

Quadratic Programming Feature Selection for Multicorrelated Signal Decoding with Partial Least Squares[☆]

R.V. Isachenko^{a,*}, V.V. Strijov^b

^a *Moscow Institute of Physics and Technology, 9 Institutskiy Per., Dolgoprudny, Moscow Region 141700, Russian Federation*

^b *A. A. Dorodnicyn Computing Centre, Federal Research Center "Computer Science and Control" of the Russian Academy of Sciences, 40 Vavilov Str., Moscow 119333, Russian Federation*

Abstract

This paper is devoted to the dimensionality reduction problem in signal decoding. The challenge of this investigation is the redundancy in the data description. High correlations among measurements lead to correlations in the input space. This study considers the multivariate problem, and the target variable is a vector. In this case, the correlations occur in both the input and target spaces. Dimensionality reduction and feature selection are used to build simple and stable model.

The partial least squares (PLS) regression is used as the base model for the dimensionality reduction. The model projects the input and target data into the joint latent space and maximizes the covariances between the projections. To obtain the sparse model, the feature selection is applied. The majority of feature selection methods ignore the dependencies in the target space. The study suggests a novel approach to feature selection using a multivariate regression. The proposed approach extends the ideas of the quadratic programming feature selection (QPFS) algorithm. The QPFS algorithm selects the noncorrelated features that are relevant to the targets. The proposed methods take into account the dependencies in the target space and select the features that are informative to all targets jointly.

The computational experiment was carried out using the electrocorticogram (ECOG) dataset. The proposed algorithms were compared using different criteria such as their stability and predictive performance. The algorithms give significantly better results compared to the baseline strategy. The QPFS approach was compared with the partial least squares (PLS) regression. The best result is obtained by a combination of the QPFS and PLS algorithms.

Keywords: partial least squares, quadratic programming feature selection, signal decoding

[☆]The research was made possible by Government of the Russian Federation (Agreement 05.Y09.21.0018)

*Corresponding author

Email addresses: roman.isachenko@phystech.edu (R.V. Isachenko), strijov@ccas.ru (V.V. Strijov)

1. Introduction

The initial data in the fields of chemometrics [1, 2] and signal decoding [3, 4] are high-dimensional and extremely redundant. The models that are built on such data are instable. In addition, the redundant data description requires excessive computations, which lead to extended analysis times. To overcome this problem, dimensionality reduction [5, 6] and feature selection [7, 8] methods are used for high-dimensional data modeling.

The partial least squares (PLS) is a widely used algorithm for dimensionality reduction [9, 10, 11, 12]. The PLS finds the optimal combinations of the initial features and uses these combinations as the model features. The algorithm projects the features and the targets onto the joint latent space and maximizes the covariances between the projected vectors. It allows researchers to save the information about the initial input and target matrices and find their relations. The dimensionality of the latent space is much less than the size of the initial data description. It leads to a stable linear model built on a small number of features. An overview of the advances in the PLS regression is given in [13, 14]. For this model, we obtain the linear model with a small latent dimension. However, the final model uses the whole range of the initial features, and it does not allow for the removal of useless features.

Feature selection is a special case of dimensionality reduction when the latent representation is a subset of the initial data description. Here, the model is built on the subsets of the features. One of the approaches to feature selection is to maximize the feature relevances and minimize the pairwise feature redundancy. This approach was recently proposed and investigated in [15, 16]. Quadratic programmic feature selection (QPFS) [17] uses this approach to construct the optimization problem. It was shown in [18] that the QPFS algorithm outperforms many existing feature selection methods for the univariate regression problem. The QPFS algorithm introduces two functions: Sim and Rel. Sim estimates the redundancy between features, and Rel contains the relevances between each feature and the target vector. QPFS minimizes the Sim function and maximizes the Rel function simultaneously. The algorithm solves the following optimization problem:

$$(1 - \alpha) \cdot \underbrace{\mathbf{z}^T \mathbf{Q} \mathbf{z}}_{\text{Sim}(\mathbf{X})} - \alpha \cdot \underbrace{\mathbf{b}^T \mathbf{z}}_{\text{Rel}(\mathbf{X}, \boldsymbol{\nu})} \rightarrow \min_{\substack{\mathbf{z} \geq \mathbf{0}_n \\ \mathbf{1}_n^T \mathbf{z} = 1}}. \quad (1)$$

Here, the columns of matrix \mathbf{X} are the features, and $\boldsymbol{\nu}$ is the target vector. The matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$ entries measure the pairwise similarities between features. The vector $\mathbf{b} \in \mathbb{R}^n$ expresses the similarities between each feature and the target vector. The normalized vector \mathbf{z} shows the importance of each feature. Function (1) penalizes the dependent features using the Sim function and promotes the features that are relevant to the target using the Rel function. The parameter α controls the trade-off between Sim and Rel. To measure the similarity, the authors of [17] use the absolute value of the sample correlation coefficient between pairs of features for the Sim function, and between the features and the target vector for the Rel function.

Paper [19] proposes a multiway version of the QPFS algorithm for tensor ECoG-based data. It was shown that QPFS is an appropriate feature selection method for the signal decoding problem. We consider the multivariate problem, where the dependent variable is a vector. It leads to correlations in the model targets. In this situation, feature selection algorithms do not take into account these dependencies. Hence, the selected feature subset is not optimal in terms of its prediction. We propose methods that take into account the dependencies in both the input and target spaces. It allows us to form a stable sparse model. We refer to the original QPFS algorithm as our baseline for the computational experiment.

The main drawback of the QPFS algorithm is its computational costs. However, the original paper [17] suggests a way to solve the quadratic problem (1) efficiently. Additionally, in [20], the sequential minimal optimization framework is proposed for solving (1).

The experiments were carried out using the ECoG dataset [21]. We compared the proposed methods for multivariate feature selection with the baseline strategy and the PLS algorithm. The stability of the proposed methods was investigated by measuring how the feature selection solution changes with data bootstrapping. The proposed algorithms outperform the baseline algorithm given the same number of features. The combination of the feature selection procedure and the PLS algorithm gives the best performance.

The main contributions of this paper are as follows:

- addressing the dimensionality reduction problem for high-dimensional data,
- proposing new feature selection methods for multivariate regression with the analysis of the input and target spaces structures,
- comparing the proposed methods using a real ECoG dataset, and showing that the proposed methods give better feature subsets than the baseline method.

2. Multivariate regression

The goal is to forecast a dependent variable $\mathbf{y} \in \mathbb{R}^r$ with r targets from an independent input object $\mathbf{x} \in \mathbb{R}^n$ with n features. We assume that there is a linear dependence between the object \mathbf{x} and the target variable \mathbf{y} as

$$\mathbf{y} = \Theta \mathbf{x} + \boldsymbol{\varepsilon}, \quad (2)$$

where $\Theta \in \mathbb{R}^{r \times n}$ is the matrix of the model parameters, and $\boldsymbol{\varepsilon} \in \mathbb{R}^r$ is a residual vector. One has to find the matrix of the model parameters Θ given a dataset (\mathbf{X}, \mathbf{Y}) , where $\mathbf{X} \in \mathbb{R}^{m \times n}$ is a design matrix and $\mathbf{Y} \in \mathbb{R}^{m \times r}$ is a target matrix:

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]^\top = [\boldsymbol{\chi}_1, \dots, \boldsymbol{\chi}_n]; \quad \mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_m]^\top = [\boldsymbol{\nu}_1, \dots, \boldsymbol{\nu}_r].$$

The columns $\boldsymbol{\chi}_j$ of \mathbf{X} correspond to the object features, and the columns $\boldsymbol{\nu}_j$ of \mathbf{Y} correspond to the targets.

The optimal parameters are determined by the minimization of an error function. We define the quadratic loss function as follows:

$$\mathcal{L}(\Theta|\mathbf{X}, \mathbf{Y}) = \left\| \begin{matrix} \mathbf{Y} & - & \mathbf{X} & \cdot & \Theta^\top \\ m \times r & & m \times n & & r \times n \end{matrix} \right\|_2^2 \rightarrow \min_{\Theta}. \quad (3)$$

The solution of (3) is given by

$$\Theta = \mathbf{Y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}.$$

The linear dependent columns of \mathbf{X} lead to an instable solution for the optimization problem (3). If there is a vector $\alpha \neq \mathbf{0}_n$ such that $\mathbf{X}\alpha = \mathbf{0}_m$, then adding α to any column of Θ does not change the value of the loss function $\mathcal{L}(\Theta|\mathbf{X}, \mathbf{Y})$. In this case, the matrix $\mathbf{X}^\top \mathbf{X}$ is close to singular and is not invertible. To avoid strong linear dependence, dimensionality reduction and feature selection are used.

3. Feature selection

The feature selection goal is to find the boolean vector $\mathbf{a} = \{0, 1\}^n$ in which the components indicate whether the feature is selected. To obtain the optimal vector \mathbf{a} among all possible $2^n - 1$ options, we introduce the feature selection error function $S(\mathbf{a}|\mathbf{X}, \mathbf{Y})$. We state the feature selection problem as follows:

$$\mathbf{a} = \arg \min_{\mathbf{a}' \in \{0, 1\}^n} S(\mathbf{a}'|\mathbf{X}, \mathbf{Y}). \quad (4)$$

The goal of feature selection is to construct the appropriate function $S(\mathbf{a}|\mathbf{X}, \mathbf{Y})$. The particular examples for the considered feature selection algorithms are given below and summarized in Table 1.

Problem (4) is hard to solve due to the discrete binary domain $\{0, 1\}^n$. We relax problem (4) to the continuous domain $[0, 1]^n$. The relaxed feature selection problem is

$$\mathbf{z} = \arg \min_{\mathbf{z}' \in [0, 1]^n} S(\mathbf{z}'|\mathbf{X}, \mathbf{Y}). \quad (5)$$

Here, the vector \mathbf{z} entries are the normalized feature importances. First, we solve problem (5) to obtain the feature importances \mathbf{z} . Then, the solution of (4) is recovered by thresholding as follows:

$$\mathbf{a} = [a_j]_{j=1}^n, \quad a_j = \begin{cases} 1, & z_j > \tau; \\ 0, & \text{otherwise.} \end{cases}$$

τ is a hyperparameter that is defined manually or chosen by cross-validation.

Once the solution \mathbf{a} of (4) is known, problem (3) becomes

$$\mathcal{L}(\Theta_{\mathbf{a}}|\mathbf{X}_{\mathbf{a}}, \mathbf{Y}) = \left\| \mathbf{Y} - \mathbf{X}_{\mathbf{a}} \Theta_{\mathbf{a}}^\top \right\|_2^2 \rightarrow \min_{\Theta_{\mathbf{a}}},$$

where subscript \mathbf{a} indicates the sub matrix with the columns in which the components of \mathbf{a} equal 1.

3.1. Quadratic Programming Feature Selection

Paper [18] shows that QPFS outperforms many existing feature selection algorithms using different quality criteria. The QPFS algorithm selects the noncorrelated features that are relevant to the target vector $\boldsymbol{\nu}$ for the linear regression problem where $r = 1$ as follows:

$$\|\boldsymbol{\nu} - \mathbf{X}\boldsymbol{\theta}\|_2^2 \rightarrow \min_{\boldsymbol{\theta} \in \mathbb{R}^n}.$$

The authors of the original QPFS paper [17] suggested the following way to select α for (1) and make $\text{Sim}(\mathbf{X})$ and $\text{Rel}(\mathbf{X}, \boldsymbol{\nu})$ have the same impacts:

$$\alpha = \frac{\overline{\mathbf{Q}}}{\overline{\mathbf{Q}} + \overline{\mathbf{b}}}, \quad \overline{\mathbf{Q}} = \text{mean}(\mathbf{Q}), \quad \overline{\mathbf{b}} = \text{mean}(\mathbf{b}).$$

The QPFS parameters are defined as follows:

$$\mathbf{Q} = [|\text{corr}(\boldsymbol{\chi}_i, \boldsymbol{\chi}_j)|]_{i,j=1}^n, \quad \mathbf{b} = [|\text{corr}(\boldsymbol{\chi}_i, \boldsymbol{\nu})|]_{i=1}^n. \quad (6)$$

Here $\text{corr}(\cdot, \cdot)$, is the absolute value of the sample Pearson correlation coefficient:

$$\text{corr}(\boldsymbol{\chi}, \boldsymbol{\nu}) = \frac{\sum_{i=1}^m (\boldsymbol{\chi}_i - \overline{\boldsymbol{\chi}})(\boldsymbol{\nu}_i - \overline{\boldsymbol{\nu}})}{\sqrt{\sum_{i=1}^m (\boldsymbol{\chi}_i - \overline{\boldsymbol{\chi}})^2 \sum_{i=1}^m (\boldsymbol{\nu}_i - \overline{\boldsymbol{\nu}})^2}}.$$

Other ways to define \mathbf{Q} and \mathbf{b} are considered in [18].

Problem (1) is convex if the matrix \mathbf{Q} is positive semidefinite. In general, this is not always true. To satisfy this condition, the matrix \mathbf{Q} spectrum is shifted and matrix \mathbf{Q} is replaced by $\mathbf{Q} - \lambda_{\min} \mathbf{I}$, where λ_{\min} is the minimum eigenvalue of \mathbf{Q} .

3.2. Multivariate QPFS

Here, we propose the algorithms for feature selection in the multivariate case. If the target space is multidimensional, it is prone to redundancy and correlations between the targets. In this section, we propose the algorithms that take into account the dependencies in both the input and target spaces.

Relevance aggregation (RelAgg). In [19], in order to apply the QPFS algorithm to the multivariate case ($r > 1$), feature relevances are aggregated through all r components. The term $\text{Sim}(\mathbf{X})$ is still the same, and matrix \mathbf{Q} is defined by (6). The vector \mathbf{b} is aggregated across all targets and is defined as

$$\mathbf{b} = \left[\sum_{k=1}^r |\text{corr}(\boldsymbol{\chi}_i, \boldsymbol{\nu}_k)| \right]_{i=1}^n.$$

The drawback of this approach is its insensitivity to the dependencies in the columns of \mathbf{Y} . Observe the following example:

$$\mathbf{X} = [\boldsymbol{\chi}_1, \boldsymbol{\chi}_2, \boldsymbol{\chi}_3], \quad \mathbf{Y} = [\underbrace{\boldsymbol{\nu}_1, \boldsymbol{\nu}_1, \dots, \boldsymbol{\nu}_1}_{r-1}, \boldsymbol{\nu}_2].$$

We have 3 features and r targets, where the first $r - 1$ targets are identical. The pairwise features similarities are given by matrix \mathbf{Q} . Matrix \mathbf{B} entries give the pairwise features relevances to the targets. Vector \mathbf{b} is obtained by the summation of matrix \mathbf{B} over the columns

$$\mathbf{Q} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0.8 \\ 0 & 0.8 & 1 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 0.4 & \dots & 0.4 & 0 \\ 0.5 & \dots & 0.5 & 0.8 \\ 0.8 & \dots & 0.8 & 0.1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} (r-1) \cdot 0.4 + 0 \\ (r-1) \cdot 0.5 + 0.8 \\ (r-1) \cdot 0.8 + 0.1 \end{bmatrix}. \quad (7)$$

We would like to select only two features. For such a configuration, the best feature subset is $[\chi_1, \chi_2]$. Feature χ_2 predicts the second target ν_2 and the combination of features χ_1, χ_2 predicts the first component. The QPFS algorithm for $r = 2$ gives the solution $\mathbf{z} = [0.37, 0.61, 0.02]$. It coincides with our knowledge. However, if we add the collinear columns to matrix \mathbf{Y} and increase r to 5, the QPFS solution will be $\mathbf{z} = [0.40, 0.17, 0.43]$. Here, we lose the relevant feature χ_2 and select the redundant feature χ_3 . The following subsections propose extensions to the QPFS algorithm that overcome the challenge of this example.

Symmetric importances (SymImp). To take into account the dependencies in the columns of matrix \mathbf{Y} , we extend the QPFS function (1) to the multivariate case. We add the term $\text{Sim}(\mathbf{Y})$ and modify the term $\text{Rel}(\mathbf{X}, \mathbf{Y})$ as follows:

$$\alpha_1 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x}_{\text{Sim}(\mathbf{X})} - \alpha_2 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y}_{\text{Rel}(\mathbf{X}, \mathbf{Y})} + \alpha_3 \cdot \underbrace{\mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y}_{\text{Sim}(\mathbf{Y})} \rightarrow \min_{\substack{\mathbf{z}_x \geq \mathbf{0}_n, \mathbf{1}_n^\top \mathbf{z}_x = 1 \\ \mathbf{z}_y \geq \mathbf{0}_r, \mathbf{1}_r^\top \mathbf{z}_y = 1}}. \quad (8)$$

We determine the entries of matrices $\mathbf{Q}_x \in \mathbb{R}^{n \times n}$, $\mathbf{Q}_y \in \mathbb{R}^{r \times r}$, and $\mathbf{B} \in \mathbb{R}^{n \times r}$ in the following way:

$$\mathbf{Q}_x = [|\text{corr}(\chi_i, \chi_j)|]_{i,j=1}^n, \quad \mathbf{Q}_y = [|\text{corr}(\nu_i, \nu_j)|]_{i,j=1}^r, \quad \mathbf{B} = [|\text{corr}(\chi_i, \nu_j)|]_{i=1, \dots, n, j=1, \dots, r}.$$

Vector \mathbf{z}_x shows the features' importances, while \mathbf{z}_y is a vector of the targets importances. The correlated targets will be penalized by $\text{Sim}(\mathbf{Y})$ and have lower importances.

The coefficients α_1 , α_2 , and α_3 control the influence of each term on function (8) and satisfy the following conditions:

$$\alpha_1 + \alpha_2 + \alpha_3 = 1, \quad \alpha_i \geq 0, \quad i = 1, 2, 3.$$

Proposition 1. *The balance between the terms $\text{Sim}(\mathbf{X})$, $\text{Rel}(\mathbf{X}, \mathbf{Y})$, and $\text{Sim}(\mathbf{Y})$ for the problem (8) is achieved by the following coefficients:*

$$\alpha_1 \propto \overline{\mathbf{Q}_y} \overline{\mathbf{B}}; \quad \alpha_2 \propto \overline{\mathbf{Q}_x} \overline{\mathbf{Q}_y}; \quad \alpha_3 \propto \overline{\mathbf{Q}_x} \overline{\mathbf{B}}. \quad (9)$$

Proof. The desired values of α_1 , α_2 , and α_3 are given by solving of the following equations:

$$\begin{aligned} \alpha_1 + \alpha_2 + \alpha_3 &= 1; \\ \alpha_1 \overline{\mathbf{Q}_x} &= \alpha_2 \overline{\mathbf{B}} = \alpha_3 \overline{\mathbf{Q}_y}. \end{aligned}$$

Here, the mean values $\overline{\mathbf{Q}_x}$, $\overline{\mathbf{B}}$, and $\overline{\mathbf{Q}_y}$ of the corresponding matrices \mathbf{Q}_x , \mathbf{B} , and \mathbf{Q}_y are the mean values of the terms $\text{Sim}(\mathbf{X})$, $\text{Rel}(\mathbf{X}, \mathbf{Y})$, and $\text{Sim}(\mathbf{Y})$. \square

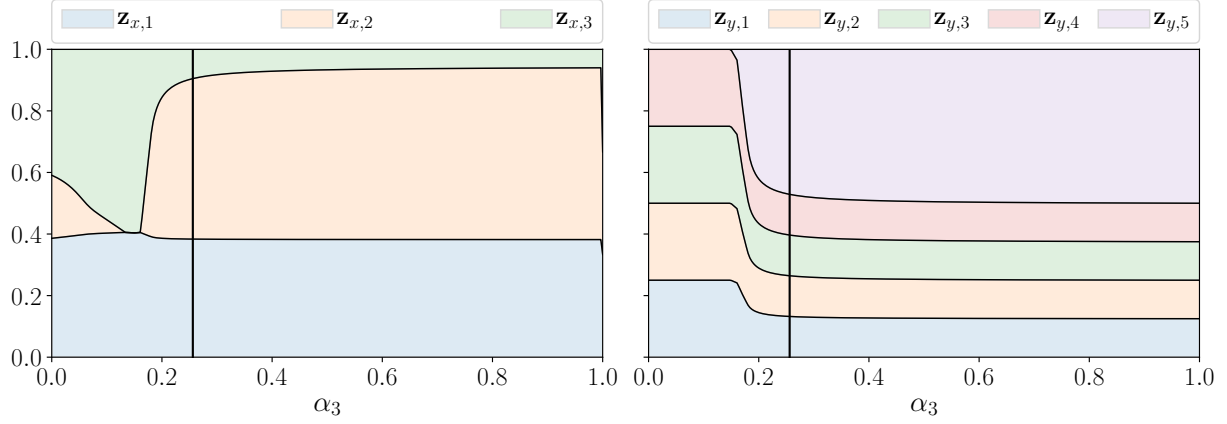


Figure 1: Feature importances \mathbf{z}_x and \mathbf{z}_y with respect to α_3 for the considered example

To investigate the impact of $\text{Sim}(\mathbf{Y})$ on function (8), we balance the terms $\text{Sim}(\mathbf{X})$ and $\text{Rel}(\mathbf{X}, \mathbf{Y})$ by fixing the proportion between α_1 and α_2 :

$$\alpha_1 = \frac{(1 - \alpha_3)\bar{\mathbf{B}}}{\bar{\mathbf{Q}}_x + \bar{\mathbf{B}}}; \quad \alpha_2 = \frac{(1 - \alpha_3)\bar{\mathbf{Q}}_x}{\bar{\mathbf{Q}}_x + \bar{\mathbf{B}}}; \quad \alpha_3 \in [0, 1]. \quad (10)$$

We apply the proposed algorithm to the discussed example (7). The given matrix \mathbf{Q} corresponds to matrix \mathbf{Q}_x . We additionally define matrix \mathbf{Q}_y by setting $\text{corr}(\boldsymbol{\nu}_1, \boldsymbol{\nu}_2) = 0.2$ and all others entries to one. Figure 1 shows the importances of features \mathbf{z}_x and targets \mathbf{z}_y with respect to α_3 . If α_3 is small, the impacts of all targets are almost identical and feature χ_3 dominates feature χ_2 . When α_3 becomes larger than 0.2, the importance $\mathbf{z}_{y,5}$ of target $\boldsymbol{\nu}_5$ increases along with the importance of feature χ_2 .

Minimax QPFS (MinMax). Function (8) is symmetric with respect to \mathbf{z}_x and \mathbf{z}_y . It penalizes the features that are correlated and irrelevant to the targets. In addition, it penalizes the targets that are correlated and are not sufficiently explained by the features. It leads to small importances for the targets that are weakly correlated with the features and large importances for the targets that are strongly correlated with the features. This result contradicts the intuition. Our goal is to predict all targets, especially those that are difficult to explain, using the selected relevant and noncorrelated features. We express this as two related problems:

$$\alpha_1 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x}_{\text{Sim}(\mathbf{X})} - \alpha_2 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y}_{\text{Rel}(\mathbf{X}, \mathbf{Y})} \rightarrow \min_{\substack{\mathbf{z}_x \geq \mathbf{0}_n, \\ \mathbf{1}_n^\top \mathbf{z}_x = 1}}; \quad (11)$$

$$\alpha_3 \cdot \underbrace{\mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y}_{\text{Sim}(\mathbf{Y})} + \alpha_2 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y}_{\text{Rel}(\mathbf{X}, \mathbf{Y})} \rightarrow \min_{\substack{\mathbf{z}_y \geq \mathbf{0}_r, \\ \mathbf{1}_r^\top \mathbf{z}_y = 1}}. \quad (12)$$

The difference between (11) and (12) is the sign of Rel. In the input space, the nonrelevant components should have smaller importances. Meanwhile, the targets that are not relevant

to the features should have larger importances. Problems (11) and (12) are merged into the joint min-max or max-min formulation

$$\min_{\substack{\mathbf{z}_x \geq \mathbf{0}_n \\ \mathbf{1}_n^\top \mathbf{z}_x = 1}} \max_{\substack{\mathbf{z}_y \geq \mathbf{0}_r \\ \mathbf{1}_r^\top \mathbf{z}_y = 1}} f(\mathbf{z}_x, \mathbf{z}_y), \quad \left(\text{or } \max_{\substack{\mathbf{z}_y \geq \mathbf{0}_r \\ \mathbf{1}_r^\top \mathbf{z}_y = 1}} \min_{\substack{\mathbf{z}_x \geq \mathbf{0}_n \\ \mathbf{1}_n^\top \mathbf{z}_x = 1}} f(\mathbf{z}_x, \mathbf{z}_y) \right), \quad (13)$$

where

$$f(\mathbf{z}_x, \mathbf{z}_y) = \alpha_1 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x}_{\text{Sim}(\mathbf{X})} - \alpha_2 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y}_{\text{Rel}(\mathbf{X}, \mathbf{Y})} - \alpha_3 \cdot \underbrace{\mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y}_{\text{Sim}(\mathbf{Y})}.$$

Theorem 1. *For positive definite matrices \mathbf{Q}_x and \mathbf{Q}_y , the max-min and min-max problems (13) have the same optimal value.*

Proof. We denote the following:

$$\mathbb{C}^n = \{\mathbf{z} : \mathbf{z} \geq \mathbf{0}_n, \mathbf{1}_n^\top \mathbf{z} = 1\}, \quad \mathbb{C}^r = \{\mathbf{z} : \mathbf{z} \geq \mathbf{0}_r, \mathbf{1}_r^\top \mathbf{z} = 1\}.$$

The sets \mathbb{C}^n and \mathbb{C}^r are compact and convex. The function $f : \mathbb{C}^n \times \mathbb{C}^r \rightarrow \mathbb{R}$ is a continuous function. If \mathbf{Q}_x and \mathbf{Q}_y are positive definite matrices, function f is convex-concave. I.e., $f(\cdot, \mathbf{z}_y) : \mathbb{C}^n \rightarrow \mathbb{R}$ is convex for a fixed \mathbf{z}_y , and $f(\mathbf{z}_x, \cdot) : \mathbb{C}^r \rightarrow \mathbb{R}$ is concave for a fixed \mathbf{z}_x . In this case, Neumann's minimax theorem states that

$$\min_{\mathbf{z}_x \in \mathbb{C}^n} \max_{\mathbf{z}_y \in \mathbb{C}^r} f(\mathbf{z}_x, \mathbf{z}_y) = \max_{\mathbf{z}_y \in \mathbb{C}^r} \min_{\mathbf{z}_x \in \mathbb{C}^n} f(\mathbf{z}_x, \mathbf{z}_y).$$

□

To solve the min-max problem (13), we fix some $\mathbf{z}_x \in \mathbb{C}^n$. For a fixed vector \mathbf{z}_x , we solve the problem

$$\max_{\substack{\mathbf{z}_y \in \mathbb{C}^r \\ \mathbf{1}_r^\top \mathbf{z}_y = 1}} f(\mathbf{z}_x, \mathbf{z}_y) = \max_{\substack{\mathbf{z}_y \geq \mathbf{0}_r \\ \mathbf{1}_r^\top \mathbf{z}_y = 1}} [\alpha_1 \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha_2 \cdot \mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y - \alpha_3 \cdot \mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y]. \quad (14)$$

The Lagrangian for this problem is

$$L(\mathbf{z}_x, \mathbf{z}_y, \lambda, \boldsymbol{\mu}) = \alpha_1 \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha_2 \cdot \mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y - \alpha_3 \cdot \mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y + \lambda \cdot (\mathbf{1}_r^\top \mathbf{z}_y - 1) + \boldsymbol{\mu}^\top \mathbf{z}_y.$$

Here, the Lagrange multipliers $\boldsymbol{\mu}$ that correspond to the inequality constraints $\mathbf{z}_y \geq \mathbf{0}_r$ are restricted to being nonnegative. The dual problem is

$$\min_{\lambda, \boldsymbol{\mu} \geq \mathbf{0}_r} g(\mathbf{z}_x, \lambda, \boldsymbol{\mu}) = \min_{\lambda, \boldsymbol{\mu} \geq \mathbf{0}_r} \left[\max_{\mathbf{z}_y \in \mathbb{R}^r} L(\mathbf{z}_x, \mathbf{z}_y, \lambda, \boldsymbol{\mu}) \right]. \quad (15)$$

The strong duality holds for quadratic problem (14) with the positive definite matrices \mathbf{Q}_x and \mathbf{Q}_y . Therefore, the optimal value for (14) equals the optimal value for (15). It allows us to solve the problem

$$\min_{\mathbf{z}_x \in \mathbb{C}^n, \lambda, \boldsymbol{\mu} \geq \mathbf{0}_r} g(\mathbf{z}_x, \lambda, \boldsymbol{\mu}) \quad (16)$$

instead of (13).

By setting the gradient of the Langrangian $\nabla_{\mathbf{z}_y} L(\mathbf{z}_x, \mathbf{z}_y, \lambda, \boldsymbol{\mu})$ equal to zero, we obtain an optimal value for \mathbf{z}_y :

$$\mathbf{z}_y = \frac{1}{2\alpha_3} \mathbf{Q}_y^{-1} (-\alpha_2 \cdot \mathbf{B}^\top \mathbf{z}_x + \lambda \cdot \mathbf{1}_r + \boldsymbol{\mu}). \quad (17)$$

The dual function is equal to

$$\begin{aligned} g(\mathbf{z}_x, \lambda, \boldsymbol{\mu}) = \max_{\mathbf{z}_y \in \mathbb{R}^r} L(\mathbf{z}_x, \mathbf{z}_y, \lambda, \boldsymbol{\mu}) &= \mathbf{z}_x^\top \left(-\frac{\alpha_2^2}{4\alpha_3} \cdot \mathbf{B} \mathbf{Q}_y^{-1} \mathbf{B}^\top - \alpha_1 \cdot \mathbf{Q}_x \right) \mathbf{z}_x \\ &\quad - \frac{1}{4\alpha_3} \lambda^2 \cdot \mathbf{1}_r^\top \mathbf{Q}_y^{-1} \mathbf{1}_r - \frac{1}{4\alpha_3} \cdot \boldsymbol{\mu}^\top \mathbf{Q}_y^{-1} \boldsymbol{\mu} + \frac{\alpha_2}{2\alpha_3} \lambda \cdot \mathbf{1}_r^\top \mathbf{Q}_y^{-1} \mathbf{B}^\top \mathbf{z}_x \\ &\quad - \frac{1}{2\alpha_3} \lambda \cdot \mathbf{1}_r^\top \mathbf{Q}_y^{-1} \boldsymbol{\mu} + \frac{\alpha_2}{2\alpha_3} \cdot \boldsymbol{\mu}^\top \mathbf{Q}_y^{-1} \mathbf{B}^\top \mathbf{z}_x + \lambda. \end{aligned} \quad (18)$$

It represents the quadratic problem (16) with $n + r + 1$ variables.

Asymmetric Importance (AsymImp). The natural way to overcome the problem of the SymImp strategy is to add penalties for targets that are correlated with features. We add the term $\mathbf{b}^\top \mathbf{z}_y$ to the term $\text{Rel}(\mathbf{X}, \mathbf{Y})$ as follows:

$$\alpha_1 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x}_{\text{Sim}(\mathbf{X})} - \alpha_2 \cdot \underbrace{(\mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y - \mathbf{b}^\top \mathbf{z}_y)}_{\text{Rel}(\mathbf{X}, \mathbf{Y})} + \alpha_3 \cdot \underbrace{\mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y}_{\text{Sim}(\mathbf{Y})} \rightarrow \min_{\substack{\mathbf{z}_x \geq \mathbf{0}_n, \mathbf{1}_n^\top \mathbf{z}_x = 1 \\ \mathbf{z}_y \geq \mathbf{0}_r, \mathbf{1}_r^\top \mathbf{z}_y = 1}}. \quad (19)$$

Proposition 2. *Let vector \mathbf{b} equal*

$$b_j = \max_{i=1, \dots, n} [\mathbf{B}]_{i,j}.$$

Then, the importance coefficients for vector \mathbf{z}_y will be nonnegative in $\text{Rel}(\mathbf{X}, \mathbf{Y})$ for problem (19).

Proof. The proposition follows from the fact that

$$\sum_{i=1}^n z_i b_{ij} \leq \left(\sum_{i=1}^n z_i \right) \max_{i=1, \dots, n} b_{ij} = \max_{i=1, \dots, n} b_{ij},$$

where $z_i \geq 0$ and $\sum_{i=1}^n z_i = 1$. □

Hence, function (19) encourages the features that are relevant to the targets and encourages the targets that are not sufficiently correlated with the features.

Proposition 3. *The balance between the terms $\text{Sim}(\mathbf{X})$, $\text{Rel}(\mathbf{X}, \mathbf{Y})$, and $\text{Sim}(\mathbf{Y})$ for the problem (19) is achieved by the following coefficients:*

$$\alpha_1 \propto \overline{\mathbf{Q}}_y (\overline{\mathbf{b}} - \overline{\mathbf{B}}); \quad \alpha_2 \propto \overline{\mathbf{Q}}_x \overline{\mathbf{Q}}_y; \quad \alpha_3 \propto \overline{\mathbf{Q}}_x \overline{\mathbf{B}}.$$

Table 1: Overview of the proposed multivariate QPFS algorithms

| Algorithm | Idea | Error function $S(\mathbf{z} \mathbf{X}, \mathbf{Y})$ |
|-----------|--|---|
| RelAgg | $\min [\text{Sim}(\mathbf{X}) - \text{Rel}(\mathbf{X}, \mathbf{Y})]$ | $\min_{\mathbf{z}_x} [(1 - \alpha) \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha \cdot \mathbf{z}_x^\top \mathbf{B} \mathbf{1}_r]$ |
| SymImp | $\min [\text{Sim}(\mathbf{X}) - \text{Rel}(\mathbf{X}, \mathbf{Y}) + \text{Sim}(\mathbf{Y})]$ | $\min_{\mathbf{z}_x, \mathbf{z}_y} [\alpha_1 \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha_2 \cdot \mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y + \alpha_3 \cdot \mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y]$ |
| MinMax | $\min [\text{Sim}(\mathbf{X}) - \text{Rel}(\mathbf{X}, \mathbf{Y})]$ $\max [\text{Rel}(\mathbf{X}, \mathbf{Y}) + \text{Sim}(\mathbf{Y})]$ | $\min_{\mathbf{z}_x} \max_{\mathbf{z}_y} [\alpha_1 \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha_2 \cdot \mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y - \alpha_3 \cdot \mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y]$ |
| AsymImp | $\min [\text{Sim}(\mathbf{X}) - \text{Rel}(\mathbf{X}, \mathbf{Y})]$ $\max [\text{Rel}(\mathbf{X}, \mathbf{Y}) + \text{Sim}(\mathbf{Y})]$ | $\min_{\mathbf{z}_x, \mathbf{z}_y} [\alpha_1 \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha_2 \cdot (\mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y - \mathbf{b}^\top \mathbf{z}_y) + \alpha_3 \cdot \mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y]$ |

Proof. The desired values of α_1 , α_2 , and α_3 are given by the solutions to the following equations:

$$\alpha_1 + \alpha_2 + \alpha_3 = 1; \quad (20)$$

$$\alpha_1 \overline{\mathbf{Q}}_x = \alpha_2 \overline{\mathbf{B}}; \quad (21)$$

$$\alpha_2 (\overline{\mathbf{b}} - \overline{\mathbf{B}}) = \alpha_3 \overline{\mathbf{Q}}_y. \quad (22)$$

Here, we balance $\text{Sim}(\mathbf{X})$ with the first term of $\text{Rel}(\mathbf{X}, \mathbf{Y})$ using (21) and $\text{Sim}(\mathbf{Y})$ with the full $\text{Rel}(\mathbf{X}, \mathbf{Y})$ using (22). \square

Proposition 4. *For the case of $r = 1$, the proposed functions (8), (13), and (19) coincide with the original QPFS algorithm (1).*

Proof. If r is equal to 1, then $\mathbf{Q}_y = q_y$ is a scalar, $\mathbf{z}_y = 1$, and $\mathbf{B} = \mathbf{b}$. It reduces problems (8), (13), and (19) to

$$\alpha_1 \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha_2 \cdot \mathbf{z}_x^\top \mathbf{b} \rightarrow \min_{\mathbf{z}_x \geq \mathbf{0}_n, \mathbf{1}_n^\top \mathbf{z}_x = 1}.$$

Setting $\alpha = \frac{\alpha_2}{\alpha_1 + \alpha_2}$ represents the original QPFS problem (1). \square

Table 1 shows the core ideas and error functions for each method and summarizes all the proposed strategies for multivariate feature selection. RelAgg is the baseline strategy, and it does not consider the target space correlations. SymImp penalizes the pairwise target correlations. MinMax are more sensitive to the targets that are difficult to predict. The AsymImp strategy adds the term to the SymImp function to make the features and targets have asymmetric influences.

3.3. Dimensionality reduction

To eliminate the linear dependence and reduce the dimensionality of the input space, principal components analysis (PCA) is a widely used algorithm. The main disadvantage

of the PCA method is that it is insensitive to the interrelation between the features and the targets. The partial least squares algorithm projects the design matrix \mathbf{X} and the target matrix \mathbf{Y} to the latent space with low dimensionality ($l < n$). The PLS algorithm finds the latent space matrices $\mathbf{T}, \mathbf{U} \in \mathbb{R}^{m \times l}$ that best describe the original matrices \mathbf{X} and \mathbf{Y} .

The design matrix \mathbf{X} and the target matrix \mathbf{Y} are projected into the latent space in the following way:

$$\underset{m \times n}{\mathbf{X}} = \underset{m \times l}{\mathbf{T}} \cdot \underset{l \times n}{\mathbf{P}^\top} + \underset{m \times n}{\mathbf{F}} = \sum_{k=1}^l \underset{m \times 1}{\mathbf{t}_k} \cdot \underset{1 \times n}{\mathbf{p}_k^\top} + \underset{m \times n}{\mathbf{F}}, \quad (23)$$

$$\underset{m \times r}{\mathbf{Y}} = \underset{m \times l}{\mathbf{U}} \cdot \underset{l \times r}{\mathbf{Q}^\top} + \underset{m \times r}{\mathbf{E}} = \sum_{k=1}^l \underset{m \times 1}{\mathbf{u}_k} \cdot \underset{1 \times r}{\mathbf{q}_k^\top} + \underset{m \times r}{\mathbf{E}}, \quad (24)$$

where \mathbf{T} and \mathbf{U} are the scores matrices in the latent space, \mathbf{P} and \mathbf{Q} are the loading matrices, \mathbf{E} and \mathbf{F} are residual matrices. The PLS maximizes the linear relation between the columns of matrices \mathbf{T} and \mathbf{U} as

$$\mathbf{U} \approx \mathbf{T}\mathbf{B}, \quad \mathbf{B} = \text{diag}(\beta_k), \quad \beta_k = \mathbf{u}_k^\top \mathbf{t}_k / (\mathbf{t}_k^\top \mathbf{t}_k).$$

We use the PLS algorithm as the dimensionality reduction algorithm in this research.

To obtain the model prediction and find the model parameters, we multiply both sides of (23) by \mathbf{W} . Since the residual matrix \mathbf{E} rows are orthogonal to the columns of \mathbf{W} , we have

$$\mathbf{X}\mathbf{W} = \mathbf{T}\mathbf{P}^\top\mathbf{W}.$$

The linear transformation between objects in the input and latent spaces is the following

$$\mathbf{T} = \mathbf{X}\mathbf{W}^*, \quad \text{where } \mathbf{W}^* = \mathbf{W}(\mathbf{P}^\top\mathbf{W})^{-1}. \quad (25)$$

The matrix of the model parameters (2) could be found from equations (24) and (25) as

$$\mathbf{Y} = \mathbf{U}\mathbf{Q}^\top + \mathbf{E} \approx \mathbf{T}\mathbf{B}\mathbf{Q}^\top + \mathbf{E} = \mathbf{X}\mathbf{W}^*\mathbf{B}\mathbf{Q}^\top + \mathbf{E} = \mathbf{X}\boldsymbol{\Theta} + \mathbf{E}. \quad (26)$$

Thus, the model parameters (2) are equal to

$$\boldsymbol{\Theta} = \mathbf{W}(\mathbf{P}^\top\mathbf{W})^{-1}\mathbf{B}\mathbf{Q}^\top.$$

The final model (26) is a linear model that are low-dimensional in the latent space. It reduces the data redundancy and increases the model stability.

4. Experiment

To evaluate the selected feature subset, we introduce criteria that estimate the quality of feature selection. We measure the smulticorrelation using the mean value of miltiple correlation coefficient as follows:

$$R^2 = \frac{1}{r} \text{tr}(\mathbf{C}^\top \mathbf{R}^{-1} \mathbf{C}); \quad \text{where } \mathbf{C} = [\text{corr}(\boldsymbol{\chi}_i, \boldsymbol{\nu}_j)]_{\substack{i=1, \dots, n \\ j=1, \dots, r}}, \quad \mathbf{R} = [\text{corr}(\boldsymbol{\chi}_i, \boldsymbol{\chi}_j)]_{i,j=1}^n.$$

This coefficient lies between 0 and 1. A bigger R^2 means that we have a better feature subset.

The model stability is given by the logarithmic ratio between the minimum eigenvalue λ_{\min} and maximum eigenvalue λ_{\max} of matrix $\mathbf{X}^\top \mathbf{X}$:

$$\text{Stability} = \ln \frac{\lambda_{\min}}{\lambda_{\max}}.$$

A smaller Stability value indicates less multicollinearity in matrix \mathbf{X} .

The scaled Root Mean Squared Error (sRMSE) shows the quality of the model prediction. We estimate the sRMSE using train and test data.

$$\text{sRMSE}(\mathbf{Y}, \hat{\mathbf{Y}}_{\mathbf{a}}) = \sqrt{\frac{\text{MSE}(\mathbf{Y}, \hat{\mathbf{Y}}_{\mathbf{a}})}{\text{MSE}(\mathbf{Y}, \bar{\mathbf{Y}})}} = \frac{\|\mathbf{Y} - \hat{\mathbf{Y}}_{\mathbf{a}}\|_2}{\|\mathbf{Y} - \bar{\mathbf{Y}}\|_2}.$$

Here, $\hat{\mathbf{Y}}_{\mathbf{a}} = \mathbf{X}_{\mathbf{a}} \Theta_{\mathbf{a}}^\top$ is the model prediction and $\bar{\mathbf{Y}}$ is the constant prediction obtained by averaging the targets across all objects. The error on the test set should be as small as possible.

The Bayesian Information Criteria (BIC) incorporates a trade-off between the prediction quality and the size of selected subset $\|\mathbf{a}\|_0 = \#\{j : a_j \neq 0\} = \sum_{j=1}^n a_j$:

$$\text{BIC} = m \ln \left(\text{MSE}(\mathbf{Y}, \hat{\mathbf{Y}}_{\mathbf{a}}) \right) + \|\mathbf{a}\|_0 \cdot \ln m,$$

A smaller value of BIC means a better feature subset.

4.1. Data

We conducted a computational experiment with the ECoG data from the NeuroTycho project [21]. The ECoG data consist of brain voltage signals recorded over 32 channels. The goal is to predict 3D hand positions in subsequent moments given the input signal. The initial voltage signals are transformed to the spatial-temporal representation using the wavelet transformation with the Morlet mother wavelet. The procedure of extracting the feature representation from the raw data is described in detail in [22, 23]. We unfold the data and feature description at each time moment has dimension of size 32 (channels) \times 27 (frequencies) = 864. Each object is the representation of the local historical time segment with duration $\Delta t = 1s$. The time step between objects is $\delta t = 0.05s$. The final matrices are $\mathbf{X} \in \mathbb{R}^{18900 \times 864}$ and $\mathbf{Y} \in \mathbb{R}^{18900 \times 3k}$, where k is the number of timestamps that we predict. We split our data into train and test parts with the ratio 0.67. Example of the initial brain signals and the corresponding hand trajectory is shown in Figure 2.

4.2. Results

Figure 3 shows the dependencies in the matrices \mathbf{X} and \mathbf{Y} for the ECoG data. The frequencies in the matrix \mathbf{X} are highly correlated. In the target matrix \mathbf{Y} , the correlations between axes are not significant in comparison with the correlations between consequent moments and these correlations decay with time.

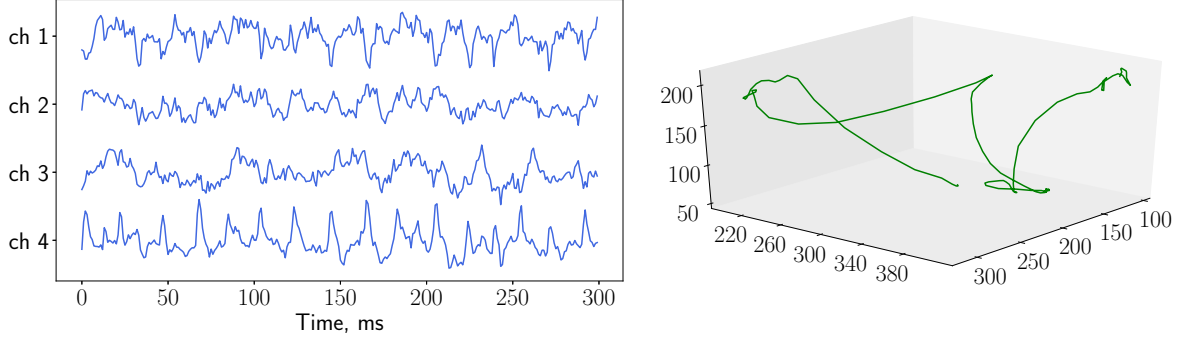


Figure 2: Brain signals (left plot) and 3D hand coordinates (right plot)

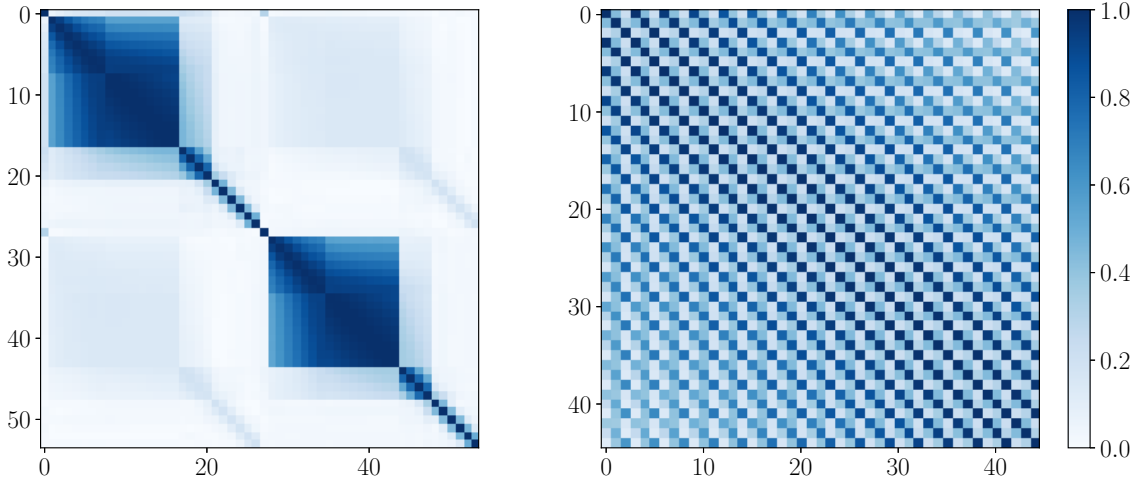


Figure 3: Correlation matrices for \mathbf{X} and \mathbf{Y}

We apply the QPFS algorithm with the SymImp strategy for different values of α_3 according to formula (10). The dependencies between target importances \mathbf{z}_y with respect to α_3 for different values of k are shown in Figure 4. The targets importances are almost the same for the predicted wrist coordinates with only one timestamp $k = 1$, which reflects the independence between the x , y , and z coordinates. For $k = 2$ and $k = 3$, the importances of some targets become zero when α_3 increases. The vertical lines correspond to the optimal value of α_3 obtained by (9). The target importances \mathbf{z}_y for this value of α_3 are similar. Thus, the algorithm does not distinguish the targets for $k = 1, 2, 3$.

We compare the proposed strategies of the multivariate QPFS that are given in Table 1 for the ECoG dataset. First, we apply all the methods to obtain the feature importances. Then, we fit a linear model with an increasing number of features. For each method, the features are sorted by their obtained importances. We show how the described quality criteria change with the increasing feature set size. Figure 5 illustrates the results for the prediction of $k = 30$ timestamps. Here, the feature importance threshold τ is represented by colored ticks. These thresholds are larger for the proposed methods in comparison to

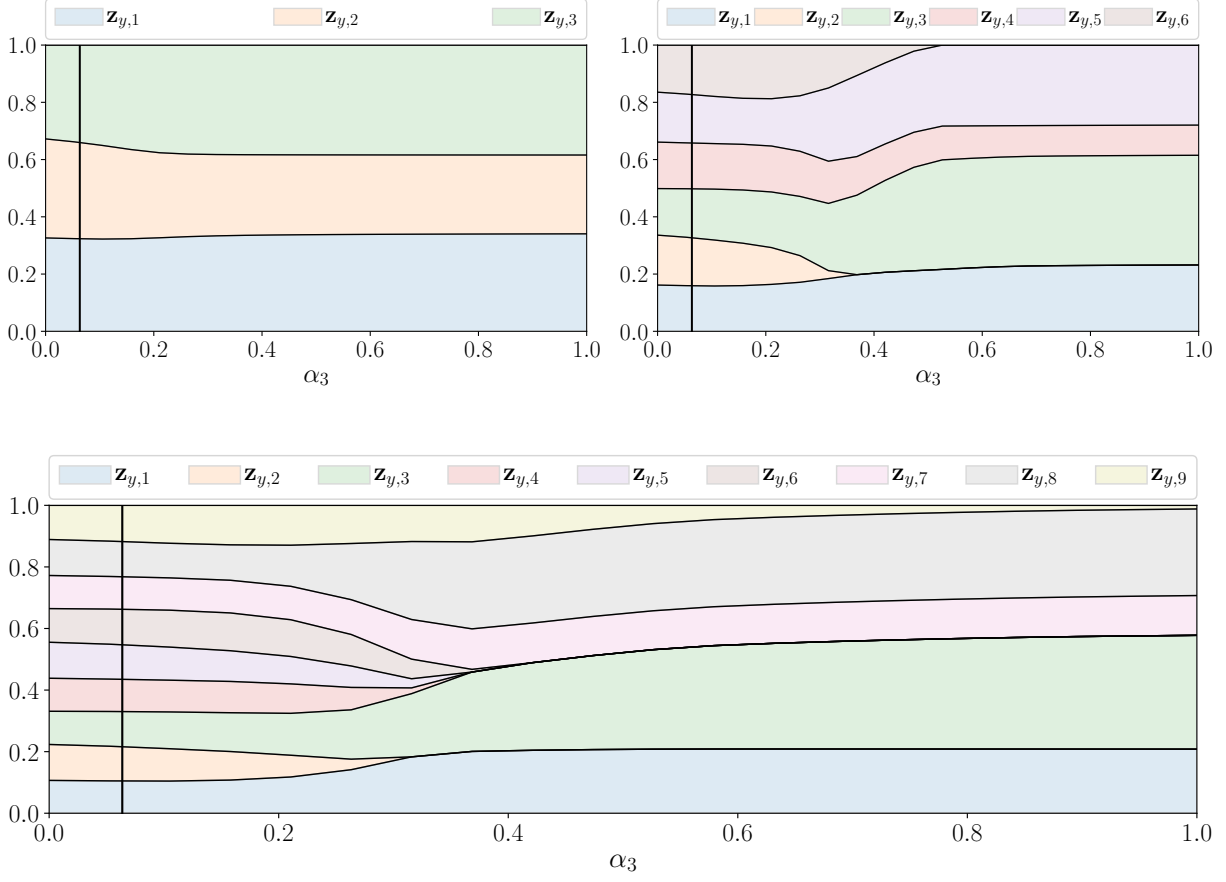


Figure 4: Target importances \mathbf{z}_y with respect to α_3 for QPFS with Symmetric Importance

the baseline RelAgg strategy. The SymImp strategy has the largest threshold, and it does not allow one to obtain a small feature subset. However, this strategy shows the best performance in terms of the sRMSE using the test data. The second performance is given by AsymImp. All proposed algorithms give smaller test errors compared to the RelAgg strategy. The Stability criteria is also increased for the proposed algorithms. Here, we consider the AsymImp strategy as the best in terms of the predictive quality and the size of selected feature subset.

To compare the structure of the selected feature subsets and investigate the stability of the selection procedure, we use the bootstrap approach. First, the bootstrap data are generated. Then, we solve the feature selection problem for each pair of the design and target matrices. The obtained feature importances are compared. We calculate the average pairwise Spearman correlation coefficient and the ℓ_2 distance as the measures of the algorithms stability. Table 2 shows the average error, the size of the subset and the described statistics for each method. The error was calculated by fitting the linear model using the 50 features with the largest importances. AsymImp gives the least error on the test data. The size of the selected feature subsets is overestimated using the threshold

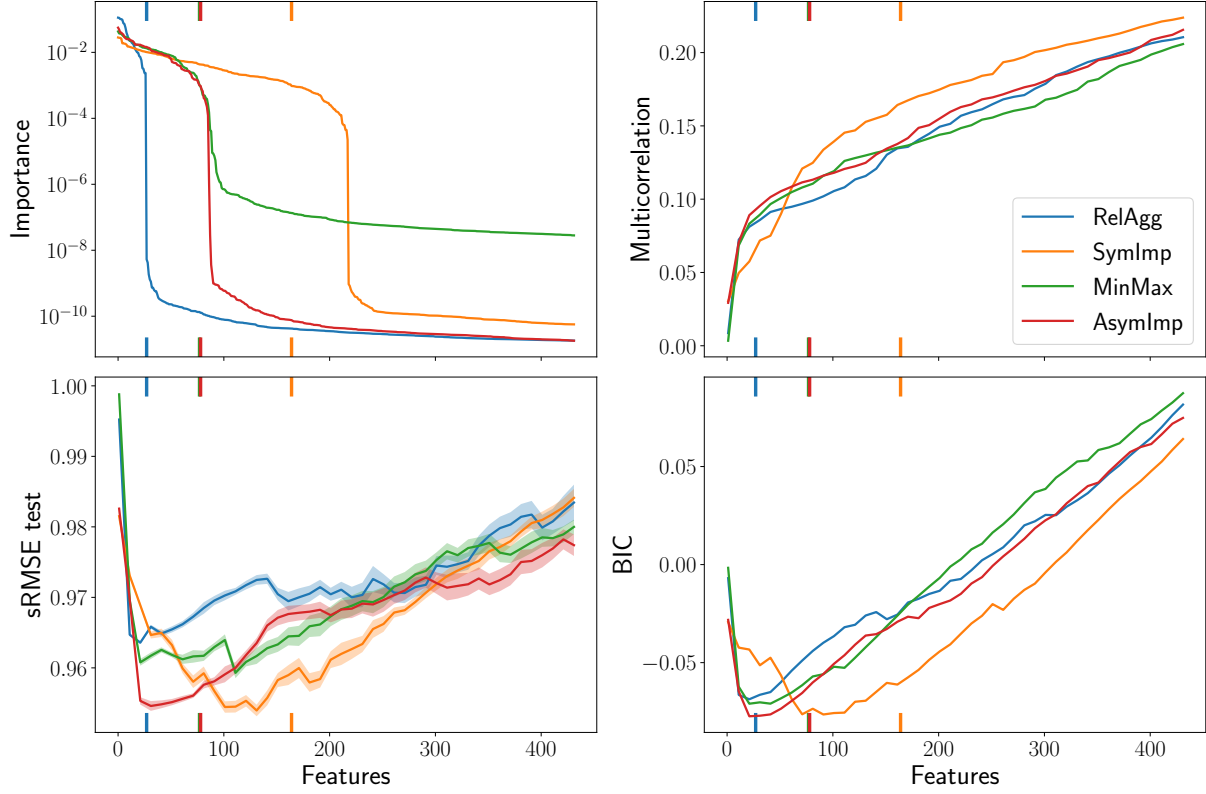


Figure 5: Feature selection algorithms evaluation for the ECoG data and the prediction of $k = 30$ times-tamps

$\tau = 10^{-4}$. The value of τ could be cross-validated to get the optimal threshold and feature subset size.

We fit the PLS regression model to the data to compare the dimensionality reduction and feature selection. Figure 6 demonstrates the scaled RMSE on the train and test data with respect to the dimensionality of the latent space l . The test error reaches its minimum at $l = 11$. The PLS regression is a more flexible approach compared to the linear model built on the subset of features. It results in the smallest error, but the model is not sparse.

Figure 7 compares 3 models: the linear regression and the PLS regression built on 100 features given QPFS and the PLS regression with all features. We do not include the linear

Table 2: The stability of the selected feature subset

| | sRMSE | $\ \mathbf{a}\ _0$ | Spearman ρ | ℓ_2 dist |
|---------|-------------------|--------------------|-------------------|-------------------|
| RelAgg | 0.965 ± 0.002 | 26.8 ± 3.8 | 0.915 ± 0.016 | 0.145 ± 0.018 |
| SymImp | 0.961 ± 0.001 | 224.4 ± 9.0 | 0.910 ± 0.017 | 0.025 ± 0.002 |
| MinMax | 0.961 ± 0.002 | 101.0 ± 2.1 | 0.932 ± 0.009 | 0.059 ± 0.004 |
| AsymImp | 0.955 ± 0.001 | 85.8 ± 10.2 | 0.926 ± 0.011 | 0.078 ± 0.007 |

regression with all features because its results are close to the constant prediction. It also provides the result for the Lasso and Elastic Net algorithms that are widely used for feature selection. We use the AsymImp strategy for QPFS in this experiment. The number of PLS latent dimension is $l = 15$. Here, the PLS regression is significantly better than the linear regression with the QPFS features. It means that the latter model is not flexible enough. However, the best result is by the PLS regression model combined with the QPFS features. This model is sparse since it uses only 100 QPFS features. The ability of the PLS model to find the optimal latent data representation improves the model performance.

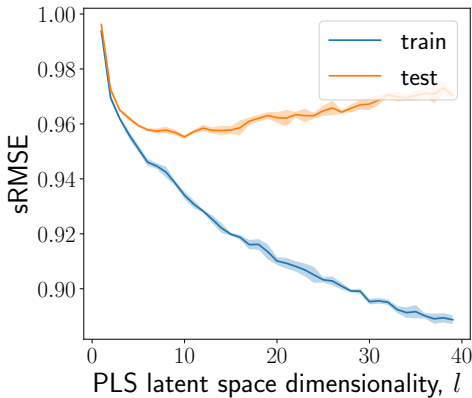


Figure 6: Test scaled RMSE for PLS regression models

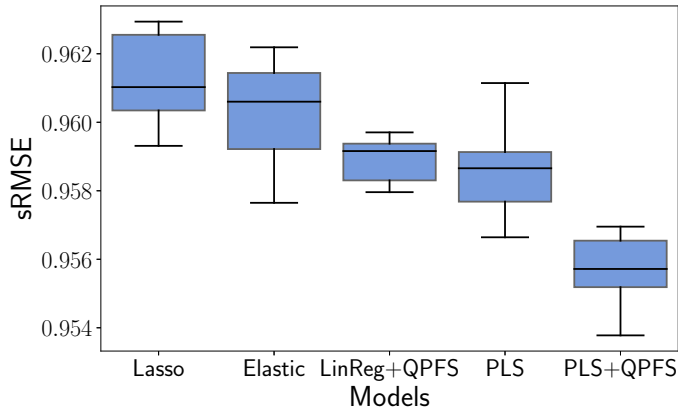


Figure 7: sRMSE box plots for different models

5. Conclusion

The study investigates the problem of signal decoding in which the data are highly redundant. To build a stable, adequate model, we reduced the dimensionality of the problem using the dependencies in both the input and target spaces. The PLS regression is considered as a linear model for dimensionality reduction. The quadratic programming approach is investigated as a feature selection algorithm. The algorithm solves feature selection in a single quadratic programming optimization problem. The multivariate extensions for the QPFS algorithms are proposed. The resulting feature subset includes noncorrelated features that are relevant to the most difficult targets.

The computational experiments were carried out using the ECoG data. The resulting model predicts the limb position of an exoskeleton using brain signals. The proposed algorithms outperform the baseline algorithm and reduce the problem dimension significantly. The combination of feature selection for sparsifying the model and the dimensionality reduction for increasing the model stability give the best result.

References

- [1] S. Karimi, M. Farrokhnia, Leukemia and small round blue-cell tumor cancer detection using microarray gene expression data set: Combining data dimension reduction and

- variable selection technique, *Chemometrics and Intelligent Laboratory Systems* 139 (2014) 6–14.
- [2] Y.-W. Lin, B.-C. Deng, Q.-S. Xu, Y.-H. Yun, Y.-Z. Liang, The equivalence of partial least squares and principal component regression in the sufficient dimension reduction framework, *Chemometrics and Intelligent Laboratory Systems* 150 (2016) 58–64.
 - [3] A. Eliseyev, T. Aksenova, Stable and artifact-resistant decoding of 3d hand trajectories from ecog signals using the generalized additive model, *Journal of neural engineering* 11 (6) (2014) 066005.
 - [4] A. Eliseyev, C. Moro, J. Faber, A. Wyss, N. Torres, C. Mestais, A. L. Benabid, T. Aksenova, L1-penalized n-way pls for subset of electrodes selection in bci experiments, *Journal of neural engineering* 9 (4) (2012) 045010.
 - [5] H. Chun, S. Keleş, Sparse partial least squares regression for simultaneous dimension reduction and variable selection, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 72 (1) (2010) 3–25.
 - [6] T. Mehmood, K. H. Liland, L. Snipen, S. Sæbø, A review of variable selection methods in partial least squares regression, *Chemometrics and Intelligent Laboratory Systems* 118 (2012) 62–69.
 - [7] A. M. Katrutsa, V. V. Strijov, Stress test procedure for feature selection algorithms, *Chemometrics and Intelligent Laboratory Systems* 142 (2015) 172–183.
 - [8] J. Li, K. Cheng, S. Wang, F. Morstatter, R. P. Trevino, J. Tang, H. Liu, Feature selection: A data perspective, *ACM Computing Surveys (CSUR)* 50 (6) (2017) 94.
 - [9] J. Lauzon-Gauthier, P. Manolescu, C. Duchesne, The sequential multi-block pls algorithm (smb-pls): Comparison of performance and interpretability, *Chemometrics and Intelligent Laboratory Systems*.
 - [10] S. Engel, T. Aksenova, A. Eliseyev, Kernel-based npls for continuous trajectory decoding from ecog data for bci applications, in: *International Conference on Latent Variable Analysis and Signal Separation*, Springer, 2017, pp. 417–426.
 - [11] A. Biancolillo, T. Næs, R. Bro, I. Måge, Extension of so-pls to multi-way arrays: So-n-pls, *Chemometrics and Intelligent Laboratory Systems* 164 (2017) 113–126.
 - [12] D. Hervás, J. Prats-Montalbán, A. Lahoz, A. Ferrer, Sparse n-way partial least squares with r package snpls, *Chemometrics and Intelligent Laboratory Systems* 179 (2018) 54 – 63.
 - [13] R. Rosipal, N. Kramer, Overview and Recent Advances in Partial Least Squares, C. Saunders et al. (Eds.): *SLSFS 2005*, LNCS 3940 (2006) 34–51doi:10.1007/11752790_2.

- [14] R. Rosipal, Nonlinear partial least squares an overview, in: Chemoinformatics and advanced machine learning perspectives: complex computational methods and collaborative techniques, IGI Global, 2011, pp. 169–189.
- [15] C. Ding, H. Peng, Minimum redundancy feature selection from microarray gene expression data, *Journal of bioinformatics and computational biology* 3 (02) (2005) 185–205.
- [16] M. Yamada, W. Jitkrittum, L. Sigal, E. P. Xing, M. Sugiyama, High-dimensional feature selection by feature-wise kernelized lasso, *Neural computation* 26 (1) (2014) 185–207.
- [17] I. Rodriguez-Lujan, R. Huerta, C. Elkan, C. S. Cruz, Quadratic programming feature selection, *Journal of Machine Learning Research* 11 (Apr) (2010) 1491–1516.
- [18] A. Katrutsa, V. Strijov, Comprehensive study of feature selection methods to solve multicollinearity problem according to evaluation criteria, *Expert Systems with Applications* 76 (2017) 1–11.
- [19] A. Motrenko, V. Strijov, Multi-way feature selection for ecog-based brain-computer interface, *Expert Systems with Applications*.
- [20] Y. Prasad, K. Biswas, P. Singla, Scaling-up quadratic programming feature selection., in: *AAAI (Late-Breaking Developments)*, 2013.
- [21] K. Shimoda, Y. Nagasaka, Z. C. Chao, N. Fujii, Decoding continuous three-dimensional hand trajectories from epidural electrocorticographic signals in japanese macaques, *Journal of neural engineering* 9 (3) (2012) 036015.
- [22] Z. C. Chao, Y. Nagasaka, N. Fujii, Long-term asynchronous decoding of arm motion using electrocorticographic signals in monkey, *Frontiers in neuroengineering* 3 (2010) 3.
- [23] A. Eliseyev, T. Aksenova, Penalized multi-way partial least squares for smooth trajectory decoding from electrocorticographic (ecog) recording, *PloS one* 11 (5) (2016) e0154878.