 ± 0.65	18.0 ± 1.92	21.9 ± 2.10	22.3 ± 1.90
R^2 test score	0.7332 ± 0.0069	0.7318 ± 0.0071	0.7274 ± 0.0079	0.7265 ± 0.0072	0.7267 ± 0.0076
GenLinCFA					
Hyperparameter	$\epsilon = 0.76$	$\epsilon = 0.77$	$\epsilon = 0.78$	$\epsilon = 0.79$	$\epsilon = 0.80$
Number reduced features d	7.3 ± 1.08	3.4 ± 0.63	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0
R^2 test score	0.7326 ± 0.0069	0.7325 ± 0.0072	0.7332 ± 0.0069	0.7332 ± 0.0069	0.7332 ± 0.0069

Table 4: 95% confidence intervals, linear synthetic setting: different hyperparameters, binary target.

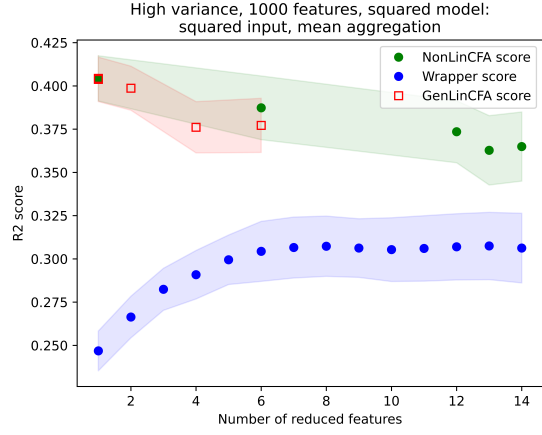
Noise variance $\sigma = 10$, number of features $D = 100$					
GenLinCFA					
Hyperparameter	$\epsilon = 0.71$	$\epsilon = 0.72$	$\epsilon = 0.73$	$\epsilon = 0.75$	$\epsilon = 0.77$
Number reduced features d	25.2 ± 1.59	19.4 ± 1.69	15.6 ± 1.39	4.3 ± 1.21	1.0 ± 0.0
$Accuracy$ test score	0.8928 ± 0.0064	0.8947 ± 0.0066	0.8956 ± 0.0065	0.8958 ± 0.0069	0.8975 ± 0.0054
Noise variance $\sigma = 100$, number of features $D = 1000$					
GenLinCFA					
Hyperparameter	$\epsilon = 0.71$	$\epsilon = 0.72$	$\epsilon = 0.73$	$\epsilon = 0.75$	$\epsilon = 0.77$
Number reduced features d	20.0 ± 3.54	11.1 ± 2.06	5.7 ± 0.88	1.0 ± 0.0	1.0 ± 0.0
$Accuracy$ test score	0.8462 ± 0.0067	0.8453 ± 0.0075	0.8429 ± 0.0064	0.8520 ± 0.0048	0.8520 ± 0.0048

variable as a linear combination of the squared value of the input features, with additive Gaussian noise. In this setting, a wrapper feature selection has been again considered as baseline, considering the squared of the inputs as candidate features to select. NonLinCFA and GenLinCFA have been applied in two different ways: considering as features the squared values of the original features ($\phi_i(x) = x_i^2$) and selecting the mean as aggregation function, or considering the original features as inputs ($\phi_i(x) = x_i$) and performing the squared sum as aggregation function. In this way, the two methods have been tested considering nonlinear transformations of features or nonlinear aggregation functions.

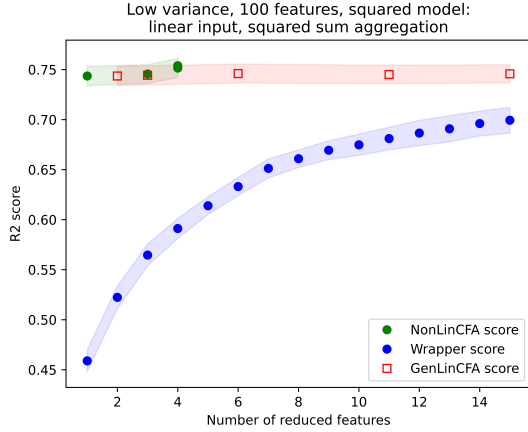
Figure 2 shows the results in terms of coefficient of determination for the two regression problems with low and high number of features and variance. In all the four cases the hyperparameter ϵ of NonLinCFA has been considered in the range $\{0.01, 0.001, 0.0001, 1e-05, 1e-06\}$. Moreover, for GenLinCFA it has been considered in the range $\{0.68, 0.7, 0.72, 0.74, 0.76\}$. All experiments have been repeated ten times to produce confidence intervals. Figure 2a and Figure 2c show the test regression performance considering respectively squared input features and the mean and aggregation function and linear input features and the squared sum as aggregation function, in the low dimensional and low variance scenario. The same results obtained with the same hyperparameters are reported in Figure 2b and Figure 2d for the high dimensional with high variance regression setting. As already discussed in the main paper for linear settings, the linear regression applied on the features aggregated by the two algorithms perform better than the linear regression performed on the features selected by the wrapper approach, that needs more features to obtain similar performances. Moreover, it is possible to observe that GenLinCFA is more prone to aggregate features in high dimensional and noisy problems. On the contrary, NonLinCFA, is less prone to aggregate the features in high dimensional and noisy problems. Finally, considering linear inputs and squared sum aggregation or



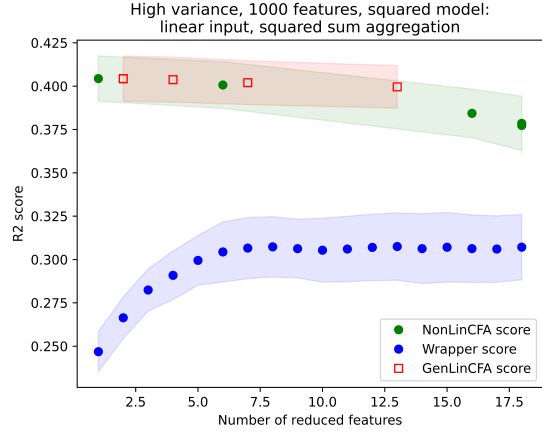
(a) Results considering standard deviation $\sigma = 10$, $D = 100$ features, squared inputs and mean aggregations.



(b) Results considering standard deviation $\sigma = 100$, $D = 1000$ features, squared inputs and mean aggregations.



(c) Results considering standard deviation $\sigma = 10$, $D = 100$ features, linear inputs and sum of square aggregations.

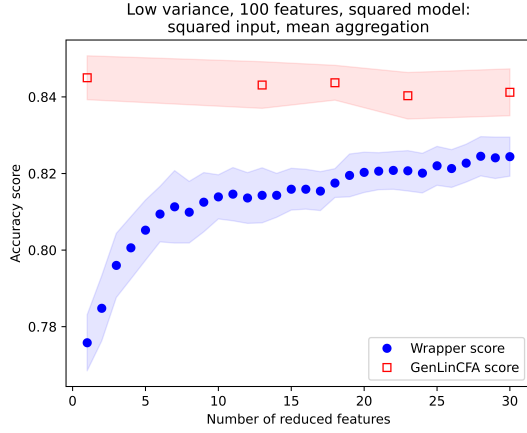


(d) Results considering standard deviation $\sigma = 100$, $D = 1000$ features, linear inputs and sum of square aggregations.

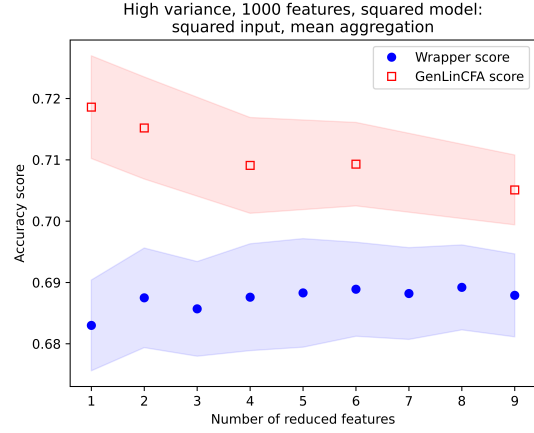
Figure 2: Application of the two proposed algorithm. Test performances and number of reduced features considering different hyperparameters, transformations of features and aggregation functions in a quadratic setting.

squared inputs and linear aggregations the two algorithms have similar performances, showing that they can be applied in both contexts. As a further comparison, LinCFA algorithm has also been applied, considering the square of the features as inputs. The linear regression on the reduced features have similar performances in the two settings, with a test score respectively of 0.7503 ± 0.0080 and 0.3619 ± 0.0186 for the low and high dimensional problems, but a larger number of reduced features ($d = 34.3 \pm 0.91$ and $d = 120.5 \pm 1.69$).

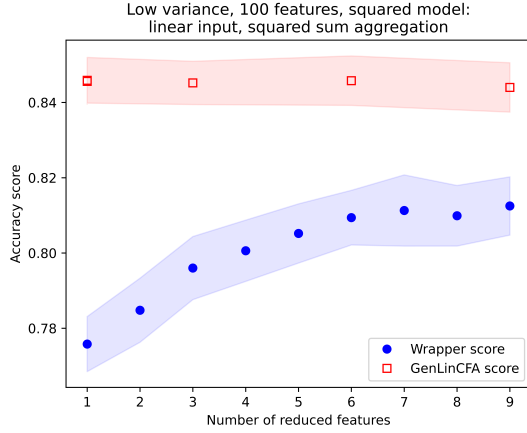
To further analyse GenLinCFA algorithm in classification contexts, the two experiments that consider a quadratic relationship between the features and the target have been repeated in classification, applying the sign function to the target. Again, mean aggregations of the squared features and squared sums of the features have been considered. The accuracy scores and number of reduced features are reported in Figure 3. The figure compares again the performance of the algorithm with a wrapper feature selection, logistic regression is the supervised model applied to produce the scores. Again, the performances of the proposed algorithm are similar or better w.r.t. the wrapper baseline and, considering the same values of hiperparameters in the two settings, the algorithm is more prone to aggregate in the more complex and noisy setting. The detailed results and confidence intervals for different hyperparameters that have been discussed in this setting can be found in Table 5 and Table 6, respectively for regression and classification.



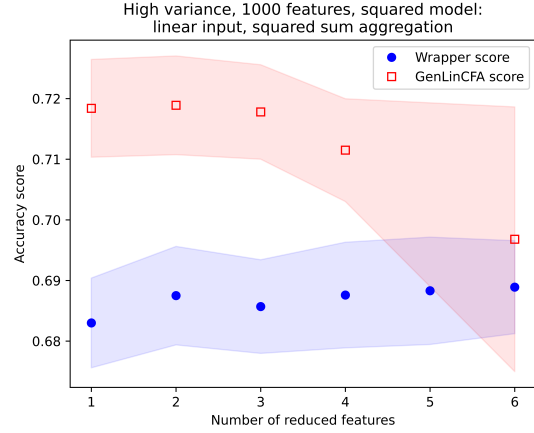
(a) Classification results with standard deviation $\sigma = 10$, $D = 100$ features, squared inputs and mean aggregations.



(b) Classification results with standard deviation $\sigma = 100$, $D = 1000$ features, squared inputs, mean aggregations.



(c) Classification results with standard deviation $\sigma = 10$, $D = 100$ features, linear inputs, sum of square aggregations.



(d) Classification results with standard deviation $\sigma = 100$, $D = 1000$ features, linear inputs, sum of square aggregations.

Figure 3: Application of GenLinCFA in classification. Test performances and number of reduced features considering different hyperparameters, transformations of features and aggregation functions.

D.1.1 Real-World experiments

Table 7 and Table 8 report the confidence intervals associated to five different repetitions of the experiments, in terms of coefficient of determination for regression and accuracy for classification. The *Finance* dataset has been retrieved from Kaggle (<https://www.kaggle.com/datasets/dgawlik/nyse>). The *Bankruptcy* and *Parkinson* datasets have been retrieved from the UCI ML repository (<https://archive.ics.uci.edu/ml/datasets/Parkinson%27s+Disease+Classification>, <https://archive.ics.uci.edu/ml/datasets/Polish+companies+bankruptcy+data>).

The climate datasets are composed of continuous climatological features and a scalar target which represents the state of vegetation of a basin of Po river. The first one (*Climate*, *Climate(Class.)*) also considers the state of the vegetation of neighbouring basins as inputs, while the second one is a more difficult problem since it tries to predict it only from temperature and precipitation features. These datasets have been composed by the authors merging different sources for the vegetation index, temperature and precipitation over different

Table 5: 95% confidence intervals for quadratic synthetic experiments: different hyperparameters, continuous target.

Noise variance $\sigma = 10$, number of features $D = 100$					
NonLinCFA, squared input, mean aggregation					
Hyperparameter	$\epsilon = 0.01$	$\epsilon = 0.001$	$\epsilon = 0.0001$	$\epsilon = 0.00001$	$\epsilon = 0.000001$
Number reduced features d	1.0 ± 0.0	7.1 ± 0.85	10.3 ± 1.64	12.9 ± 1.58	12.4 ± 1.39
R^2 test score	0.7436 ± 0.0096	0.7469 ± 0.0097	0.7483 ± 0.0099	0.7479 ± 0.0099	0.7475 ± 0.0096
GenLinCFA, squared input, mean aggregation					
Hyperparameter	$\epsilon = 0.68$	$\epsilon = 0.70$	$\epsilon = 0.72$	$\epsilon = 0.74$	$\epsilon = 0.76$
Number reduced features d	17.6 ± 0.85	14.5 ± 1.64	7.2 ± 1.58	4.9 ± 1.39	3.0 ± 1.39
R^2 test score	0.7462 ± 0.0104	0.7459 ± 0.0101	0.7453 ± 0.0093	0.7455 ± 0.0094	0.7455 ± 0.0096
NonLinCFA, linear input, square sum aggregation					
Hyperparameter	$\epsilon = 0.01$	$\epsilon = 0.001$	$\epsilon = 0.0001$	$\epsilon = 0.00001$	$\epsilon = 0.000001$
Number reduced features d	1.0 ± 0.0	2.8 ± 0.46	3.8 ± 1.37	4.1 ± 1.43	4.1 ± 1.43
R^2 test score	0.7457 ± 0.0090	0.7450 ± 0.0091	0.7460 ± 0.0094	0.7442 ± 0.0096	0.7435 ± 0.0096
GenLinCFA, linear input, square sum aggregation					
Hyperparameter	$\epsilon = 0.68$	$\epsilon = 0.70$	$\epsilon = 0.72$	$\epsilon = 0.74$	$\epsilon = 0.76$
Number reduced features d	15.4 ± 0.46	10.9 ± 0.37	6.3 ± 0.41	3.3 ± 0.43	2.1 ± 0.43
R^2 test score	0.7462 ± 0.0104	0.7459 ± 0.0101	0.7453 ± 0.0093	0.7455 ± 0.0094	0.7455 ± 0.0096
Noise variance $\sigma = 100$, number of features $D = 1000$					
NonLinCFA, squared input, mean aggregation					
Hyperparameter	$\epsilon = 0.01$	$\epsilon = 0.001$	$\epsilon = 0.0001$	$\epsilon = 0.00001$	$\epsilon = 0.000001$
Number reduced features d	1.0 ± 0.0	5.9 ± 0.75	11.6 ± 0.69	13.2 ± 0.61	13.7 ± 0.62
R^2 test score	0.4044 ± 0.0130	0.3873 ± 0.0184	0.3735 ± 0.0179	0.3627 ± 0.0201	0.3649 ± 0.0199
GenLinCFA, squared input, mean aggregation					
Hyperparameter	$\epsilon = 0.68$	$\epsilon = 0.70$	$\epsilon = 0.72$	$\epsilon = 0.74$	$\epsilon = 0.76$
Number reduced features d	5.5 ± 0.76	3.5 ± 0.69	2.0 ± 0.60	1.2 ± 0.63	1.1 ± 0.69
R^2 test score	0.3772 ± 0.0156	0.3760 ± 0.0148	0.3986 ± 0.0127	0.4037 ± 0.0128	0.4043 ± 0.0130
NonLinCFA, linear input, square sum aggregation					
Hyperparameter	$\epsilon = 0.01$	$\epsilon = 0.001$	$\epsilon = 0.0001$	$\epsilon = 0.00001$	$\epsilon = 0.000001$
Number reduced features d	1.0 ± 0.0	6.3 ± 0.78	15.5 ± 1.55	17.7 ± 1.92	18.0 ± 1.84
R^2 test score	0.4043 ± 0.0130	0.4006 ± 0.0134	0.3843 ± 0.0140	0.3785 ± 0.0156	0.3772 ± 0.0166
GenLinCFA, linear input, square sum aggregation					
Hyperparameter	$\epsilon = 0.68$	$\epsilon = 0.70$	$\epsilon = 0.72$	$\epsilon = 0.74$	$\epsilon = 0.76$
Number reduced features d	12.7 ± 0.78	7.2 ± 1.55	4.1 ± 1.92	2.4 ± 1.83	1.7 ± 1.61
R^2 test score	0.3995 ± 0.0122	0.4020 ± 0.0126	0.4037 ± 0.0129	0.4043 ± 0.0130	0.4042 ± 0.0130

basins (see (Didan, 2015; Cornes et al., 2018; Zellner, 2022)), and they are available in the repository of this work.

The two tables show the performances of the NonLinCFA algorithm, considering five different hyperparameters ($\epsilon \in \{0.01, 0.001, 0.0001, 0.00001, 0.000001\}$). The same holds for the GenLinCFA algorithm, where the hyperparameters have been selected in order to show some aggregations that are not too small (everything is aggregated) or is too large (more than 50 reduced features). The state of the art methods applied have all been considered identifying the best validation performance between 1 and 50 reduced features, when possible. The results, discussed also in the main paper, show better or competitive results with respect to the baselines.

Table 6: 95% confidence intervals for quadratic synthetic experiments: different hyperparameters, binary target.

Noise variance $\sigma = 10$, number of features $D = 100$					
GenLinCFA, squared input, mean aggregation					
Hyperparameter	$\epsilon = 0.77$	$\epsilon = 0.79$	$\epsilon = 0.81$	$\epsilon = 0.85$	$\epsilon = 0.95$
Number reduced features d	30.2 ± 1.97	22.9 ± 2.01	18.5 ± 0.84	13.0 ± 1.17	1.0 ± 0.0
<i>Accuracy</i> test score	0.8412 ± 0.0061	0.8403 ± 0.0060	0.8436 ± 0.0045	0.8431 ± 0.0061	0.8450 ± 0.0057
GenLinCFA, linear input, square sum aggregation					
Hyperparameter	$\epsilon = 0.77$	$\epsilon = 0.79$	$\epsilon = 0.81$	$\epsilon = 0.85$	$\epsilon = 0.95$
Number reduced features d	8.8 ± 1.38	5.5 ± 1.57	3.2 ± 1.03	1.4 ± 0.41	1.0 ± 0.0
<i>Accuracy</i> test score	0.8440 ± 0.0065	0.8457 ± 0.0066	0.8452 ± 0.0057	0.8459 ± 0.0061	0.8455 ± 0.0060
Noise variance $\sigma = 100$, number of features $D = 1000$					
GenLinCFA, squared input, mean aggregation					
Hyperparameter	$\epsilon = 0.77$	$\epsilon = 0.79$	$\epsilon = 0.81$	$\epsilon = 0.85$	$\epsilon = 0.95$
Number reduced features d	9.4 ± 1.15	5.9 ± 0.65	3.5 ± 0.50	1.6 ± 0.41	1.0 ± 0.0
<i>Accuracy</i> test score	0.7051 ± 0.0056	0.7093 ± 0.0068	0.7091 ± 0.0078	0.7152 ± 0.0083	0.7186 ± 0.0083
GenLinCFA, linear input, square sum aggregation					
Hyperparameter	$\epsilon = 0.77$	$\epsilon = 0.79$	$\epsilon = 0.81$	$\epsilon = 0.85$	$\epsilon = 0.95$
Number reduced features d	6.1 ± 2.22	3.7 ± 1.33	2.6 ± 0.63	1.6 ± 0.30	1.3 ± 0.28
<i>Accuracy</i> test score	0.6968 ± 0.0218	0.7114 ± 0.0084	0.7178 ± 0.0077	0.7189 ± 0.0081	0.7184 ± 0.0080

Table 7: Experiments on climate datasets. The total number of samples n has been divided into train (66% of data) and test (33% of data) sets.

Quantity	Climate	Climate (Class.)	Climate II	Climate II (class.)
# samples n	981	981	867	867
# features D	1991	1991	2408	2408
Reduced Dimension				
NonLinCFA ($\epsilon = 0.01$)	7.4 ± 1.3	NA	7.0 ± 0.9	NA
NonLinCFA ($\epsilon = 0.001$)	16.0 ± 0.8	NA	11.4 ± 0.9	NA
NonLinCFA ($\epsilon = 0.0001$)	19.4 ± 1.4	NA	12.2 ± 0.6	NA
NonLinCFA ($\epsilon = 0.00001$)	19.6 ± 0.9	NA	12.4 ± 0.7	NA
NonLinCFA ($\epsilon = 0.000001$)	19.6 ± 1.3	NA	12.4 ± 0.7	NA
GenLinCFA ($\epsilon = \epsilon_1$)	31.6 ± 10.9	33.75 ± 11.5	13.6 ± 4.8	34.0 ± 6.1
GenLinCFA ($\epsilon = \epsilon_2$)	26.0 ± 9.1	27.0 ± 8.2	7.2 ± 1.1	26.0 ± 6.1
GenLinCFA ($\epsilon = \epsilon_3$)	18.6 ± 5.4	22.3 ± 5.7	4.6 ± 0.9	18.7 ± 5.4
GenLinCFA ($\epsilon = \epsilon_4$)	15.0 ± 4.3	21.3 ± 5.2	2.4 ± 0.4	11.0 ± 4.5
GenLinCFA ($\epsilon = \epsilon_5$)	13.8 ± 3.0	17.5 ± 3.3	2 ± 0	2.5 ± 0.4
LinCFA	38.2 ± 1.6	NA	222.0 ± 2.7	NA
PCA	18.2 ± 0.7	18.2 ± 0.7	29.4 ± 0.1	29.4 ± 0.4
LDA	NA	1	NA	1
Kernel PCA	35.6 ± 7.9	27.0 ± 4.8	21.8 ± 9.5	11.2 ± 2.3
Isomap	2.6 ± 0.7	12.6 ± 10.8	42.4 ± 13.3	16.6 ± 9.2
LLE	15.4 ± 13.8	28.2 ± 12.9	35.0 ± 10.7	25.8 ± 13.1
Supervised PCA	41.8 ± 2.5	37.0 ± 6.6	7.4 ± 3.1	13.8 ± 4.9
NCA	NA	24.2 ± 2.9	NA	7.0 ± 0.6
Test performance	R² score	Accuracy score	R² score	Accuracy score
NonLinCFA ($\epsilon = 0.01$)	0.8524 ± 0.0407	NA	0.2949 ± 0.0156	NA
NonLinCFA ($\epsilon = 0.001$)	0.9395 ± 0.0125	NA	0.2547 ± 0.0182	NA
NonLinCFA ($\epsilon = 0.0001$)	0.9124 ± 0.0055	NA	0.2541 ± 0.0121	NA
NonLinCFA ($\epsilon = 0.00001$)	0.9113 ± 0.0056	NA	0.2529 ± 0.0148	NA
NonLinCFA ($\epsilon = 0.000001$)	0.9121 ± 0.0047	NA	0.2530 ± 0.0146	NA
GenLinCFA ($\epsilon = \epsilon_1$)	0.9194 ± 0.0039	0.9068 ± 0.0037	0.2439 ± 0.0246	0.6820 ± 0.0199
GenLinCFA ($\epsilon = \epsilon_2$)	0.9226 ± 0.0026	0.9075 ± 0.0029	0.2841 ± 0.0051	0.7127 ± 0.0159
GenLinCFA ($\epsilon = \epsilon_3$)	0.9207 ± 0.0052	0.9056 ± 0.0039	0.2764 ± 0.0069	0.7094 ± 0.0181
GenLinCFA ($\epsilon = \epsilon_4$)	0.9269 ± 0.0012	0.9062 ± 0.0036	0.2493 ± 0.0125	0.7061 ± 0.0105
GenLinCFA ($\epsilon = \epsilon_5$)	0.9275 ± 0.0004	0.9107 ± 0.0022	0.2269 ± 0.0157	0.6776 ± 0.0159
LinCFA	0.9007 ± 0.031	NA	-1.2861 ± 0.2322	NA
PCA	0.7536 ± 0.019	0.8515 ± 0.0054	0.1917 ± 0.0395	0.6868 ± 0.0147
LDA	NA	0.7357 ± 0.0188	NA	0.5526 ± 0.0224
Kernel PCA	0.7990 ± 0.0061	0.8698 ± 0.0081	0.3889 ± 0.0199	0.7640 ± 0.0062
Isomap	0.1354 ± 0.0118	0.6443 ± 0.0041	0.3216 ± 0.0146	0.7360 ± 0.0145
LLE	0.1149 ± 0.0139	0.6444 ± 0.0052	0.3102 ± 0.0367	0.7456 ± 0.0140
Supervised PCA	0.8454 ± 0.0049	0.8827 ± 0.0098	0.3835 ± 0.0230	0.7482 ± 0.0067
NCA	NA	0.8776 ± 0.0086	NA	0.7638 ± 0.0051

Table 8: Experiments on real world datasets. The total number of samples n has been divided into train (66% of data) and test (33% of data) sets.

Quantity	Finance	Bankruptcy (class.)	Parkinson (class.)
# samples n	1299	1084	384
# features D	75	65	753
Reduced Dimension			
NonLinCFA ($\epsilon = 0.01$)	5.6 ± 0.4	NA	NA
NonLinCFA ($\epsilon = 0.001$)	7.2 ± 0.7	NA	NA
NonLinCFA ($\epsilon = 0.0001$)	7.4 ± 0.4	NA	NA
NonLinCFA ($\epsilon = 0.00001$)	7.4 ± 0.4	NA	NA
NonLinCFA ($\epsilon = 0.000001$)	7.4 ± 0.4	NA	NA
GenLinCFA ($\epsilon = \epsilon_1$)	15.2 ± 1.2	27.6 ± 5.6	53.4 ± 5.9
GenLinCFA ($\epsilon = \epsilon_2$)	8.0 ± 1.4	16.4 ± 3.4	26.4 ± 1.2
GenLinCFA ($\epsilon = \epsilon_3$)	4.8 ± 0.4	9.8 ± 1.5	23.4 ± 1.1
GenLinCFA ($\epsilon = \epsilon_4$)	4.2 ± 0.4	8.2 ± 0.4	20.2 ± 1.5
GenLinCFA ($\epsilon = \epsilon_5$)	3.8 ± 0.3	7.0 ± 0.6	14.2 ± 1.3
LinCFA	11.4 ± 0.7	NA	NA
PCA	26.6 ± 0.4	4.8 ± 0.9	77.6 ± 3.8
LDA	NA	1	1
Kernel PCA	36.0 ± 10.8	13.8 ± 5.1	32.4 ± 7.7
Isomap	27.2 ± 6.8	21.2 ± 9.8	14.8 ± 9.3
LLE	44.0 ± 7.3	1	39.6 ± 4.6
Supervised PCA	31.0 ± 13.8	16.2 ± 7.5	21.7 ± 3.4
NCA	NA	12.0 ± 9.5	28.2 ± 15.3
Test performance	R² score	Accuracy score	Accuracy score
NonLinCFA ($\epsilon = 0.01$)	0.8131 ± 0.0032	NA	NA
NonLinCFA ($\epsilon = 0.001$)	0.8061 ± 0.0076	NA	NA
NonLinCFA ($\epsilon = 0.0001$)	0.8133 ± 0.0037	NA	NA
NonLinCFA ($\epsilon = 0.00001$)	0.8133 ± 0.0037	NA	NA
NonLinCFA ($\epsilon = 0.000001$)	0.8136 ± 0.0036	NA	NA
GenLinCFA ($\epsilon = \epsilon_1$)	0.8104 ± 0.0015	0.7503 ± 0.0012	0.7984 ± 0.0112
GenLinCFA ($\epsilon = \epsilon_2$)	0.8119 ± 0.0010	0.7480 ± 0.0018	0.7647 ± 0.0130
GenLinCFA ($\epsilon = \epsilon_3$)	0.8101 ± 0.0003	0.7430 ± 0.0044	0.8016 ± 0.0069
GenLinCFA ($\epsilon = \epsilon_4$)	0.8114 ± 0.0003	0.7413 ± 0.0025	0.7808 ± 0.0071
GenLinCFA ($\epsilon = \epsilon_5$)	0.8107 ± 0.0009	0.7408 ± 0.0024	0.7712 ± 0.0095
LinCFA	0.8010 ± 0.0128	NA	NA
PCA	0.7559 ± 0.0027	0.7413 ± 0.0037	0.7840 ± 0.0117
LDA	NA	0.7587 ± 0.0065	0.7632 ± 0.0489
Kernel PCA	0.7764 ± 0.0118	0.7592 ± 0.0061	0.7920 ± 0.0177
Isomap	$0.2610 \pm .0458$	0.7463 ± 0.0039	0.7856 ± 0.0112
LLE	0.7281 ± 0.0267	0.7402 ± 0	0.7696 ± 0.0525
Supervised PCA	0.7731 ± 0.0126	0.7508 ± 0.0018	0.7913 ± 0.0069
NCA	NA	0.7637 ± 0.0079	0.7952 ± 0.0181