

Multivariate Quadratic Programming Feature Selection for ECoG signal decoding

R. V. Isachenko, V. V. Strijov

Abstract: The paper is devoted to the problem of signal decoding for Brain Computer Interface. The goal is to build a model which predicts a limb position by brain signals. The challenge of the investigation is redundancy in data description. High correlation in measurements leads to correlation in both input and target spaces. To overcome multi-correlation in feature representation, feature selection are used. However, the majority of feature selection methods ignore the dependencies in the target space. The authors suggest a novel approach to feature selection in multivariate regression. To take into account the correlations in the target matrix, the proposed approach extend the ideas of the Quadratic Programming Feature Selection algorithm. The overall algorithm select non-correlated features, which are relevant to the targets. The computational experiment shows performance of the proposed algorithms in the ECoG data.

Keywords: multivariate regression, quadratic programming feature selection, multi-correlation

1 Introduction

The paper investigates the problem of signal decoding for Brain Computer Interface (BCI) [1]. The BCI aims to develop systems that help people with a severe motor control disability to recover mobility. The minimally-invasive implant records cortical signals and the model decodes them on real time to predict the limb coordinates for an exoskeleton [2, 3]. The subject placed inside the exoskeleton can drive it by imagining movements as if they were making the movement by themselves.

The challenge is redundancy in initial data description. The features are highly multi-correlated. The correlation comes from spatial nature of the data. The brain sensors are close to each other. It leads to the redundant measurements. In this case the final model is unstable. In addition, the redundant data description requires redundant computations which lead to real-time delay. To overcome this problem feature selection methods are used [4, 5].

One of the approach to the feature selection is to maximize feature relevances and minimize pairwise feature redundancy. This approach was proposed in [6]. The Quadratic Programming Feature Selection (QPFS) [7, 8] algorithm introduces two functions: Sim and Rel. The Sim measures the redundancy between features, the Rel contains relevances between each feature and the target vector. We want to minimize the function Sim and maximize the Rel simultaneously. QPFS offers the explicit way to construct the functions Sim and Rel. The method minimizes the following function

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

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. The function (1) penalizes the dependent features by the function Sim and encourages features relevant to the target by the function Rel. The parameter α controls the trade-off between the functions Sim and the Rel. To measure similarity the authors use the absolute value of sample correlation coefficient or sample mutual information coefficient between pairs of features for the function Sim, and between features and the target vector for the function Rel.

We consider the multivariate problem, where the dependent variable is a vector. It refers to the prediction of limb position for not just one moment, but for some period of time. The subsequent hand positions are correlated. It leads to correlations in the model output. In this situation feature selection algorithms do not take into account these dependencies. Hence, the selected feature subset is not optimal. We propose methods to take into account the dependencies in both input and output spaces. It allows to get the stable model with fewer variables.

The experiment was carried out in the ECoG data from the NeuroTycho project [9]. The proposed algorithms outperforms the original methods.

2 Problem statement

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 there is a linear dependence

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

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

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

The columns $\boldsymbol{\chi}_j$ of the matrix \mathbf{X} respond to object features.

62 The optimal parameters are determined by minimization of an error function. Define
 63 the quadratic loss function:

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

64 The solution of (3) is given by

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

65 The linear dependent columns of the matrix \mathbf{X} leads to an instable solution for the
 66 optimization problem (3). If there is a vector $\alpha \neq \mathbf{0}_n$ such that $\mathbf{X}\alpha = \mathbf{0}_m$, then adding
 67 the vector α to any column of the matrix Θ does not change the value of the loss function
 68 $\mathcal{L}(\Theta|\mathbf{X}, \mathbf{Y})$. In this case the matrix $\mathbf{X}^T \mathbf{X}$ is not invertible. To avoid the strong linear
 69 dependence, feature selection and dimensionality reduction techniques are used.

70 3 Feature selection

71 The feature selection goal is to find the boolean vector $\mathbf{a} = \{0, 1\}^n$, which components
 72 indicates whether the feature are selected. To obtain the optimal vector \mathbf{a} among all
 73 possible $2^n - 1$ options, introduce the feature selection error function

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

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

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

$$\mathcal{L}(\Theta_{\mathbf{a}}|\mathbf{X}_{\mathbf{a}}, \mathbf{Y}) = \left\| \mathbf{Y} - \mathbf{X}_{\mathbf{a}} \Theta_{\mathbf{a}}^T \right\|_2^2 \rightarrow \min_{\Theta_{\mathbf{a}}}, \quad (5)$$

78 where the subscript \mathbf{a} indicates the submatrix with the columns for which components of \mathbf{a}
 79 equal 1.

80 3.1 Quadratic Programming Feature Selection

The QPFS algorithm selects non-correlated features, which are relevant to the target vector ν for the linear regression problem with $r = 1$

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

81 The domain of QPFS function (1) is $[0, 1]^n$. It refers to the relaxed error function $S(\mathbf{a}|\mathbf{X}, \nu)$,
 82 which has boolean domain $\{0, 1\}^n$. The link between the boolean vector \mathbf{a} from (4) and

83 the QPFS vector \mathbf{z} from (1) is given by:

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

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

$$\alpha = \frac{\overline{\mathbf{Q}}}{\overline{\mathbf{Q}} + \overline{\mathbf{b}}},$$

84 where $\overline{\mathbf{Q}}$, $\overline{\mathbf{b}}$ are the mean values of \mathbf{Q} and \mathbf{b} respectively. We use the absolute value of
85 sample correlation coefficient as similarity measure:

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

86 The other choices to define \mathbf{Q} and \mathbf{b} , such as mutual information and normalized feature
87 significance, are considered in [8].

88 The problem (1) is convex if the matrix \mathbf{Q} is positive semidefinite. In general it is not
89 always true. To satisfy this condition, the matrix \mathbf{Q} spectrum is shifted and the matrix \mathbf{Q}
90 is replaced by $\mathbf{Q} - \lambda_{\min} \mathbf{I}$, where λ_{\min} is a \mathbf{Q} minimal eigenvalue.

91 3.2 Multivariate QPFS

92 We are aimed to propose the algorithms which suitable for feature selection in multivariate
93 case. If the target space is multidimensional it prone to redundancy and correlations
94 between the targets. In this section we consider the algorithms that take into account the
95 probable dependencies in both input and target spaces.

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

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

96 The drawback of this approach is that it does not use the dependencies in the columns
97 of the matrix \mathbf{Y} . Observe the following example:

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

98 We have three features and r targets, where first $r - 1$ target are identical. The pairwise
99 features similarities are given by the matrix \mathbf{Q} . The matrix \mathbf{B} entries show pairwise features

100 relevances to the targets. The vector \mathbf{b} is obtained by summation of the matrix \mathbf{B} over
 101 columns.

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

102

103 We would like to select only two features. For such configuration the best feature sub-
 104 set is $[\chi_1, \chi_2]$. The feature χ_2 predicts the second target ν_2 and the combination of
 105 features χ_1, χ_2 predicts the first component. The QPFS algorithm for $r = 2$ gives
 106 the solution $\mathbf{z} = [0.37, 0.61, 0.02]$. It coincides with our knowledge. However, if we
 107 add the collinear columns to the matrix \mathbf{Y} and increase r to 5, the QPFS solution will
 108 be $\mathbf{z} = [0.40, 0.17, 0.43]$. Here we lose the relevant feature χ_2 and select the redundant
 109 feature χ_3 .

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

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

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

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

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

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

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

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

$$\alpha_1 = \frac{\overline{\mathbf{Q}_y} \overline{\mathbf{B}}}{\overline{\mathbf{Q}_y} \overline{\mathbf{B}} + \overline{\mathbf{Q}_x} \overline{\mathbf{Q}_y} + \overline{\mathbf{Q}_x} \overline{\mathbf{B}}}, \quad \alpha_2 = \frac{\overline{\mathbf{Q}_x} \overline{\mathbf{Q}_y}}{\overline{\mathbf{Q}_y} \overline{\mathbf{B}} + \overline{\mathbf{Q}_x} \overline{\mathbf{Q}_y} + \overline{\mathbf{Q}_x} \overline{\mathbf{B}}}, \quad \alpha_3 = \frac{\overline{\mathbf{Q}_x} \overline{\mathbf{B}}}{\overline{\mathbf{Q}_y} \overline{\mathbf{B}} + \overline{\mathbf{Q}_x} \overline{\mathbf{Q}_y} + \overline{\mathbf{Q}_x} \overline{\mathbf{B}}},$$

120 where $\overline{\mathbf{Q}_x}$, $\overline{\mathbf{B}}$, $\overline{\mathbf{Q}_y}$ are mean values of \mathbf{Q}_x , \mathbf{B} , and \mathbf{Q}_y , respectively.

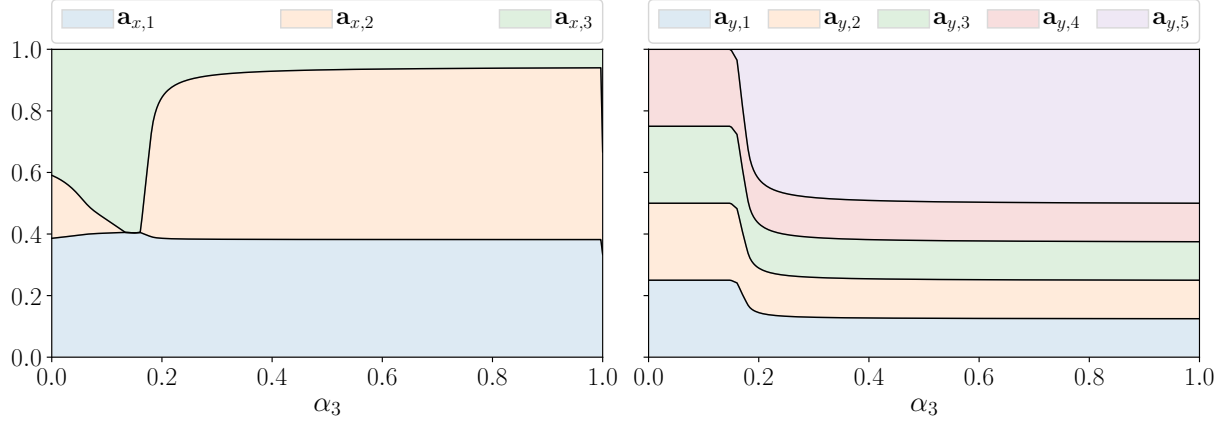


Figure 1: Feature importances \mathbf{z}_x and \mathbf{z}_y with respect to the α_3 coefficient

Proof. The desired values of α_1 , α_2 , and α_3 are given by solution 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 $\overline{\mathbf{Q}}_x$, $\overline{\mathbf{B}}$, $\overline{\mathbf{Q}}_y$ 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 $\text{Rel}(\mathbf{X}, \mathbf{Y})$. \square

To investigate the impact of the term $\text{Sim}(\mathbf{Y})$ on the 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 :

$$\alpha_1 = \frac{(1 - \alpha_3)\overline{\mathbf{B}}}{\overline{\mathbf{Q}}_x + \overline{\mathbf{B}}}; \quad \alpha_2 = \frac{(1 - \alpha_3)\overline{\mathbf{Q}}_x}{\overline{\mathbf{Q}}_x + \overline{\mathbf{B}}}; \quad \alpha_3 \in [0, 1]. \quad (9)$$

We apply the proposed algorithm to the discussed example (7). The given matrix \mathbf{Q} corresponds to the matrix \mathbf{Q}_x . We additionally define the matrix \mathbf{Q}_y by setting $\text{corr}(\nu_1, \nu_2) = 0.2$ 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 coefficient. If α_3 is small, the impact of all targets are almost equal and the feature χ_3 dominates the feature χ_2 . When α_3 becomes larger than 0.2, the importance $(\mathbf{z}_y)_5$ of the target ϕ_5 grows up along with the importance of the feature χ_2 .

Minimax QPFS (MinMax and MaxMin). The function (8) is symmetric with respect to \mathbf{z}_x and \mathbf{z}_y . It penalizes features that are correlated and do not relevant to targets. At the same time it penalizes targets that are correlated and are not sufficiently explained by the features. It leads to small importances for targets which are difficult to predict by features and large importances for targets which are strongly correlated with features. It contradicts with the intuition. Our goal is to predict all targets, especially which are difficult to explain, by selected relevant and non-correlated features. We express this into

two related problems.

$$\alpha_1 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x}_{\text{Sim}(\mathbf{X})} - \alpha_2 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y}_{\text{Rel}(\mathbf{X}, \mathbf{Y})} \rightarrow \min_{\substack{\mathbf{z}_x \geq \mathbf{0}_n, \\ \mathbf{1}_n^\top \mathbf{z}_x = 1}}; \quad (10)$$

$$\alpha_3 \cdot \underbrace{\mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y}_{\text{Sim}(\mathbf{Y})} + \alpha_2 \cdot \underbrace{\mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y}_{\text{Rel}(\mathbf{X}, \mathbf{Y})} \rightarrow \min_{\substack{\mathbf{z}_y \geq \mathbf{0}_r, \\ \mathbf{1}_r^\top \mathbf{z}_y = 1}}. \quad (11)$$

131 The difference is in the term Rel. In feature space the non-relevant components should
 132 have smaller scores. Meanwhile, the targets that are not relevant to the features should
 133 have larger scores. The problems (10) and (11) are merged into the joint min-max or
 134 max-min formulation

$$\min_{\substack{\mathbf{z}_x \geq \mathbf{0}_n \\ \mathbf{1}_n^\top \mathbf{z}_x = 1}} \max_{\substack{\mathbf{z}_y \geq \mathbf{0}_r \\ \mathbf{1}_r^\top \mathbf{z}_y = 1}} f(\mathbf{z}_x, \mathbf{z}_y), \quad \left(\text{or } \max_{\substack{\mathbf{z}_y \geq \mathbf{0}_r \\ \mathbf{1}_r^\top \mathbf{z}_y = 1}} \min_{\substack{\mathbf{z}_x \geq \mathbf{0}_n \\ \mathbf{1}_n^\top \mathbf{z}_x = 1}} f(\mathbf{z}_x, \mathbf{z}_y) \right), \quad (12)$$

135 where

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

136 **Theorem 1.** For positive definite matrices \mathbf{Q}_x and \mathbf{Q}_y the max-min and min-max prob-
 137 lems (12) have the same optimal value.

Proof. Denote

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

138 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
 139 function. If \mathbf{Q}_x and \mathbf{Q}_y are positive definite matrices, the function f is convex-concave,
 140 i.e. $f(\cdot, \mathbf{z}_y) : \mathbb{C}^n \rightarrow \mathbb{R}$ is convex for fixed \mathbf{z}_y , and $f(\mathbf{z}_x, \cdot) : \mathbb{C}^r \rightarrow \mathbb{R}$ is concave for fixed \mathbf{z}_x .
 141 In this case Neumann's minimax theorem states

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

142

□

143 To solve the min-max problem (12), fix some $\mathbf{z}_x \in \mathbb{C}^n$. For fixed vector \mathbf{z}_x we solve the
 144 problem

$$\max_{\mathbf{z}_y \in \mathbb{C}^r} f(\mathbf{z}_x, \mathbf{z}_y) = \max_{\substack{\mathbf{z}_y \geq \mathbf{0}_r \\ \mathbf{1}_r^\top \mathbf{z}_y = 1}} [\alpha_1 \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha_2 \cdot \mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y - \alpha_3 \cdot \mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y]. \quad (13)$$

145 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 $\boldsymbol{\mu}$, corresponding to the inequality constraints $\mathbf{z}_y \geq \mathbf{0}_r$, are restricted to be non-negative. The dual problem is

$$\min_{\lambda, \boldsymbol{\mu} \geq \mathbf{0}_r} g(\mathbf{z}_x, \lambda, \boldsymbol{\mu}) = \min_{\lambda, \boldsymbol{\mu} \geq \mathbf{0}_r} \left[\max_{\mathbf{z}_y \in \mathbb{R}^r} L(\mathbf{z}_x, \mathbf{z}_y, \lambda, \boldsymbol{\mu}) \right]. \quad (14)$$

The strong duality holds for (13). Therefore, the optimal value for (13) equals the optimal value for (14). It allows to solve the problem

$$\min_{\mathbf{z}_x \in \mathbb{C}^n, \lambda, \boldsymbol{\mu} \geq \mathbf{0}_r} g(\mathbf{z}_y, \lambda, \boldsymbol{\mu}) \quad (15)$$

instead of (12).

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

$$\mathbf{z}_y = \frac{1}{2\alpha_3} \mathbf{Q}_y^{-1} (-\alpha_2 \cdot \mathbf{B}^\top \mathbf{z}_x + \lambda \cdot \mathbf{1}_r + \boldsymbol{\mu}). \quad (16)$$

The dual function is equal to

$$\begin{aligned} g(\mathbf{z}_x, \lambda, \boldsymbol{\mu}) = \max_{\mathbf{z}_y \in \mathbb{R}^r} L(\mathbf{z}_x, \mathbf{z}_y, \lambda, \boldsymbol{\mu}) &= \mathbf{z}_x^\top \left(-\frac{\alpha_2^2}{4\alpha_3} \cdot \mathbf{B} \mathbf{Q}_y^{-1} \mathbf{B}^\top - \alpha_1 \cdot \mathbf{Q}_x \right) \mathbf{z}_x \\ &\quad - \frac{1}{4\alpha_3} \lambda^2 \cdot \mathbf{1}_r^\top \mathbf{Q}_y^{-1} \mathbf{1}_r - \frac{1}{4\alpha_3} \cdot \boldsymbol{\mu}^\top \mathbf{Q}_y^{-1} \boldsymbol{\mu} + \frac{\alpha_2}{2\alpha_3} \lambda \cdot \mathbf{1}_r^\top \mathbf{Q}_y^{-1} \mathbf{B}^\top \mathbf{z}_x \\ &\quad - \frac{1}{2\alpha_3} \lambda \cdot \mathbf{1}_r^\top \mathbf{Q}_y^{-1} \boldsymbol{\mu} + \frac{\alpha_2}{2\alpha_3} \cdot \boldsymbol{\mu}^\top \mathbf{Q}_y^{-1} \mathbf{B}^\top \mathbf{z}_x + \lambda. \end{aligned} \quad (17)$$

It brings to the quadratic problem (15) with $n + r + 1$ variables.

Minimax Relevances (MaxRel). The problem (15) is not convex. If we shift the spectrum for the matrix of quadratic form (17), the optimality is lost. To overcome this problem, we drop the term $\text{Sim}(\mathbf{Y})$.

$$\min_{\substack{\mathbf{z}_x \geq \mathbf{0}_n \\ \mathbf{1}_n^\top \mathbf{z}_x = 1}} \max_{\substack{\mathbf{z}_y \geq \mathbf{0}_r \\ \mathbf{1}_r^\top \mathbf{z}_y = 1}} [(1 - \alpha) \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha \cdot \mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y]. \quad (18)$$

The Lagrangian for the problem (18) with the fixed vector \mathbf{z}_x is

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

Setting the gradient of the Langrangian $\nabla_{\mathbf{z}_y} L(\mathbf{z}_x, \mathbf{z}_y, \lambda, \boldsymbol{\mu})$ to zero, we obtain:

$$\alpha \cdot \mathbf{B}^\top \mathbf{z}_x = \lambda \cdot \mathbf{1}_r + \boldsymbol{\mu}.$$

The dual function is equal to

$$g(\mathbf{z}_x, \lambda, \boldsymbol{\mu}) = \begin{cases} (1 - \alpha) \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \lambda, & \alpha \cdot \mathbf{B}^\top \mathbf{z}_x = \lambda \cdot \mathbf{1}_r + \boldsymbol{\mu}; \\ +\infty, & \text{otherwise.} \end{cases} \quad (19)$$

In this case the feature scores are the solution of (15).

Algorithm	Strategy	Error function $S(\mathbf{a} \mathbf{X}, \mathbf{Y})$
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	$\min [\text{Sim}(\mathbf{X}) - \text{Rel}(\mathbf{X}, \mathbf{Y})]$ $\max [\text{Rel}(\mathbf{X}, \mathbf{Y}) + \text{Sim}(\mathbf{Y})]$	$\min_{\mathbf{z}_x} \max_{\mathbf{z}_y} [\alpha_1 \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha_2 \cdot \mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y - \alpha_3 \cdot \mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y]$
MaxMin	$\max [\text{Rel}(\mathbf{X}, \mathbf{Y}) + \text{Sim}(\mathbf{Y})]$ $\min [\text{Sim}(\mathbf{X}) - \text{Rel}(\mathbf{X}, \mathbf{Y})]$	$\max_{\mathbf{z}_y} \min_{\mathbf{z}_x} [\alpha_1 \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha_2 \cdot \mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y - \alpha_3 \cdot \mathbf{z}_y^\top \mathbf{Q}_y \mathbf{z}_y]$
MaxRel	$\min [\text{Sim}(\mathbf{X}) - \text{Rel}(\mathbf{X}, \mathbf{Y})]$ $\max [\text{Rel}(\mathbf{X}, \mathbf{Y})]$	$\min_{\mathbf{z}_x} \max_{\mathbf{z}_y} [(1 - \alpha) \cdot \mathbf{z}_x^\top \mathbf{Q}_x \mathbf{z}_x - \alpha \cdot \mathbf{z}_x^\top \mathbf{B} \mathbf{z}_y]$

Table 1: Overview of proposed multivariate QPFS algorithms

Proposition 2. For the case $r = 1$ the proposed functions (8), (12), and (18) coincide with the original QPFS algorithm (1).

Proof. If r is equal to 1, then $\mathbf{Q}_y = q_y$ is a scalar, $\mathbf{z}_y = 1$, $\mathbf{B} = \mathbf{b}$. It reduces the problems (8), (12), and (18) 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 $\alpha = \frac{\alpha_2}{\alpha_1 + \alpha_2}$ brings to the original QPFS problem (1). \square

To summarize all proposed strategies for multivariate feature selection, Table 1 shows the core ideas and error functions for each method.

4 Metrics

To evaluate the selected subset \mathcal{A} we introduce criteria that estimates the quality of feature selection procedure. We measure multicorrelation by mean value of multiple 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. The bigger R^2 means the better feature subset we have.

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

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

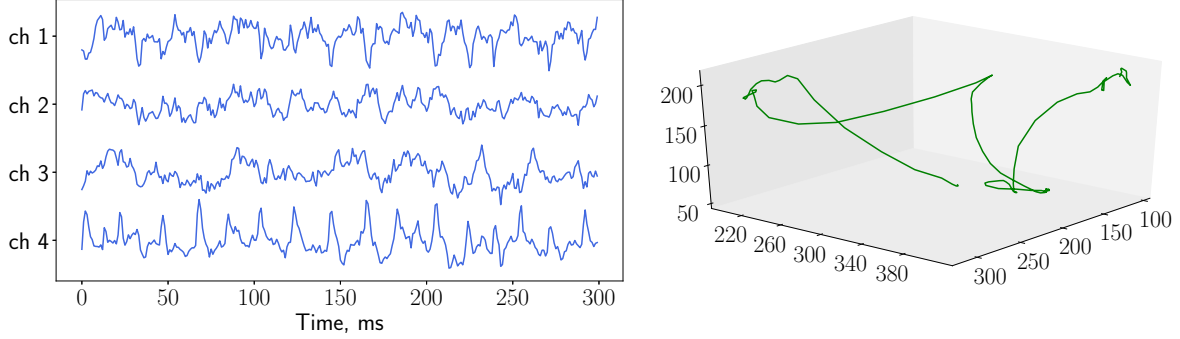


Figure 2: Brain signals and the corresponding hand position

175 The Root Mean Squared Error (RMSE) shows the quality of the model prediction. We
 176 estimate RMSE on train and test data.

$$\text{RMSE}(\mathbf{Y}, \hat{\mathbf{Y}}_{\mathcal{A}}) = \sqrt{\text{MSE}(\mathbf{Y}, \hat{\mathbf{Y}}_{\mathcal{A}})} = \|\mathbf{Y} - \hat{\mathbf{Y}}_{\mathcal{A}}\|_2, \quad \text{where} \quad \hat{\mathbf{Y}}_{\mathcal{A}} = \mathbf{X}_{\mathcal{A}} \boldsymbol{\Theta}_{\mathcal{A}}^{\top}.$$

177 Akaike Information Criteria (AIC) is a trade-off between prediction quality and the size of
 178 selected subset \mathcal{A} :

$$\text{AIC} = m \ln \left(\frac{\text{MSE}(\mathbf{Y}, \hat{\mathbf{Y}}_{\mathcal{A}})}{m} \right) + 2|\mathcal{A}|.$$

179 5 Experiment

180 We carried out computational experiment with ECoG data from the NeuroTycho project.
 181 The input data consists of brain voltage signals recorded from 32 channels. The goal is
 182 to predict 3D hand position in the next moments given the input signal. The example of
 183 input signals and the 3D wrist coordinates are shown in Figure 2. The initial voltage signals
 184 are transformed to the spatial-temporal representation using wavelet transformation. The
 185 procedure of extracting feature representation from the raw data are described in details
 186 in [10, 11]. Feature description at each time moment has dimension equals to 32 (channels)
 187 $\times 27$ (frequencies) = 864. Each object is the representation of local history time segment
 188 with duration $\Delta t = 1s$. The time step between objects is $\delta t = 0.05s$. The final matrices
 189 are $\mathbf{X} \in \mathbb{R}^{18900 \times 864}$ and $\mathbf{Y} \in \mathbb{R}^{18900 \times 3k}$, where k is a number of timestamps that we predict.
 190 We split our data into train and test parts with the ratio 0.67.

191 Figures 3 and 4 show the result of the QPFS algorithm, where we use the Relevance
 192 Aggregation strategy and $k = 1$. QPFS scores \mathbf{z}_x decrease sharply. Only about one hundred
 193 features have scores significantly greater than zero. The test error stops to decrease using
 194 more than this amount of features. It confirms that the initial data representation is highly
 195 redundant.

196 Figure 5 shows the dependencies in the matrices \mathbf{X} and \mathbf{Y} . Frequencies in the matrix \mathbf{X}
 197 are highly correlated. The correlations between axes are not significant in comparison with
 198 the correlations between consequent moments.

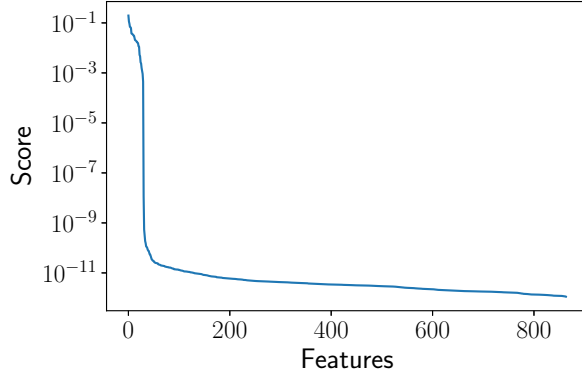


Figure 3: Sorted feature importances for the QPFS algorithm

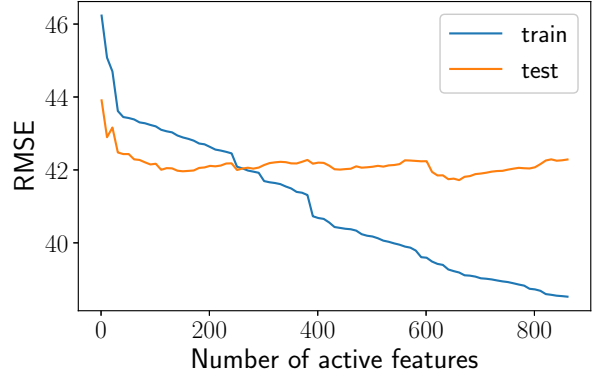


Figure 4: RMSE w.r.t. size of active set, features are ranked by QPFS algorithm

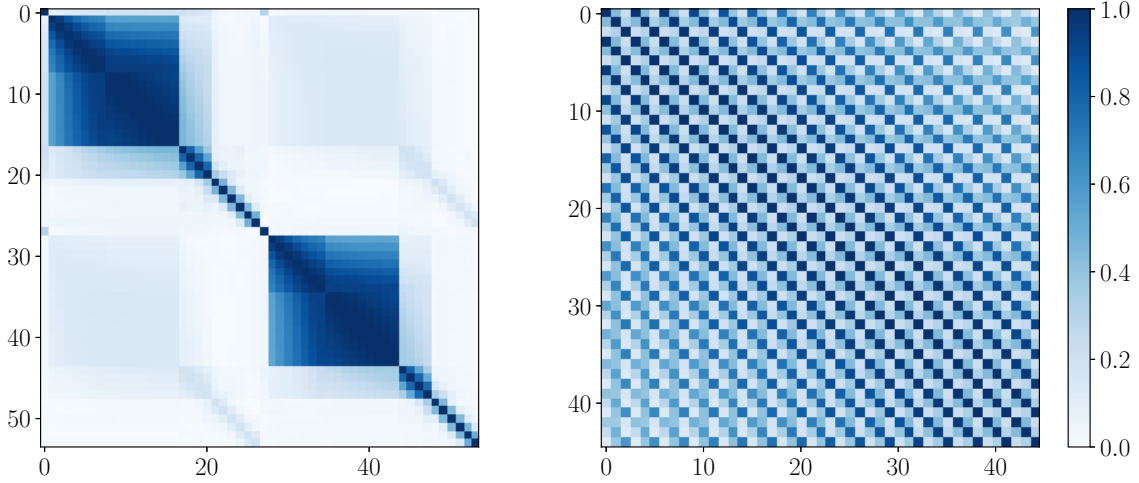


Figure 5: Correlation matrices for \mathbf{X} and \mathbf{Y}

We apply the QPFS algorithm with Relevance Aggregation strategy for different values of α_3 coefficient according to formulas (9). The dependence between target scores \mathbf{z}_y with respect to α_3 for different values of k is shown in Figure 6. If we predict wrist coordinates only for one timestamp $k = 1$, targets scores are almost the same. It tells about the independence between x , y , and z coordinates. For $k = 2$ and $k = 3$ the scores of some targets become zero when α_3 increases.

We compare the proposed strategies of multivariate QPFS that are given in Table 1 for the ECoG dataset. Firstly, we apply all methods to get feature scores. Then we fit linear model with increasing number of used features. For each method the features are sorted by the obtained scores. We show how the described metrics are changed with the increasing feature set size. Figure 7 illustrates the results for $k = 1$. Here all metrics values for ReIAgg, SymImp, and MaxRel are quite similar. However, MaxMin and MinMax algorithms

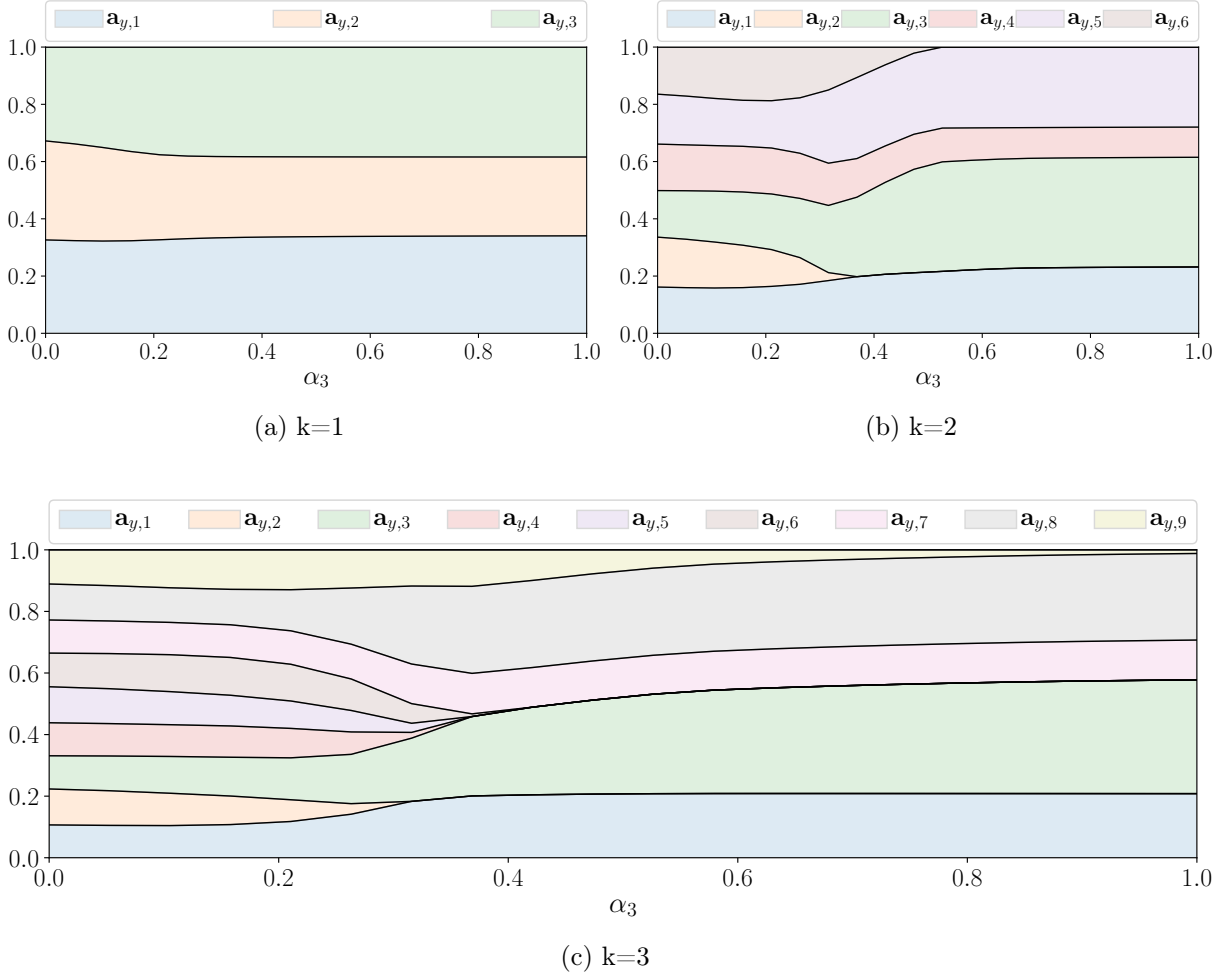


Figure 6: Target importances \mathbf{z}_y with respect to α_3 for QPFS with Relevance Aggregation

show worse performance. This behaviour is possibly due to unreasonable penalty on the matrix \mathbf{Y} .

The situation changes for predicting several timestamps. In this case the correlation in \mathbf{Y} matrix is crucial. Figure 8 shows algorithms performance for $k = 15$. The test error is minimal for SymImp, MinMax, and MaxRel strategies. The RelAgg strategy shows almost the worst error rate.

6 Conclusion

The paper investigates the problem of signal decoding in relation to Brain Computer Interface. To build a stable edequate model, the authors proposed multivariate feature selection strategies. The algorithms incorporates dependencies in both input and output spaces. The final model are more robust and effective. The computational experiments

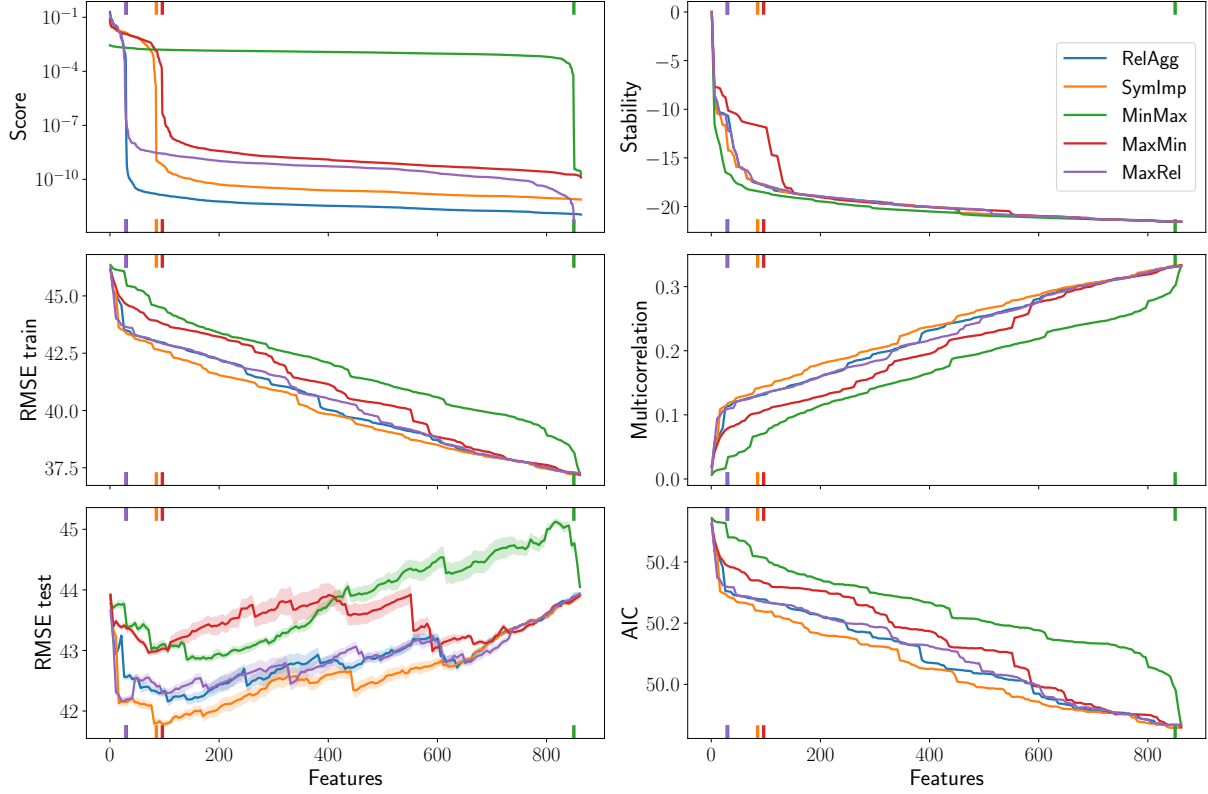


Figure 7: Metrics values for ECoG data with $k = 1$

shows that the proposed algorithms outperforms the baseline strategy by multiple metrics.

References

- [1] Thomas Costecalde, Tetiana Aksenova, Napoleon Torres-Martinez, Andriy Eliseyev, Corinne Mestais, Cecile Moro, and Alim Louis Benabid. A long-term bci study with ecog recordings in freely moving rats. *Neuromodulation: Technology at the Neural Interface*, 21(2):149–159, 2018.
- [2] Corinne S Mestais, Guillaume Charvet, Fabien Sauter-Starace, Michael Foerster, David Ratel, and Alim Louis Benabid. Wimage: Wireless 64-channel ecog recording

	RMSE	$\ \mathbf{a}\ _0$	Spearman ρ	ℓ_2 dist
RelAgg	41.9 ± 0.1	27.0 ± 2.4	0.941 ± 0.005	0.14 ± 0.02
SymImp	41.1 ± 0.1	198.6 ± 3.8	0.942 ± 0.008	0.03 ± 0.00
MinMax	41.4 ± 0.2	92.4 ± 8.2	0.933 ± 0.007	0.10 ± 0.01
MaxMin	41.7 ± 0.1	97.0 ± 4.4	0.950 ± 0.004	0.07 ± 0.01
MaxRel	41.7 ± 0.0	37.6 ± 1.6	0.893 ± 0.012	0.17 ± 0.02

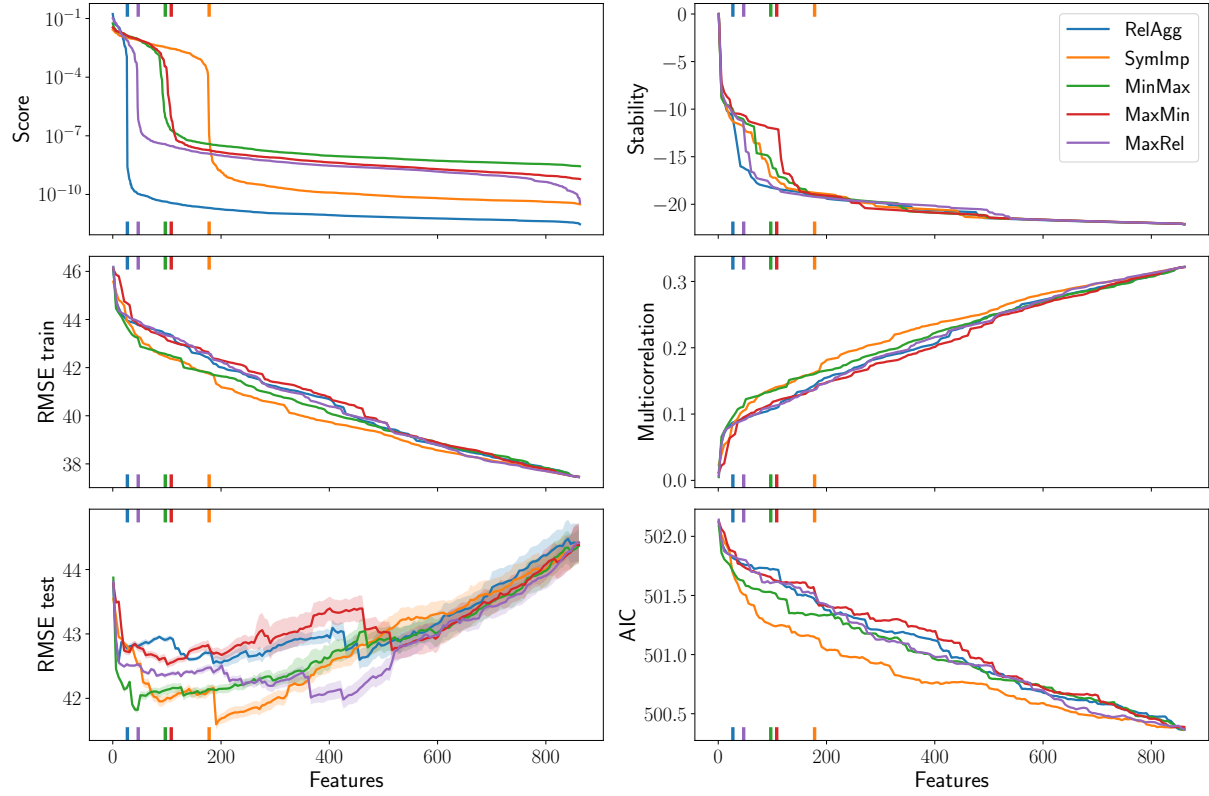


Figure 8: Metrics values for ECoG data with $k = 15$

implant for long term clinical applications. *IEEE transactions on neural systems and rehabilitation engineering*, 23(1):10–21, 2015.

- [3] Andrey Eliseyev, Corinne Mestais, Guillaume Charvet, Fabien Sauter, Neil Abroug, Nana Arizumi, Serpil Cokgungor, Thomas Costecalde, Michael Foerster, Louis Karczowski, et al. Clinatex® bci platform based on the ecog-recording implant wimagine® and the innovative signal-processing: preclinical results. In *Engineering in Medicine and Biology Society (EMBC), 2014 36th Annual International Conference of the IEEE*, pages 1222–1225. IEEE, 2014.

- [4] A. M. Katrutsa and V. V. Strijov. Stress test procedure for feature selection algorithms. *Chemometrics and Intelligent Laboratory Systems*, 142:172–183, 2015.

- [5] Jundong Li, Kewei Cheng, Suhang Wang, Fred Morstatter, Robert P Trevino, Jiliang Tang, and Huan Liu. Feature selection: A data perspective. *ACM Computing Surveys (CSUR)*, 50(6):94, 2017.

- [6] Chris Ding and Hanchuan Peng. Minimum redundancy feature selection from microarray gene expression data. *Journal of bioinformatics and computational biology*, 3(02):185–205, 2005.

- 246 [7] Irene Rodriguez-Lujan, Ramon Huerta, Charles Elkan, and Carlos Santa Cruz.
 247 Quadratic programming feature selection. *Journal of Machine Learning Research*,
 248 11(Apr):1491–1516, 2010.
- 249 [8] Alexandr Katrutsa and Vadim Strijov. Comprehensive study of feature selection meth-
 250 ods to solve multicollinearity problem according to evaluation criteria. *Expert Systems*
 251 *with Applications*, 76:1–11, 2017.
- 252 [9] Project tycho <http://neurotycho.org/food-tracking-task>.
- 253 [10] Zenas C Chao, Yasuo Nagasaka, and Naotaka Fujii. Long-term asynchronous decoding
 254 of arm motion using electrocorticographic signals in monkey. *Frontiers in neuroengi-*
 255 *neering*, 3:3, 2010.
- 256 [11] Andrey Eliseyev and Tetiana Aksenova. Penalized multi-way partial least squares for
 257 smooth trajectory decoding from electrocorticographic (ecog) recording. *PloS one*,
 258 11(5):e0154878, 2016.