

# 1 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

$$\mathbf{y} = \mathbf{\Theta}\mathbf{x} + \varepsilon \quad (1)$$

between the objects  $\mathbf{x}$  and the target variable  $\mathbf{y}$ , where  $\mathbf{\Theta} \in \mathbb{R}^{r \times n}$  is the matrix of model parameters,  $\varepsilon \in \mathbb{R}^r$  is the residual vector. The task is 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

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]^T = [\chi_1, \dots, \chi_n]; \quad \mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_m]^T = [\phi_1, \dots, \phi_r].$$

The columns  $\chi_j$  of the matrix  $\mathbf{X}$  respond to object features. The examples of how to construct the dataset for a particular application task are described in Section Computational experiment.

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

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

The solution of the problem (2) is given by

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

The linear dependent columns of the matrix  $\mathbf{X}$  leads to an instable solution for the optimization problem (2). If there is a vector  $\boldsymbol{\alpha} \neq 0$  such that  $\mathbf{X}\boldsymbol{\alpha} = 0$ , then adding the vector  $\boldsymbol{\alpha}$  to any column of the matrix  $\mathbf{\Theta}$  does not change the error function  $S(\mathbf{\Theta}|\mathbf{X}, \mathbf{Y})$ . In this case the matrix  $\mathbf{X}^T \mathbf{X}$  is not invertible. To avoid the strong linear dependence, feature selection and dimensionality reduction techniques are used.

## 2 Feature selection

The feature selection goal is to find the index set  $\mathcal{A} = \{1, \dots, n\}$  of the matrix  $\mathbf{X}$  columns. To select the set  $\mathcal{A}$  among all possible  $2^n - 1$  subsets, introduce the feature selection quality criteria

$$\mathcal{A} = \arg \max_{\mathcal{A}' \subseteq \{1, \dots, n\}} Q(\mathcal{A}'|\mathbf{X}, \mathbf{Y}). \quad (3)$$

Once the solution  $\mathcal{A}$  for the problem (3) is known, the problem (2) becomes

$$S(\mathbf{\Theta}_{\mathcal{A}}|\mathbf{X}_{\mathcal{A}}, \mathbf{Y}) = \left\| \mathbf{Y} - \mathbf{X}_{\mathcal{A}} \mathbf{\Theta}_{\mathcal{A}}^T \right\|_2^2 \rightarrow \min_{\mathbf{\Theta}_{\mathcal{A}}}, \quad (4)$$

where the subscript  $\mathcal{A}$  indicates columns with indices from the set  $\mathcal{A}$ .

## 2.1 Quadratic Programming Feature Selection

One of the approach to the feature selection is to maximize feature relevances and minimize pairwise feature redundancy. The QPFS algorithm selects non-correlated features, which are relevant to the target vector  $\phi$  for the linear regression problem ( $r = 1$ )

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

Introduce two functions:  $\text{Sim}(\mathbf{X})$  and  $\text{Rel}(\mathbf{X}, \phi)$ . The  $\text{Sim}(\mathbf{X})$  measures the redundancy between features, the  $\text{Rel}(\mathbf{X}, \phi)$  contains relevances between each feature and the target vector  $\phi$ . We want to minimize the function  $\text{Sim}$  and maximize the  $\text{Rel}$  simultaneously.

QPFS offers the explicit way to construct the functions  $\text{Sim}$  and  $\text{Rel}$ . The method minimizes the following functional

$$(1 - \alpha) \cdot \underbrace{\mathbf{a}^T \mathbf{Q} \mathbf{a}}_{\text{Sim}} - \alpha \cdot \underbrace{\mathbf{b}^T \mathbf{a}}_{\text{Rel}} \rightarrow \min_{\substack{\mathbf{a} \in \mathbb{R}_+^n \\ \|\mathbf{a}\|_1=1}}. \quad (5)$$

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 matrix  $\mathbf{b}$ . The normalized vector  $\mathbf{a}$  shows the importance of each feature. The functional (5) penalizes the dependent features by the function  $\text{Sim}$  and encourages features relevant to the target by the function  $\text{Rel}$ . The parameter  $\alpha$  allows to control the trade-off between the functions  $\text{Sim}$  and the  $\text{Rel}$ . The authors of the original QPFS paper suggested the way to select  $\alpha$  and make  $\text{Sim}(\mathbf{X})$  and  $\text{Rel}(\mathbf{X}, \phi)$  impact the same

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

where  $\overline{\mathbf{Q}}$ ,  $\overline{\mathbf{b}}$  are the mean values of  $\mathbf{Q}$  and  $\mathbf{b}$  respectively. Apply the thresholding for  $\mathbf{a}$  to find the optimal feature subset:

$$j \in \mathcal{A} \Leftrightarrow a_j > \tau.$$

To measure similarity the authors use the absolute value of sample correlation coefficient between pairs of features for the function  $\text{Sim}$ , and between features and the target vector  $\phi$  for the function  $\text{Rel}$

$$\mathbf{Q} = \{|\text{corr}(\mathbf{x}_i, \mathbf{x}_j)|\}_{i,j=1}^n, \quad \mathbf{b} = \{|\text{corr}(\mathbf{x}_i, \phi)|\}_{i=1}^n. \quad (6)$$

The problem (5) is convex if the matrix  $\mathbf{Q}$  is positive semidefinite. In general it is not always true. To satisfy this condition, the matrix  $\mathbf{Q}$  spectrum is shifted and the matrix  $\mathbf{Q}$  is replaced by  $\mathbf{Q} - \lambda_{\min} \mathbf{I}$ , where  $\lambda_{\min}$  is a  $\mathbf{Q}$  minimal eigenvalue.

The functional (5) corresponds to the quality criteria  $Q(\mathcal{A}|\mathbf{X}, \phi)$

$$\mathcal{A} = \arg \max_{\mathcal{A}' \subseteq \{1, \dots, n\}} Q(\mathcal{A}'|\mathbf{X}, \phi) \Leftrightarrow \arg \min_{\mathbf{a} \in \mathbb{R}_+^n, \|\mathbf{a}\|_1=1} [\mathbf{a}^T \mathbf{Q} \mathbf{a} - \alpha \cdot \mathbf{b}^T \mathbf{a}]. \quad (7)$$

## 2.2 Multivariate QPFS

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, and the matrix  $\mathbf{Q}$  and the vector  $\mathbf{b}$  are equal to

$$\mathbf{Q} = \{|\text{corr}(\chi_i, \chi_j)|\}_{i,j=1}^n, \quad \mathbf{b} = \left\{ \sum_{k=1}^r |\text{corr}(\chi_i, \phi_k)| \right\}_{i=1}^n.$$

This approach does not use the dependencies in the columns of the matrix  $\mathbf{Y}$ . Let consider the following example:

$$\mathbf{X} = [\chi_1, \chi_2, \chi_3], \quad \mathbf{Y} = [\underbrace{\phi_1, \phi_1, \dots, \phi_1}_{r-1}, \phi_2],$$

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

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

We would like to select only two features. For such configuration the best feature subset is  $[\chi_1, \chi_2]$ . The feature  $\chi_2$  predicts the second target  $\phi_2$  and the combination of features  $\chi_1, \chi_2$  predicts the first component. The QPFS algorithm for  $r = 2$  gives the solution  $\mathbf{a} = [0.37, 0.61, 0.02]$ . It coincides with our knowledge. However, if we add the collinear columns to the matrix  $\mathbf{Y}$  and increase  $r$  to 5, the QPFS solution will be  $\mathbf{a} = [0.40, 0.17, 0.43]$ . Here we lost the relevant feature  $\chi_2$  and select the redundant feature  $\chi_3$ .

To take into account the dependencies in the columns of the matrix  $\mathbf{Y}$  we extend the QPFS functional (5) to the multivariate case. We add the term  $\text{Sim}(\mathbf{Y})$  and extend the term  $\text{Rel}(\mathbf{X}, \mathbf{Y})$ :

$$\alpha_1 \cdot \underbrace{\mathbf{a}_x^T \mathbf{Q}_x \mathbf{a}_x}_{\text{Sim}(\mathbf{X})} - \alpha_2 \cdot \underbrace{\mathbf{a}_x^T \mathbf{B} \mathbf{a}_y}_{\text{Rel}(\mathbf{X}, \mathbf{Y})} + \alpha_3 \cdot \underbrace{\mathbf{a}_y^T \mathbf{Q}_y \mathbf{a}_y}_{\text{Sim}(\mathbf{Y})} \rightarrow \min_{\substack{\mathbf{a}_x \in \mathbb{R}_+^n, \|\mathbf{a}_x\|_1=1 \\ \mathbf{a}_y \in \mathbb{R}_+^r, \|\mathbf{a}_y\|_1=1}}. \quad (8)$$

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

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

The vector  $\mathbf{a}_x$  shows the feature importances, while  $\mathbf{a}_y$  is a vector with the importance of each target. The targets which are correlated will be penalized by  $\text{Sim}(\mathbf{Y})$  and have the lower importances.

60 **Statement 1.** For the case  $r = 1$  the proposed functional (8) coincides with the original  
 61 QPFS algorithm (5).

62 *Proof.* If  $r$  is equal to 1, then  $\mathbf{Q}_y = \mathbf{1}$ ,  $\mathbf{a}_y = \mathbf{1}$ ,  $\mathbf{B} = \mathbf{b}$ . It reduces the problem (8) to

$$\alpha_1 \cdot \mathbf{a}_x^T \mathbf{Q}_x \mathbf{a}_x - \alpha_2 \cdot \mathbf{a}_x^T \mathbf{b} \rightarrow \min_{\mathbf{a}_x \in \mathbb{R}_+^n, \|\mathbf{a}_x\|_1=1}.$$

63 Setting  $\alpha = \frac{\alpha_2}{\alpha_1 + \alpha_2}$  brings to the original QPFS problem (5).  $\square$

64 The coefficients  $\alpha_1$ ,  $\alpha_2$ , and  $\alpha_3$  control the influence of each term to the functional (8)  
 65 and satisfy the conditions:

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

66 We balance the terms  $\text{Sim}(\mathbf{X})$  and  $\text{Rel}(\mathbf{X}, \mathbf{Y})$  by fixing the proportion between  $\alpha_1$  and  $\alpha_2$ .

67 **Statement 2.** Balance between the terms  $\text{Sim}(\mathbf{X})$  and  $\text{Rel}(\mathbf{X}, \mathbf{Y})$  for the problem (8) is  
 68 achieved by the following coefficients:

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

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

70 *Proof.* The impact of these terms are equal if  $\alpha_1 \cdot \text{Sim}(\mathbf{X}) = \alpha_2 \cdot \text{Rel}(\mathbf{X}, \mathbf{Y})$ . The mean  
 71 values of the terms  $\text{Sim}(\mathbf{X})$  and  $\text{Rel}(\mathbf{X}, \mathbf{Y})$  are given by the mean values  $\overline{\mathbf{Q}_x}$  and  $\overline{\mathbf{B}}$  of the  
 72 corresponding matrices  $\mathbf{Q}_x$  and  $\mathbf{B}$ . Since  $\alpha_1 + \alpha_2 + \alpha_3 = 1$ , we obtain  $(1 - \alpha_3 - \alpha_2)\overline{\mathbf{Q}_x} = \alpha_2\overline{\mathbf{B}}$ .  
 73 Express  $\alpha_2$  to get

$$\alpha_2 = \frac{(1 - \alpha_3)\overline{\mathbf{Q}_x}}{\overline{\mathbf{Q}_x} + \overline{\mathbf{B}}}.$$

74 The value for  $\alpha_1$  is derived from the  $\alpha_1 + \alpha_2 + \alpha_3 = 1$ .  $\square$

75 We apply the proposed algorithm to the discussed example. The given matrix  $\mathbf{Q}$  corre-  
 76 sponds to the matrix  $\mathbf{Q}_x$ . We additionally define the matrix  $\mathbf{Q}_y$  by setting  $\text{corr}(\phi_1, \phi_2) =$   
 77 0.2 and all others entries to one. Figure 1 shows the importances of features  $\mathbf{a}_x$  and tar-  
 78 gets  $\mathbf{a}_y$  with respect to  $\alpha_3$  coefficient. If  $\alpha_3$  is small, the impact of all targets are almost  
 79 equal and the feature  $\chi_3$  dominates the feature  $\chi_2$ . When  $\alpha_3$  becomes larger than 0.2, the  
 80 importance  $(\mathbf{a}_y)_5$  of the target  $\phi_5$  grows up along with the importance of the feature  $\chi_2$ .

### 81 3 Feature categorization

82 Feature selection algorithms eliminate features which are not relevant to the target variable.  
 83 To determine whether the feature is relevant the t-test could be applied for the correlation  
 84 coefficient.

$$r = \text{corr}(\chi, \phi), \quad t = \frac{r\sqrt{m-2}}{1-r^2} \sim \text{St}(m-2).$$

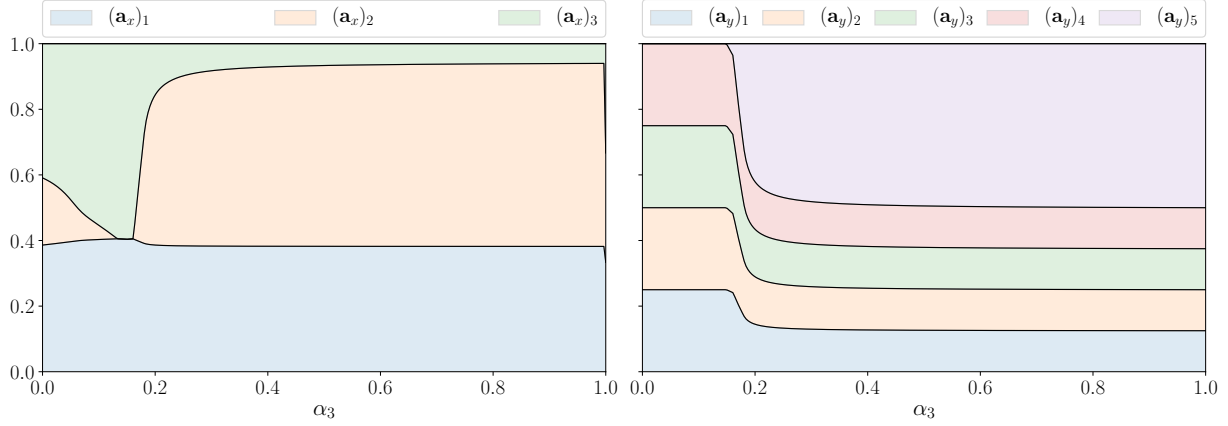


Figure 1: Feature importances  $\mathbf{a}_x$  and  $\mathbf{a}_y$  with respect to the  $\alpha_3$  coefficient

$$H_0 : r = 0$$

$$H_1 : r \neq 0$$

85 If features are relevant, but correlated, feature selection methods pick the subset of them  
 86 to reduce the multicollinearity and redundancy. The goal is to find relevant, non-correlated  
 87 features. However, in this case the correlations between targets in matrix  $\mathbf{Y}$  are crucial.  
 88 To measure the dependence of each feature or target, the Variance Inflation Factor is  
 89 computed

$$\text{VIF}(\mathbf{x}_j) = \frac{1}{1 - R_j^2}, \quad \text{VIF}(\phi_k) = \frac{1}{1 - R_k^2},$$

90 where  $R_j^2(R_k^2)$  are coefficients of determination for the regression of  $\mathbf{x}_j(\phi_k)$  on the other  
 91 features(targets).

92 On that basis, we categorize features into 5 disjoint groups:

93 1. non-relevant features

$$\{j : \text{corr}(\mathbf{x}_j, \phi_k) = 0, \forall k \in \{1, \dots, r\}\};$$

94 2. non- $\mathbf{X}$ -correlated features, which are relevant to non- $\mathbf{Y}$ -correlated targets

$$\{j : (\text{VIF}(\mathbf{x}_j) < 10) \text{ and } (\text{VIF}(\phi_k) < 10, \forall k \in \{1, \dots, r\} : \text{corr}(\mathbf{x}_j, \phi_k) \neq 0)\};$$

95 3. non- $\mathbf{X}$ -correlated features, which are relevant to  $\mathbf{Y}$ -correlated targets

$$\{j : (\text{VIF}(\mathbf{x}_j) < 10) \text{ and } (\exists k \in \{1, \dots, r\} : \text{VIF}(\phi_k) > 10 \text{ \& } \text{corr}(\mathbf{x}_j, \phi_k) \neq 0)\};$$

96 4.  $\mathbf{X}$ -correlated features, which are relevant to non- $\mathbf{Y}$ -correlated targets

$$\{j : (\text{VIF}(\mathbf{x}_j) > 10) \text{ and } (\text{VIF}(\phi_k) < 10, \forall k \in \{1, \dots, r\} : \text{corr}(\mathbf{x}_j, \phi_k) \neq 0)\};$$

97 5.  $\mathbf{X}$ -correlated features, which are relevant to  $\mathbf{Y}$ -correlated targets

$$\{j : (\text{VIF}(\mathbf{x}_j) > 10) \text{ and } (\exists k \in \{1, \dots, r\} : \text{VIF}(\phi_k) > 10 \text{ \& } \text{corr}(\mathbf{x}_j, \phi_k) \neq 0)\}.$$

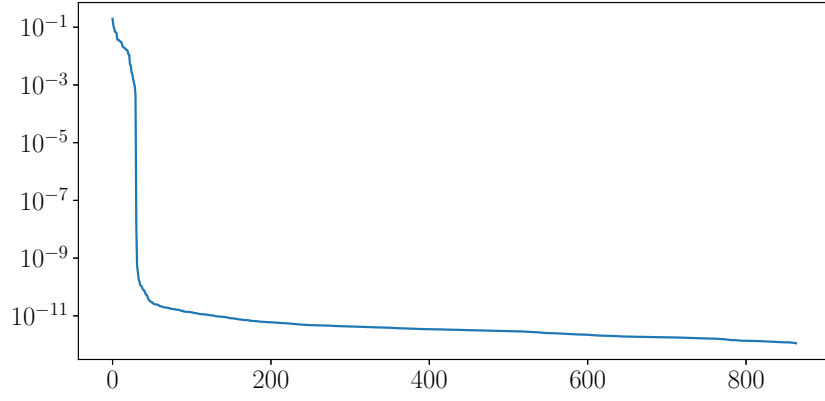


Figure 2: Sorted feature importances for the QPFS algorithm

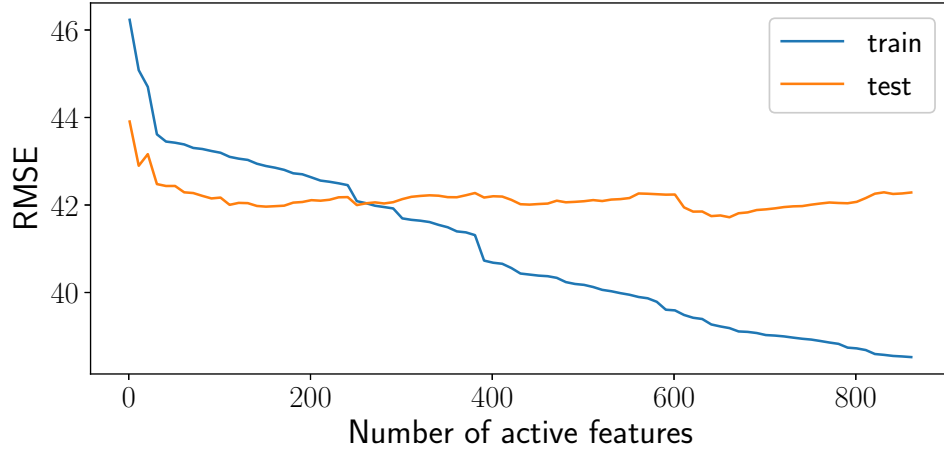


Figure 3: RMSE with respect to size of active set, features are ranked by QPFS algorithm

## 98 4 Experiment

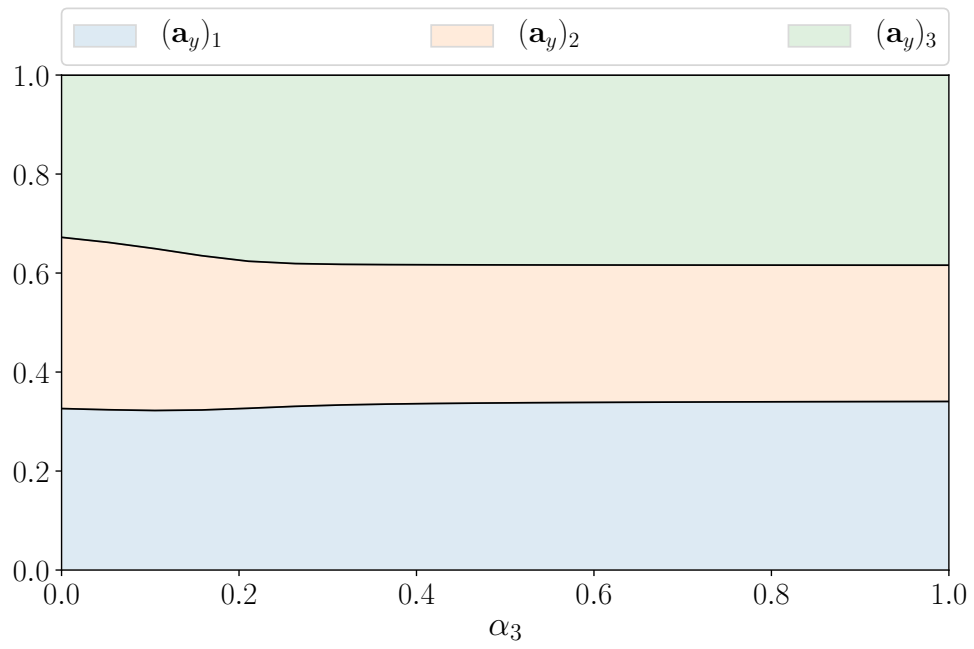


Figure 4: Targets importances for ECoG data, each component is related to one of the axis