

Orthogonal Least Squares Based Fast Feature Selection for Linear Classification

Sikai Zhang^a, Zi-Qiang Lang^{b,*}

^a*Department of Mechanical Engineering, The University of Sheffield, Sheffield, United Kingdom*

^b*Department of Automatic Control and Systems Engineering, The University of Sheffield, Sheffield, United Kingdom*

Abstract

An Orthogonal Least Squares (OLS) based feature selection method is proposed for both binomial and multinomial classification. The novel Squared Orthogonal Correlation Coefficient (SOCC) is defined based on Error Reduction Ratio (ERR) in OLS and used as the feature ranking criterion. The equivalence between the canonical correlation coefficient, Fisher's criterion, and the sum of the SOCCs is revealed, which unveils the statistical implication of ERR in OLS for the first time. It is also shown that the OLS based feature selection method has speed advantages when applied for greedy search. The proposed method is comprehensively compared with the mutual information based feature selection methods and the embedded methods using both synthetic and real world datasets. The results show that the proposed method is always in the top 5 among the 12 candidate methods. Besides, the proposed method can be directly applied to continuous features without discretisation, which is another significant advantage over mutual information based methods.

Keywords: Feature selection, orthogonal least squares, canonical correlation analysis, linear discriminant analysis, multi-label, multivariate time series, feature interaction.

1. Introduction

The aim of the feature selection for classification is to select an optimal subset of features given the candidate features, which are numerical or categorical, and the response, which is categorical. The feature selection methods can be divided into three types: filter, wrapper, and embedded methods [1, 2, 3]. The filter methods rank the individual candidate features based on certain statistics, such as the correlation coefficient and mutual information [4, 5, 6]. The wrapper methods train classifiers by ranking the subsets of candidate

*Corresponding author

Email addresses: matthew.szhang91@gmail.com (Sikai Zhang), z.lang@sheffield.ac.uk (Zi-Qiang Lang)

features based on their classification performance. The embedded methods, e.g. LASSO [7] and CART [8], select optimal features during the training process of a specific classifier.

Comparing with the other two methods, a filter method is not based on a specific type of classifiers, so a filter method is more suitable to be used in the early stage where the type of classifiers has not been decided. To rank the features by a filter method, it is desired that the features in the subset have the high relevance to the response, while the low redundancy between themselves. A straightforward way is to optimise the objective function constructed by the difference or the quotient between the relevance and the redundancy. For example, the well-known minimal-Redundancy-Maximal-Relevance (mRMR) method adopts this idea, in which the relevance and redundancy are quantified by the mutual information [9]. The second idea is to control the redundancy by orthogonalising the candidate features, and to find the maximum relevance between the orthogonalised features and the response. The second idea has been used in the term selection of time series models by Orthogonal Least Squares (OLS), where the relevance is defined by the Error Reduction Ratio (ERR) [10]. These two ideas basically evaluate the relevance between the single feature and the response, and the relevance is analysed separately with the redundancy. The third idea uses the overall relevance between the subset features and the response. The definition of the overall relevance has taken the redundancy into consideration, e.g. the multiple correlation coefficient and the canonical correlation coefficient [11]. The first idea is extensively used in mutual information based filter methods [12]. However, the idea does not take the feature interaction into consideration, which makes the filter methods have a well-known drawback compared to the wrapper and embedded methods [13]. The interaction¹ between features exists when a feature has to be combined with one or more other features to represent the response [14]. Without considering the feature interaction, the feature selection methods will fail to select the features having low relevance individually but high relevance together. In this paper, the proposed feature selection method is based on the third idea with the feature interaction issue addressed using an approach similar to the wrapper methods in order to simultaneously handle the feature relevance, feature redundancy, and feature interaction. This can, in conjunction with the second idea, achieve a faster computation speed. In addition, compared to the mutual information based methods, which can only work with the discrete or categorical features [15, 16], the proposed method is applicable to both numerical (including discrete and continuous) and categorical features. The proposed method is also closely related to two recent advanced topics in the feature selection field, which are the multi-label feature selection [17, 18, 19] and multivariate time series

¹The feature interaction can also be understood in the way of the conditional redundancy [12].

feature selection [20, 21], respectively. The fundamental issue that is addressed under the two research topics is to develop methods that can deal with the features and response that have to be represented in a matrix form. The proposed method can naturally deal with the feature and response represented by a matrix without a need to introduce any additional techniques.

Basically, the contributions of the present study are in two aspects. First, the study reveals, for the first time, the relationships between the OLS and some well-known statistics including multiple correlation coefficient, canonical correlation coefficient, and Fisher’s criterion. Second, via utilising these relationships, a novel feature selection method for classification is developed. The novel method can deal with both numerical (including continuous and discrete) and categorical features, has a much faster computation speed when used with a greedy search, and can simultaneously address issues associated with feature relevance, feature redundancy, and feature interaction. The computational efficiency and general applicability show the proposed method has potential to be widely applied to address feature selection issues in classification problems.

The rest of the paper is organised as follows. The related work about OLS is introduced in Section 2. In Section 3, based on OLS, the definition of the SOCCs is given. The relationships of the SOCCs with multiple correlation coefficient and canonical correlation coefficient are analysed. Then, via these relationships, an OLS based feature selection method is developed for binomial classification (Section 4) and multinomial classification (Section 5), respectively. After that, the speed advantage of the method in the greedy search is analysed for both binomial and multinomial classification problems. The relationship of the SOCC with Linear Discriminant Analysis (LDA) is also studied in Section 5 to demonstrate the statistical implication of the proposed SOCC in classification. In Section 6, a detailed example is provided to illustrate the procedure of the proposed method, and its relationship with Canonical Correlation Analysis (CCA) and linear discriminant analysis. Moreover, a comprehensive comparison of the proposed method with the mutual information based methods and the embedded methods is carried out on both synthetic and real world datasets. Finally, conclusions are summarised in Section 7.

2. Related work

An important basis of the new method proposed in the present study is the OLS. The OLS and associated representative works are briefly summarised in Table 1. OLS was firstly developed by Korenberg for the fast parameter estimation [22] and term selection [23] of the polynomial Nonlinear AutoRegressive with eXogenous input (NARX) model. The criterion used for term selection is the Mean-Square Error Reduction

(MSER). Given a linear regression model

$$\mathbf{y} = (\mathbf{1}, \mathbf{X}) \begin{pmatrix} \beta_0 \\ \boldsymbol{\beta} \end{pmatrix} + \mathbf{e}, \quad (1)$$

where the response vector is

$$\mathbf{y} = (y_1, \dots, y_N)^\top, \quad (2)$$

the design matrix of n independent variables with a constant term is

$$(\mathbf{1}, \mathbf{X}) = (\mathbf{1}, \mathbf{x}_1, \dots, \mathbf{x}_n) = \begin{pmatrix} 1 & x_{1,1} & \dots & x_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N,1} & \dots & x_{N,n} \end{pmatrix}, \quad (3)$$

the unknown parameter vector is

$$\begin{pmatrix} \beta_0 \\ \boldsymbol{\beta} \end{pmatrix} = (\beta_0, \beta_1, \dots, \beta_n)^\top, \quad (4)$$

the error term is

$$\mathbf{e} = (e_1, \dots, e_N)^\top, \quad (5)$$

the MSER is defined as

$$\text{MSER}_i = \hat{g}_i^{*2} \mathbf{w}_i^{*\top} \mathbf{w}_i^* \quad i = 0, \dots, n. \quad (6)$$

In (6), \mathbf{w}_i is the i^{th} column of matrix $\mathbf{W}^* = (\mathbf{w}_0^*, \mathbf{w}_1^*, \dots, \mathbf{w}_n^*)$ determined from the unnormalised reduced QR decomposition of the design matrix $(\mathbf{1}, \mathbf{X})$, such that

$$(\mathbf{1}, \mathbf{X}) = \mathbf{W}^* \mathbf{A}^* \quad (7)$$

where \mathbf{A}^* is the upper triangle matrix from the QR decomposition. In (6), \hat{g}_i^* is obtained by solving the normal equation

$$\mathbf{W}^{*\top} \mathbf{W}^* \begin{pmatrix} \hat{g}_0^* \\ \vdots \\ \hat{g}_n^* \end{pmatrix} = \mathbf{W}^{*\top} \mathbf{y}. \quad (8)$$

Chen et al. [10] extended OLS to the polynomial Nonlinear AutoRegressive Moving Average with eXogenous input (NARMAX) model, where the term ranking criterion is changed to the ERR, which is the MSER normalised by the inner product of the measured response, i.e.

$$\text{ERR}_i = \frac{\hat{g}_i^{*2} \mathbf{w}_i^{*\top} \mathbf{w}_i^*}{\mathbf{y}^\top \mathbf{y}} \quad i = 0, \dots, n. \quad (9)$$

Later, OLS became well-known in the nonlinear system identification field and the researchers explored the application of OLS in machine learning. For example, Chen et al. [24] applied the ERR to choose the centres of Radial Basis Functions (RBFs) for training RBF neural networks. Wei and Billings [25] extended the ERR to the unsupervised feature selection, where the ERR is applied to evaluate the explanation capability of the selected features to all candidate features. Recently, Solares et al. [26] applied the OLS method to the selection of the terms in the logistic-NARX model, which maps the continuous response of the NARX model to the binary response by a logistic function. The point-Biserial Correlation Coefficient between the Orthogonalised term and measured response was adopted as the term ranking criterion (OBCC). However, none of these previous studies has revealed any relationship between the OLS and well-known statistics such as canonical correlation coefficient and Fisher’s criterion. Probably because of this, the OLS has never been used as a general feature selection method to solve binomial or multinomial classification problems.

In the present study, a novel feature ranking criterion referred to as Squared Orthogonal Correlation Coefficient (SOCC) is defined using ERR associated with the OLS approach to a standard linear regression problem. For the first time, the SOCC reveals the statistic implication of ERR in OLS and unveils a significant relationship between the ERR and classic statistics such as canonical correlation coefficient and Fisher’s criterion. This, consequently, enables the development of an effective Canonical Correlation Analysis (CCA) based fast feature selection approach for both binomial and multinomial classifications.

Table 1. A summary of OLS and associated representative works.

Author	Year	Criterion	Task
Korenberg [23]	1989	MSER	Term selection for NARX model
Chen et al. [10]	1989	ERR	Term selection for NARMAX model
Chen et al. [24]	1991	ERR	RBF centre selection
Wei and Billings [25]	2006	ERR	Unsupervised feature selection
Solares et al. [26]	2019	OBCC	Term selection for the logistic-NARX model
Zhang and Lang (This paper)	2020	SOCC	Supervised feature selection for classification

3. Squared orthogonal correlation coefficients

3.1. Definition

The vectors \mathbf{y} and \mathbf{w}_i^* used in the traditional ERR (9) normally have non-zero mean values, which makes the connection between ERR and other classic statistics lost. To overcome this problem, the SOCCs between

\mathbf{X} and \mathbf{y} are defined as

$$h_i = \frac{\hat{s}_i^2 \mathbf{w}_{Ci}^\top \mathbf{w}_{Ci}}{\mathbf{y}_C^\top \mathbf{y}_C} \quad i = 1, \dots, n, \quad (10)$$

In (10), \mathbf{y}_C is the centred \mathbf{y} , \mathbf{w}_{Ci} is the i^{th} column of matrix $\mathbf{W}_C = (\mathbf{w}_{C1}, \dots, \mathbf{w}_{Cn})$ determined from the unnormalised reduced QR decomposition of the centred \mathbf{X} , denoted as \mathbf{X}_C , such that

$$\mathbf{X}_C = \mathbf{W}_C \mathbf{A} \quad (11)$$

where \mathbf{A} is the upper triangle matrix of the QR decomposition. In (10), the parameter \hat{s}_i is obtained from the normal equation

$$\mathbf{W}_C^\top \mathbf{W}_C \begin{pmatrix} \hat{g}_1 \\ \vdots \\ \hat{g}_n \end{pmatrix} = \mathbf{W}_C^\top \mathbf{y}_C. \quad (12)$$

As \mathbf{W}_C is orthogonal, the inner product $\mathbf{W}_C^\top \mathbf{W}_C$ is the diagonal matrix $\text{diag}(\mathbf{w}_{C1}^\top \mathbf{w}_{C1}, \dots, \mathbf{w}_{Cn}^\top \mathbf{w}_{Cn})$. Thus, the computation of the parameter vector \hat{g}_i can be simplified as

$$\hat{g}_i = \frac{\mathbf{w}_{Ci}^\top \mathbf{y}_C}{\mathbf{w}_{Ci}^\top \mathbf{w}_{Ci}}. \quad (13)$$

Correspondingly, the SOCCs (10) can be rewritten as

$$h_i = \frac{\mathbf{y}_C^\top \mathbf{w}_{Ci} \mathbf{w}_{Ci}^\top \mathbf{y}_C}{\mathbf{w}_{Ci}^\top \mathbf{w}_{Ci} \mathbf{y}_C^\top \mathbf{y}_C}, \quad i = 1, \dots, n. \quad (14)$$

The SOCCs have a close relationship with the Pearson correlation coefficient, multiple correlation coefficient and the canonical correlation coefficients.

3.2. Relationship with Pearson correlation coefficient

The sample Pearson correlation coefficient between $\boldsymbol{\gamma} \in \mathbb{R}^n$ and $\boldsymbol{\omega} \in \mathbb{R}^n$ is given by

$$r(\boldsymbol{\gamma}, \boldsymbol{\omega}) = \frac{\boldsymbol{\gamma}_C^\top \boldsymbol{\omega}_C}{\sqrt{\boldsymbol{\gamma}_C^\top \boldsymbol{\gamma}_C} \sqrt{\boldsymbol{\omega}_C^\top \boldsymbol{\omega}_C}}, \quad (15)$$

where $\boldsymbol{\gamma}_C \in \mathbb{R}^n$ is the vector $\boldsymbol{\gamma}$ centred by its sample mean and $\boldsymbol{\omega}_C \in \mathbb{R}^n$ is the vector $\boldsymbol{\omega}$ centred by its sample mean. Obviously, the SOCCs h_i in (14) is the squared sample Pearson correlation coefficient between \mathbf{y} and \mathbf{w}_{Ci} , i.e.

$$r^2(\mathbf{y}, \mathbf{w}_{Ci}) = h_i \quad i = 1, \dots, n. \quad (16)$$

3.3. Relationship with multiple correlation coefficient

The multiple correlation coefficient is the measure of linear association between two or more independent variables and a dependent variable. If the n columns in the design matrix \mathbf{X} are the samples of n independent variables and the response vector \mathbf{y} is the samples of a dependent variable, the association between \mathbf{X} and \mathbf{y} can be measured by the multiple correlation coefficient $R(\mathbf{X}, \mathbf{y})$ or $R(\mathbf{y}, \mathbf{X})$. The multiple correlation analysis of \mathbf{X} and \mathbf{y} is to find a projection direction, so that the Pearson correlation coefficient between \mathbf{y}_C and the projected \mathbf{X}_C is maximised. The optimal projection direction is the solution $\hat{\boldsymbol{\beta}}$ of the normal equation [27]

$$(\mathbf{X}_C^T \mathbf{X}_C) \hat{\boldsymbol{\beta}} = \mathbf{X}_C^T \mathbf{y}_C. \quad (17)$$

Then, the multiple correlation coefficient $R(\mathbf{X}, \mathbf{y})$ or $R(\mathbf{y}, \mathbf{X})$ is defined as [27]

$$R(\mathbf{X}, \mathbf{y}) = R(\mathbf{y}, \mathbf{X}) = r(\hat{\mathbf{y}}_C, \mathbf{y}_C) = \frac{\hat{\mathbf{y}}_C^T \mathbf{y}_C}{\sqrt{\hat{\mathbf{y}}_C^T \hat{\mathbf{y}}_C} \sqrt{\mathbf{y}_C^T \mathbf{y}_C}}, \quad (18)$$

where

$$\hat{\mathbf{y}}_C = \mathbf{X}_C \hat{\boldsymbol{\beta}}. \quad (19)$$

The relationship between the SOCCs and the multiple correlation coefficient is (see Appendix A for proof)

$$R^2(\mathbf{X}, \mathbf{y}) = \sum_{i=1}^n h_i. \quad (20)$$

In other words, the sum of the SOCCs between \mathbf{X} and \mathbf{y} is equal to the squared multiple correlation coefficient between \mathbf{X} and \mathbf{y} .

3.4. Relationship with canonical correlation coefficient

The canonical correlation coefficient is the measure of linear association between two or more independent variables and two or more dependent variables. Given a response matrix as

$$\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_m) = \begin{pmatrix} y_{1,1} & \dots & y_{1,m} \\ \vdots & \ddots & \vdots \\ y_{N,1} & \dots & y_{N,m} \end{pmatrix}, \quad (21)$$

if the n columns in the design matrix \mathbf{X} are the samples of n independent variables and the m columns in the response matrix \mathbf{Y} are the samples of m dependent variables, the association between \mathbf{X} and \mathbf{Y} can be measured by the canonical correlation coefficient $R(\mathbf{X}, \mathbf{Y})$. The Canonical Correlation Analysis (CCA) for \mathbf{X}

and \mathbf{Y} is to find a pair of the projection directions \mathbf{a} and \mathbf{b} , so that the Pearson correlation coefficient between $\mathbf{X}_C \mathbf{a}$ and $\mathbf{Y}_C \mathbf{b}$ is maximised, that is

$$\arg \max_{\mathbf{a}, \mathbf{b}} r(\mathbf{X}_C \mathbf{a}, \mathbf{Y}_C \mathbf{b}), \quad (22)$$

where

$$\begin{aligned} \mathbf{Y}_C &= (\mathbf{y}_{C1}, \dots, \mathbf{y}_{Cm}) \\ &= \begin{pmatrix} y_{1,1} - \bar{y}_1 & \dots & y_{1,m} - \bar{y}_m \\ \vdots & \ddots & \vdots \\ y_{N,1} - \bar{y}_1 & \dots & y_{N,m} - \bar{y}_m \end{pmatrix}, \end{aligned} \quad (23)$$

and \bar{y}_i is the sample mean of \mathbf{y}_i . The canonical correlation coefficient between \mathbf{X} and \mathbf{Y} can be computed by

$$R(\mathbf{X}, \mathbf{Y}) = r(\mathbf{X}_C \mathbf{a}, \mathbf{Y}_C \mathbf{b}) = \frac{\mathbf{a}^\top \mathbf{R}_{\mathbf{X}, \mathbf{Y}} \mathbf{b}}{\sqrt{\mathbf{a}^\top \mathbf{R}_{\mathbf{X}, \mathbf{X}} \mathbf{a}} \sqrt{\mathbf{b}^\top \mathbf{R}_{\mathbf{Y}, \mathbf{Y}} \mathbf{b}}}, \quad (24)$$

where the correlation matrices are given by

$$\mathbf{R}_{\mathbf{X}, \mathbf{Y}} = \begin{pmatrix} r(\mathbf{x}_1, \mathbf{y}_1) & \dots & r(\mathbf{x}_1, \mathbf{y}_m) \\ \vdots & \ddots & \vdots \\ r(\mathbf{x}_n, \mathbf{y}_1) & \dots & r(\mathbf{x}_n, \mathbf{y}_m) \end{pmatrix} \quad \mathbf{R}_{\mathbf{X}, \mathbf{X}} = \begin{pmatrix} r(\mathbf{x}_1, \mathbf{x}_1) & \dots & r(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ r(\mathbf{x}_n, \mathbf{x}_1) & \dots & r(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix} \quad \mathbf{R}_{\mathbf{Y}, \mathbf{Y}} = \begin{pmatrix} r(\mathbf{y}_1, \mathbf{y}_1) & \dots & r(\mathbf{y}_1, \mathbf{y}_m) \\ \vdots & \ddots & \vdots \\ r(\mathbf{y}_m, \mathbf{y}_1) & \dots & r(\mathbf{y}_m, \mathbf{y}_m) \end{pmatrix}. \quad (25)$$

The multiple correlation coefficient $R(\mathbf{X}, \mathbf{y})$ is a special case of the canonical correlation coefficient $R(\mathbf{X}, \mathbf{Y})$, when \mathbf{Y} is a column vector \mathbf{y} . The CCA can be transformed to the eigenvalue problem given by [11, p. 173]

$$\mathbf{R}_{\mathbf{X}, \mathbf{X}}^{-1} \mathbf{R}_{\mathbf{X}, \mathbf{Y}} \mathbf{R}_{\mathbf{Y}, \mathbf{Y}}^{-1} \mathbf{R}_{\mathbf{Y}, \mathbf{X}} \mathbf{a} = R^2(\mathbf{X}, \mathbf{Y}) \mathbf{a} \quad (26a)$$

$$\mathbf{R}_{\mathbf{Y}, \mathbf{Y}}^{-1} \mathbf{R}_{\mathbf{Y}, \mathbf{X}} \mathbf{R}_{\mathbf{X}, \mathbf{X}}^{-1} \mathbf{R}_{\mathbf{X}, \mathbf{Y}} \mathbf{b} = R^2(\mathbf{X}, \mathbf{Y}) \mathbf{b}. \quad (26b)$$

The two projection directions \mathbf{a} and \mathbf{b} are the eigenvectors, and the eigenvalue is the square of the canonical correlation coefficient. If \mathbf{X}_C and \mathbf{Y}_C have full column rank, the number of the non-zero solutions of (26) is not more than $n \wedge m$, where the operator \wedge returns the minimum of two values on both sides. Thus, in contrast with the multiple correlation coefficient which only has one value, there are $n \wedge m$ canonical correlation coefficients (which may contain zeros) for \mathbf{X} and \mathbf{Y} , which are denoted as $R_1(\mathbf{X}, \mathbf{Y}), \dots, R_{n \wedge m}(\mathbf{X}, \mathbf{Y})$.

To connect the SOCCs with the canonical correlation coefficients, the SOCCs between \mathbf{X} and \mathbf{Y} are defined as

$$h_{i,j} = \frac{\mathbf{v}_{Cj}^\top \mathbf{w}_{Ci} \mathbf{w}_{Ci}^\top \mathbf{v}_{Cj}}{\mathbf{w}_{Ci}^\top \mathbf{w}_{Ci} \mathbf{v}_{Cj}^\top \mathbf{v}_{Cj}}, \quad (27)$$

where \mathbf{w}_{Ci} is the i^{th} column of matrix $\mathbf{W}_C = (\mathbf{w}_{C1}, \dots, \mathbf{w}_{Cn})$ determined from the unnormalised reduced QR decomposition of \mathbf{X}_C , and \mathbf{v}_{Cj} is the j^{th} column of matrix $\mathbf{V}_C = (\mathbf{v}_{C1}, \dots, \mathbf{v}_{Cm})$ determined from the unnormalised reduced QR decomposition of \mathbf{Y}_C , i.e.

$$\begin{aligned}\mathbf{X}_C &= \mathbf{W}_C \mathbf{A} \\ \mathbf{Y}_C &= \mathbf{V}_C \mathbf{B}.\end{aligned}\tag{28}$$

The matrices \mathbf{A} and \mathbf{B} are the upper triangle matrices of the QR decomposition. The relationship is (see Appendix B for proof)

$$\sum_{k=1}^{n \wedge m} R_k^2(\mathbf{X}, \mathbf{Y}) = \sum_{i=1}^n \sum_{j=1}^m h_{i,j},\tag{29}$$

which is a natural extension of (20) to the case where the response vector \mathbf{y} becomes the response matrix \mathbf{Y} .

4. OLS based fast feature selection for binomial classification

If the N instances of \mathbf{X} belong to two classes and the n variables in \mathbf{X} represent n features, the feature selection problem for the binomial classification is to find the t features from the n features of \mathbf{X} , which is optimal to classify the N instances into the two classes. The two classes can be assigned values 0 and 1 to form a dummy response vector \mathbf{y} for the N instances. Thus, the goodness of the features for the classification can be evaluated by the multiple correlation coefficient between the features of interest and the dummy response vector. In fact, the two classes can be assigned to any distinct values to form the response vector, which has no effect on the value of the multiple correlation coefficient. However, to be consistent to the multinomial classification case, the dummy encoding is adopted. The optimal t features can be searched by comparing all $\binom{n}{t}$ feature combinations exhaustively, where

$$\binom{n}{t} = \frac{n!}{t!(n-t)!}.\tag{30}$$

In some cases, the *exhaustive search* is too expensive in computation. A realistic approach is to select only one feature in one step. In each step, the previously selected features will not be changed. For example, the three ‘optimal’ features can be selected in three steps as shown in Table 2. As each step selects the feature which maximises the multiple correlation, the search is referred to the *greedy search* [1]. The algorithm of the OLS based greedy feature selection for binomial classification can be summarised in 5 steps.

Input:

\mathbf{X} : $N \times n$ matrix containing N instances and n features.

Table 2. An example for selecting three features from n features by the greedy search for binomial classification, where $i = 1, \dots, n$ for step 1, $i = 1, 2, 4, 5, \dots, n$ for step 2, $i = 1, 2, 4, 6, 7, \dots, n$ for step 3.

	Multiple Correlation	Selected Feature
Step 1	$R(\mathbf{x}_3, \mathbf{y}) \geq R(\mathbf{x}_i, \mathbf{y})$	\mathbf{x}_3
Step 2	$R((\mathbf{x}_3, \mathbf{x}_5), \mathbf{y}) \geq R((\mathbf{x}_3, \mathbf{x}_i), \mathbf{y})$	$\mathbf{x}_3, \mathbf{x}_5$
Step 3	$R((\mathbf{x}_3, \mathbf{x}_5, \mathbf{x}_1), \mathbf{y}) \geq R((\mathbf{x}_3, \mathbf{x}_5, \mathbf{x}_i), \mathbf{y})$	$\mathbf{x}_3, \mathbf{x}_5, \mathbf{x}_1$

\mathbf{y} : $N \times 1$ vector formed by dummy encoding.

t : The number of features to be selected.

Step 1. Centre \mathbf{X} into \mathbf{X}_C , centre \mathbf{y} into \mathbf{y}_C , and let $p = 0$.

Step 2. Divide \mathbf{X} into $(\mathbf{X}_s, \mathbf{X}_r)$, where the selected feature matrix is given by

$$\mathbf{X}_s = (\mathbf{x}_{s1}, \dots, \mathbf{x}_{sp}), \quad (31)$$

and the remaining feature matrix is given by

$$\mathbf{X}_r = (\mathbf{x}_{r1}, \dots, \mathbf{x}_{rq}), \quad (32)$$

where p is the number of the selected features, and q is the number of the remaining features. Correspondingly, divide \mathbf{X}_C into $(\mathbf{X}_{Cs}, \mathbf{X}_{Cr})$, where

$$\begin{aligned} \mathbf{X}_{Cs} &= (\mathbf{x}_{Cs1}, \dots, \mathbf{x}_{Csp}) \\ \mathbf{X}_{Cr} &= (\mathbf{x}_{Cr1}, \dots, \mathbf{x}_{Crq}). \end{aligned} \quad (33)$$

Step 3. If $p = 0$, let

$$\begin{aligned} \mathbf{W}_{Cr} &= \mathbf{X}_{Cr} \\ \mathbf{w}_{Cri} &= \mathbf{x}_{Cri}, \quad i = 1, \dots, q. \end{aligned} \quad (34)$$

Otherwise, first, orthogonalise \mathbf{X}_{Cs} into \mathbf{W}_{Cs} , where

$$\mathbf{W}_{Cs} = (\mathbf{w}_{Cs1}, \dots, \mathbf{w}_{Csp}), \quad (35)$$

and $\mathbf{w}_{Csi}^\top \mathbf{w}_{Cs j} = 0$ for $i \neq j$. Then, orthogonalise each feature in \mathbf{X}_{Cr} to \mathbf{W}_{Cs} to form the matrix \mathbf{W}_{Cr} , where

$$\mathbf{W}_{Cr} = (\mathbf{w}_{Cr1}, \dots, \mathbf{w}_{Crq}), \quad (36)$$

and \mathbf{w}_{Cri} is obtained through the classical Gram-Schmidt process, which is given by

$$\mathbf{w}_{\text{Cri}} = \mathbf{x}_{\text{Cri}} - \sum_{j=1}^p \frac{\mathbf{x}_{\text{Cri}}^\top \mathbf{w}_{\text{Cs}j}}{\mathbf{w}_{\text{Cs}j}^\top \mathbf{w}_{\text{Cs}j}} \mathbf{w}_{\text{Cs}j}, \quad i = 1, \dots, q. \quad (37)$$

It should be noticed that \mathbf{w}_{Cri} is orthogonal to \mathbf{W}_{Cs} but not to \mathbf{W}_{Cr} , that is $\mathbf{w}_{\text{Cri}}^\top \mathbf{w}_{\text{Cs}j} = 0$ but $\mathbf{w}_{\text{Cri}}^\top \mathbf{w}_{\text{Cr}j} \neq 0$.

Step 4. Compute $r^2(\mathbf{w}_{\text{Cri}}, \mathbf{y}_{\text{C}})$ by

$$r^2(\mathbf{w}_{\text{Cri}}, \mathbf{y}_{\text{C}}) = h_i, \quad i = 1, \dots, q, \quad (38)$$

where

$$h_i = \frac{\mathbf{y}_{\text{C}}^\top \mathbf{w}_{\text{Cri}} \mathbf{w}_{\text{Cri}}^\top \mathbf{y}_{\text{C}}}{\mathbf{w}_{\text{Cri}}^\top \mathbf{w}_{\text{Cri}} \mathbf{y}_{\text{C}}^\top \mathbf{y}_{\text{C}}}. \quad (39)$$

Step 5. Find an i which maximises $r^2(\mathbf{w}_{\text{Cri}}, \mathbf{y}_{\text{C}})$ such that

$$i_{\max} = \arg \max_i r^2(\mathbf{w}_{\text{Cri}}, \mathbf{y}_{\text{C}}). \quad (40)$$

Then, remove $\mathbf{x}_{i_{\max}}$ from \mathbf{X}_{r} , add it into \mathbf{X}_{s} , reduce q by 1, and increase p by 1. After that, return to **Step 2** until $p = t$, when **Output** \mathbf{X}_{s} to complete the feature selection.

The pseudocode of the algorithm is given in Algorithm 1.

The multiple correlation coefficient can be obtained either using the definition (18) or the sum of the SOCCs (20). In the greedy search, the OLS based feature selection method has the computational speed advantage over the definition based feature selection method. The computation complexity of the two feature selection methods can be explicitly compared by the asymptotic upper bound notation O [28, p. 47]. At Step k of the greedy search, the $k - 1$ optimal features have been selected, and the rest of the $n - k + 1$ features are the candidates of the k^{th} optimal feature. The candidate feature matrix is a $N \times k$ matrix composed of the $k - 1$ selected features and a candidate feature. According to the normal equation (17), the definition based feature selection method is dominated by computing the inner product of the $N \times k$ centred candidate feature matrix. The computational complexity of the inner product of one centred candidate matrix is $O(k^2 N)$. There are $n - k + 1$ candidate features, so the complexity for Step k is

$$(n - k + 1)O(k^2 N) = O(k^2 n N). \quad (41)$$

Thus, the overall complexity for t features selection is given by

$$\sum_{k=1}^t O(k^2 n N) = O\left(\sum_{k=1}^t k^2 n N\right) = O(t^3 n N). \quad (42)$$

Algorithm 1: Pseudocode of the OLS based feature selection for binomial classification.

Input: \mathbf{X} , \mathbf{y} , t

Output: \mathbf{X}_s

Centre \mathbf{X} and \mathbf{y} to \mathbf{X}_C and \mathbf{y}_C ;

$p \leftarrow 0$;

while $p < t$ **do**

 Divide \mathbf{X}_C into the selected centred features \mathbf{X}_{Cs} and the remaining centred feature \mathbf{X}_{Cr} ;

if $p = 0$ **then**

$\mathbf{W}_{Cr} \leftarrow \mathbf{X}_{Cr}$, which is composed of $(\mathbf{w}_{Cr1}, \dots, \mathbf{w}_{Crq})$;

else

 Orthogonalise \mathbf{X}_{Cs} to itself to form \mathbf{W}_{Cs} , which is composed of $(\mathbf{w}_{Cs1}, \dots, \mathbf{w}_{Csp})$;

 Orthogonalise \mathbf{X}_{Cr} to \mathbf{W}_{Cs} to form \mathbf{W}_{Cr} , which is composed of $(\mathbf{w}_{Cr1}, \dots, \mathbf{w}_{Crq})$;

end

 Compute $r^2(\mathbf{w}_{Cri}, \mathbf{y}_C)$ by (14);

 Find feature index i_{\max} , such that $r^2(\mathbf{w}_{Cri}, \mathbf{y}_C)$ is maximum with $i \in \{1, \dots, n\}$;

 Select feature $\mathbf{x}_{i_{\max}}$ into \mathbf{X}_s ;

$p \leftarrow p + 1$;

end

For OLS based feature selection, as the SOCCs of the selected features (h_1 to h_{k-1}) have been computed in Step 1 to Step $k - 1$, only the SOCCs of the candidate feature (h_k) is required to compute. Thus, OLS based feature selection is dominated by the classical Gram-Schmidt orthogonalisation process. At Step k of the greedy search, the computational complexity of the orthogonalisation of one candidate feature is $O(kN)$. There are $n - k + 1$ candidate features, so the complexity for Step k is

$$(n - k + 1)O(kN) = O(knN). \quad (43)$$

Thus, the overall complexity for t features selection is given by

$$\sum_{k=1}^t O(knN) = O\left(\sum_{k=1}^t knN\right) = O(t^2 nN). \quad (44)$$

Consequently, compared to the definition based feature selection method, the OLS based feature selection method has a significant computational speed advantage in the greedy search.

5. OLS based fast feature selection for multinomial classification

If the N instances of \mathbf{X} belong to c classes, where $c \leq N$, and the n columns in \mathbf{X} represent n features, the feature selection problem for the multinomial classification is to find the t features from the n features of \mathbf{X} , which is optimal to classify the N instances into the c classes. Similar to the last section, the c classes can be encoded to certain values to form a response variable. The ordinal encoding is to assign $1, \dots, c$ to the c labels to form a vector \mathbf{y} . Then, the multiple correlation coefficient between the features and \mathbf{y} can be adopted to indicate the goodness of the features for the classification. The c labels can also be encoded to form a matrix \mathbf{Y} using, e.g. c -label dummy encoding (or called one-hot encoding), $c - 1$ -label dummy encoding, effects encoding, and contrast encoding [27, Chapter 5]. When the response is encoded as a matrix \mathbf{Y} , the canonical correlation coefficients between \mathbf{X} and \mathbf{Y} can be used as the feature selection criterion. Similar to the last section, an example of the greed search for multinomial classification is illustrated in Table 3, where the response is encoded as an $N \times c - 1$ matrix \mathbf{Y} and the ranking criterion is the sum of the squared canonical correlation coefficients.

Table 3. An example for selecting three features from n features by the greedy search for multinomial classification, where $i = 1, \dots, n$ for step 1, $i = 1, 2, 4, 5, \dots, n$ for step 2, $i = 1, 2, 4, 6, 7, \dots, n$ for step 3.

	Ranking Criterion	Selected Feature
1	$\sum_{k=1}^{1 \wedge c-1} R_k^2(\mathbf{x}_3, \mathbf{Y}) \geq \sum_{k=1}^{1 \wedge c-1} R_k^2(\mathbf{x}_i, \mathbf{Y})$	\mathbf{x}_3
2	$\sum_{k=1}^{2 \wedge c-1} R_k((\mathbf{x}_3, \mathbf{x}_5), \mathbf{Y}) \geq \sum_{k=1}^{2 \wedge c-1} R_k((\mathbf{x}_3, \mathbf{x}_i), \mathbf{Y})$	$\mathbf{x}_3, \mathbf{x}_5$
3	$\sum_{k=1}^{3 \wedge c-1} R_k((\mathbf{x}_3, \mathbf{x}_5, \mathbf{x}_1), \mathbf{Y}) \geq \sum_{k=1}^{3 \wedge c-1} R_k((\mathbf{x}_3, \mathbf{x}_5, \mathbf{x}_i), \mathbf{Y})$	$\mathbf{x}_3, \mathbf{x}_5, \mathbf{x}_1$

In the following, the relationship between CCA and a classical LDA classifier will first be demonstrated to reveal the implication of the canonical correlation coefficient in a classification problem. Then, the algorithm of the OLS based feature selection for multinomial classification is developed where the sum of the squared canonical correlation coefficients will be used as the feature ranking criterion. After that, a version of the algorithm that can be used to deal with categorical features is presented.

5.1. Relationship with linear discriminant analysis

As the feature selection is used for the multinomial classification, it is reasonable to know the performance of the features in the Linear Discriminant Analysis (LDA), where the label encoding is not required.

For the convenience of LDA, the feature matrix \mathbf{X} is decomposed into $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(c)}$, where the $N_i \times n$ matrix $\mathbf{X}^{(i)}$ represents the N_i instances belonged to the i^{th} class. The within-class scatter matrix of the samples is

$$\mathbf{S}_w = \sum_{i=1}^c \left(\mathbf{X}^{(i)} - \mathbf{1}^{(i)} \bar{\mathbf{x}}^{(i)\top} \right) \left(\mathbf{X}^{(i)} - \mathbf{1}^{(i)} \bar{\mathbf{x}}^{(i)\top} \right)^\top, \quad (45)$$

where $\bar{\mathbf{x}}^{(i)}$ is the sample mean of each feature in $\mathbf{X}^{(i)}$ given by

$$\bar{\mathbf{x}}^{(i)} = (\bar{x}_1^{(i)}, \dots, \bar{x}_n^{(i)})^\top \quad (46)$$

and $\mathbf{1}^i$ is $N_i \times 1$ vector of ones. The between-class scatter matrix of the samples is

$$\mathbf{S}_b = \sum_{i=1}^c N_i \left(\bar{\mathbf{x}}^{(i)} - \bar{\mathbf{x}} \right) \left(\bar{\mathbf{x}}^{(i)} - \bar{\mathbf{x}} \right)^\top, \quad (47)$$

where $\bar{\mathbf{x}}$ is the overall sample mean of each feature. The aim of LDA is to find a projection direction \mathbf{d} for \mathbf{X} , so that the ratio between the projected between-class scatter and the projected within-class scatter is maximised. The ratio is called Fisher's criterion, which is given by

$$J = \frac{\mathbf{d}^\top \mathbf{S}_b \mathbf{d}}{\mathbf{d}^\top \mathbf{S}_w \mathbf{d}}. \quad (48)$$

The larger Fisher's criterion J implies the better the separation of the c classes. The LDA can be transformed to the eigenvalue problem given by [11, p. 246]

$$\mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{d} = J \mathbf{d}, \quad (49)$$

where the eigenvector is the optimal projection direction \mathbf{d} and the eigenvalue is the maximised Fisher's criterion J .

The relationship between LDA and CCA can be found when \mathbf{Y} is formed by c or $c - 1$ -label dummy encoding. Under the two encoding schemes, the eigenvalue problem (26a) can be rewritten as [29]

$$(\mathbf{S}_b + \mathbf{S}_w)^{-1} \mathbf{S}_b \mathbf{a} = R^2(\mathbf{X}, \mathbf{Y}) \mathbf{a}, \quad (50)$$

or in the form of

$$\mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{a} = \frac{R^2(\mathbf{X}, \mathbf{Y})}{1 - R^2(\mathbf{X}, \mathbf{Y})} \mathbf{a}. \quad (51)$$

Comparing (49) and (51), it is found that LDA and CCA are equivalent, and Fisher's criterion of LDA can be evaluated by

$$J = \frac{R^2(\mathbf{X}, \mathbf{Y})}{1 - R^2(\mathbf{X}, \mathbf{Y})}. \quad (52)$$

indicating that CCA has a close relationship with the Fisher's criterion.

To better reveal this relationship, the $c - 1$ -label dummy encoding and the canonical correlation coefficients are adopted, instead of the ordinal encoding and the multiple correlation coefficient. The $c - 1$ -label dummy encoding constructs a $N \times c - 1$ matrix $Y = (y_{i,j})$, where

$$y_{i,j} = \begin{cases} 1 & i^{\text{th}} \text{ instance is belonged to } j^{\text{th}} \text{ class} \\ 0 & \text{otherwise.} \end{cases} \quad (53)$$

Thus, the dummy response in the last section is a special case of $c - 1$ -label dummy encoding where $c = 2$.

5.2. OLS based feature selection algorithm

The algorithm of OLS based feature selection for multinomial classification can be summarised in 5 steps.

Input:

X: $N \times n$ matrix containing N instances and n features.

Y: $N \times c - 1$ matrix formed by $c - 1$ -label dummy encoding.

t: The number of features to be selected.

Step 1. Centre **Y** into \mathbf{Y}_C , orthogonalise \mathbf{Y}_C into \mathbf{V}_C , centre **X** into \mathbf{X}_C , and let $p = 0$.

Step 2. Divide **X** into $(\mathbf{X}_s, \mathbf{X}_r)$, where the selected feature matrix is given by

$$\mathbf{X}_s = (\mathbf{x}_{s1}, \dots, \mathbf{x}_{sp}), \quad (54)$$

and the remaining feature matrix is given by

$$\mathbf{X}_r = (\mathbf{x}_{r1}, \dots, \mathbf{x}_{rq}), \quad (55)$$

where p is the number of the selected features, and q is the number of the remaining features. Correspondingly, divide \mathbf{X}_C into $(\mathbf{X}_{Cs}, \mathbf{X}_{Cr})$, where

$$\begin{aligned} \mathbf{X}_{Cs} &= (\mathbf{x}_{Cs1}, \dots, \mathbf{x}_{Csp}) \\ \mathbf{X}_{Cr} &= (\mathbf{x}_{Cr1}, \dots, \mathbf{x}_{Crq}). \end{aligned} \quad (56)$$

Step 3. If $p = 0$, let

$$\begin{aligned}\mathbf{W}_{Cr} &= \mathbf{X}_{Cr} \\ \mathbf{w}_{Cri} &= \mathbf{x}_{Cri}, \quad i = 1, \dots, q.\end{aligned}\tag{57}$$

Otherwise, first, orthogonalise \mathbf{X}_{Cs} into \mathbf{W}_{Cs} , where

$$\mathbf{W}_{Cs} = (\mathbf{w}_{Cs1}, \dots, \mathbf{w}_{Csp}),\tag{58}$$

and $\mathbf{w}_{Csi}^\top \mathbf{w}_{Csj} = 0$ for $i \neq j$. Then, orthogonalise each feature in \mathbf{X}_{Cr} to \mathbf{W}_{Cs} to form the matrix \mathbf{W}_{Cr} , where

$$\mathbf{W}_{Cr} = (\mathbf{w}_{Cr1}, \dots, \mathbf{w}_{Crq}),\tag{59}$$

and \mathbf{w}_{Cri} is obtained through the classical Gram-Schmidt process, which is given by

$$\mathbf{w}_{Cri} = \mathbf{x}_{Cri} - \sum_{j=1}^p \frac{\mathbf{x}_{Cri}^\top \mathbf{w}_{Csj}}{\mathbf{w}_{Csj}^\top \mathbf{w}_{Csj}} \mathbf{w}_{Csj}, \quad i = 1, \dots, q.\tag{60}$$

It should be noticed that \mathbf{w}_{Cri} is orthogonal to \mathbf{W}_{Cs} but not to \mathbf{W}_{Cr} , that is $\mathbf{w}_{Cri}^\top \mathbf{w}_{Csj} = 0$ but $\mathbf{w}_{Cri}^\top \mathbf{w}_{Crj} \neq 0$.

Step 4. Compute $R^2(\mathbf{w}_{Cri}, \mathbf{V}_C)$ by

$$R^2(\mathbf{w}_{Cri}, \mathbf{V}_C) = \sum_{j=1}^{c-1} h_{i,j}, \quad i = 1, \dots, q,\tag{61}$$

where

$$h_{i,j} = \frac{\mathbf{v}_{Cj}^\top \mathbf{w}_{Cri} \mathbf{w}_{Cri}^\top \mathbf{v}_{Cj}}{\mathbf{w}_{Cri}^\top \mathbf{w}_{Cri} \mathbf{v}_{Cj}^\top \mathbf{v}_{Cj}}.\tag{62}$$

Step 5. Find an i which maximises $R^2(\mathbf{w}_{Cri}, \mathbf{V}_C)$ such that

$$i_{\max} = \arg \max_i R^2(\mathbf{w}_{Cri}, \mathbf{V}_C).\tag{63}$$

Then, remove $\mathbf{x}_{i_{\max}}$ from \mathbf{X}_r , add it into \mathbf{X}_s , reduce q by 1, and increase p by 1. After that, return to **Step 2** until $p = t$, when **Output** \mathbf{X}_s to complete the feature selection.

The pseudocode of the algorithm is given in Algorithm 2.

The speed advantage of the OLS based feature selection method is reflected in **Step 4**. To evaluate the goodness of the candidate feature \mathbf{x}_{ri} , CCA requires to compute the canonical correlation coefficient

Algorithm 2: Pseudocode of the OLS based feature selection for multinomial classification.

Input: $\mathbf{X}, \mathbf{Y}, t$

Output: \mathbf{X}_s

Centre \mathbf{X} and \mathbf{Y} to \mathbf{X}_C and \mathbf{Y}_C ;

Orthogonalise \mathbf{Y}_C to itself to form \mathbf{V}_C ;

$p \leftarrow 0$;

while $p < t$ **do**

 Divide \mathbf{X}_C into the selected centred features \mathbf{X}_{Cs} and the remaining centred feature \mathbf{X}_{Cr} ;

if $p = 0$ **then**

$\mathbf{W}_{Cr} \leftarrow \mathbf{X}_{Cr}$, which is composed of $(\mathbf{w}_{Cr1}, \dots, \mathbf{w}_{Crq})$;

else

 Orthogonalise \mathbf{X}_{Cs} to itself to form \mathbf{W}_{Cs} , which is composed of $(\mathbf{w}_{Cs1}, \dots, \mathbf{w}_{Csp})$;

 Orthogonalise \mathbf{X}_{Cr} to \mathbf{W}_{Cs} to form \mathbf{W}_{Cr} , which is composed of $(\mathbf{w}_{Cr1}, \dots, \mathbf{w}_{Crq})$;

end

 Compute $R^2(\mathbf{w}_{Cr}, \mathbf{V}_C)$ by (27);

 Find feature index i_{\max} , such that $R^2(\mathbf{w}_{Cr}, \mathbf{V}_C)$ is maximum with $i \in \{1, \dots, n\}$;

 Select feature $\mathbf{x}_{i_{\max}}$ into \mathbf{X}_s ;

$p \leftarrow p + 1$;

end

$R((\mathbf{X}_s, \mathbf{x}_{ri}), \mathbf{Y})$, while OLS only needs to compute the multiple correlation coefficient $R(\mathbf{w}_{Cr}, \mathbf{V}_C)$, because

$$\begin{aligned} \sum_{k=1}^{p+1 \wedge c-1} R_k^2((\mathbf{X}_s, \mathbf{x}_{ri}), \mathbf{Y}) &= \sum_{k=1}^{p+1 \wedge c-1} R_k^2((\mathbf{W}_{Cs}, \mathbf{w}_{Cr}), \mathbf{V}_C) \\ &= \sum_{k=1}^{p \wedge c-1} R_k^2(\mathbf{W}_{Cs}, \mathbf{V}_C) + R^2(\mathbf{w}_{Cr}, \mathbf{V}_C). \end{aligned} \tag{64}$$

For each candidate feature \mathbf{x}_{ri} , $R(\mathbf{W}_{Cs}, \mathbf{V}_C)$ is the same. Thus, to find the maximal $R((\mathbf{X}_s, \mathbf{x}_{ri}), \mathbf{Y})$, only $R(\mathbf{w}_{Cr}, \mathbf{V}_C)$ is required to compute. In addition, although the multiple correlation coefficient $R(\mathbf{w}_{Ci}, \mathbf{Y})$, which is equal to $R(\mathbf{w}_{Ci}, \mathbf{V}_C)$, can be computed through the definition (18), OLS provides a faster way of computation. Equation (18) requires to solve the normal equation which is dominated by the inner product of the $N \times c - 1$ matrix \mathbf{Y}_C , whose computational complexity is $O(c^2 N)$. For OLS, as the orthogonalisation of \mathbf{Y} is only required once in **Step 1**, the dominant part of computation is (61) whose computational complexity

is only $\mathcal{O}(cN)$.

As the above introduction of the OLS based algorithm is basically conceptual, some speed optimisation steps have been omitted. For example, in **Step 3**, \mathbf{W}_{Cs} computed for selecting the i^{th} optimal feature can be reused for selecting the $i + 1^{\text{th}}$ optimal feature. The further optimisation of the OLS speed can be found in the original paper of OLS based model term selection [10].

5.3. Dealing with categorical features

When the features are categorical, the feature encoding is required for OLS based feature selection. In the previous analysis, n features are represented by n column vectors in \mathbf{X} , but some encoding methods may encode the categorical features into matrices. In these cases, the feature matrix is composed of n submatrices, that is

$$\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n), \quad (65)$$

where the matrix \mathbf{X}_i is the encoded i^{th} feature. An OLS based feature selection algorithm similar to the algorithm in Section 5.2 can be applied to the matrix encoded features. The candidate orthogonal feature matrix in **Step 3** in Section 5.2 is given by

$$\mathbf{W}_{Cr} = (\mathbf{W}_{Cr1}, \dots, \mathbf{W}_{Crq}), \quad (66)$$

where

$$\mathbf{W}_{Cri} = (\mathbf{w}_{Cri}^{[1]}, \dots, \mathbf{w}_{Cri}^{[z_i]}). \quad (67)$$

is a $N \times z_i$ matrix. Besides being orthogonal to the selected orthogonal feature matrix \mathbf{W}_{Cs} , the submatrix \mathbf{W}_{Cri} should be column-wise orthogonal via an additional orthogonalisation process. In **Step 4**, the sum of the squared canonical correlation coefficients can be computed by

$$\sum_{k=1}^{z_i \wedge c-1} R_k^2(\mathbf{W}_{Cri}, \mathbf{V}_C) = \sum_{j=1}^{c-1} \sum_{g=1}^{z_i} h_{i,j}^{[g]}, \quad i = 1, \dots, q, \quad (68)$$

where

$$h_{i,j}^{[g]} = \frac{\mathbf{v}_{Cj}^\top \mathbf{w}_{Cri}^{[g]} \mathbf{w}_{Cri}^{[g]\top} \mathbf{v}_{Cj}}{\mathbf{w}_{Cri}^{[g]\top} \mathbf{w}_{Cri}^{[g]} \mathbf{v}_{Cj}^\top \mathbf{v}_{Cj}}. \quad (69)$$

Finally, the sum of the squared canonical correlation coefficients are used to rank the features for **Step 5**.

6. Empirical study

In this section, firstly, a simple example is used to illustrate the procedure of the OLS based feature selection method when applied to the Fisher’s iris data [30]. The relationship between the SOCCs with canonical correlation coefficient and Fisher’s criterion is demonstrated via this case study. Then, the OLS based feature selection methods are compared with mutual information based filter methods and the embedded methods using both synthetic and real world datasets. For the filter methods, the features are selected via greedy search with different ranking criteria. The mutual information based feature selection methods in this comparison are summarised in [12]. The ranking criteria are the difference and quotient schemes (mRMRd and mRMRq), Mutual Information Maximisation (MIM), Joint Mutual Information (JMI), Conditional Mutual Information Maximisation (CMIM), Conditional Infomax Feature Extraction (CIFE), Interaction Capping (ICAP), and Double Input Symmetrical Relevance (DISR). The embedded methods are LASSO and elastic net (Net), and the MATLAB function `lasso` is adopted.

The OLS based method shows superiority in computation speed compared to the direct use of the definition of the canonical correlation coefficient. For example, in the two real world datasets (Dexter and Gisette), the OLS method takes 401 ms for Dexter and 5109 ms for Gisette to select 20 features on a 2.6 GHz personal laptop, while the traditional definition based method takes 13200 ms and 134922 ms, respectively. The empirical studies are implemented in MATLAB R2021a, and the code will be published in GitHub².

6.1. An illustration of the OLS based feature selection

The Fisher’s iris data are given in Table 4. The 7 instances have 4 features and 3 classes, so $N = 7$, $n = 4$, and $c = 3$. The objective of the feature selection is to find 3 optimal features for the 3 species classification.

The feature matrix is \mathbf{X} and the $c - 1$ -label dummy encoded response is

$$\mathbf{Y} = \begin{pmatrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 \end{pmatrix}^T, \quad (70)$$

where (1, 0) represents setosa, (0, 1) represents versicolor, and (0, 0) represents virginica. Following the algorithm introduced in Section 5.2, the procedure of the OLS based feature selection method is shown below.

²https://github.com/MatthewSZhang/fs_ols

Table 4. Fisher's Iris Dataset.

Sepal Length	Sepal Width	Petal Length	Petal Width	Species
5.1	3.5	1.4	0.2	setosa
4.9	3	1.4	0.2	setosa
7	3.2	4.7	1.4	versicolor
6.4	3.2	4.5	1.5	versicolor
6.3	3.3	6	2.5	virginica
5.8	2.7	5.1	1.9	virginica
7.1	3	5.9	2.1	virginica

Step 1. First, centre \mathbf{Y} into \mathbf{Y}_C . Second, orthogonalise \mathbf{Y}_C into \mathbf{V}_C . Through the classical Gram-Schmidt process, use the first column of \mathbf{Y}_C as \mathbf{v}_{C1} , then orthogonalise the second column to the first column. Thus, the centred orthogonalised response matrix is given by

$$\mathbf{V}_C = \begin{pmatrix} 0.7143 & 0.7143 & -0.2857 & -0.2857 & -0.2857 & -0.2857 & -0.2857 \\ 0.0000 & 0.0000 & 0.6000 & 0.6000 & -0.4000 & -0.4000 & -0.4000 \end{pmatrix}^T. \quad (71)$$

Third, centre \mathbf{X} into \mathbf{X}_C .

Step 2. As no feature has been selected, \mathbf{X}_s is empty and \mathbf{X}_r is the same as \mathbf{X} . Correspondingly, \mathbf{X}_{Cs} is empty and \mathbf{X}_{Cr} is the same as \mathbf{X}_C .

Step 3. In this step, the centred features in \mathbf{X}_{Cr} are required to be orthogonalised to \mathbf{W}_{Cs} . As no feature has been selected, let \mathbf{W}_{Cr} equal to \mathbf{X}_{Cr} .

Step 4. The multiple correlation coefficients between \mathbf{w}_{Cr_i} and \mathbf{V}_C are 0.7628, 0.2264, 0.9779, and 0.9604.

Step 5. The third feature (i.e. petal length) has the highest multiple correlation. Thus, the petal length is selected into \mathbf{X}_s , and the features contained in \mathbf{X}_r in order are sepal length, sepal width, and petal width.

Step 2. According to the new \mathbf{X}_s and \mathbf{X}_r , the centred matrix \mathbf{X}_C is divided into $(\mathbf{X}_{Cs}, \mathbf{X}_{Cr})$.

Step 3. As only one feature is in \mathbf{X}_{Cs} , let the orthogonalised feature \mathbf{W}_{Cs} equal to \mathbf{X}_{Cs} . Through the classical Gram-Schmidt process, the features in \mathbf{X}_{Cr} are orthogonalised to \mathbf{W}_{Cs} .

Step 4. The multiple correlation coefficients between \mathbf{w}_{Cr_i} and \mathbf{V}_C are 0.4458, 0.0841, and 0.4644.

Step 5. The third feature (i.e. petal width) has the highest multiple correlation. Thus, the features contained in \mathbf{X}_s in order are petal length and petal width, and the features contained in \mathbf{X}_r in order are sepal

length and sepal width.

Step 2. According to the new \mathbf{X}_s and \mathbf{X}_r , the centred matrix \mathbf{X}_C is divided into $(\mathbf{X}_{Cs}, \mathbf{X}_{Cr})$.

Step 3. Keep the first column of \mathbf{X}_{Cs} unchanged, and orthogonalise the second column to the first column through the classical Gram-Schmidt process. Each feature in \mathbf{X}_{Cr} is orthogonalised to \mathbf{W}_{Cs} , respectively, to obtain \mathbf{W}_{Cr} .

Step 4. The multiple correlation coefficients between \mathbf{w}_{Cri} and \mathbf{V}_C are 0.0382 and 0.1108.

Step 5. The second feature (i.e. sepal width), which has the highest multiple correlation, is selected into \mathbf{X}_s . Therefore, the 3 selected features are petal length, petal width, and sepal width.

The squared canonical correlation coefficients between the 3 features and \mathbf{Y} are given by

$$\begin{aligned} R_1^2((\mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_2), \mathbf{Y}) &= 0.9905 \\ R_2^2((\mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_2), \mathbf{Y}) &= 0.5626. \end{aligned} \tag{72}$$

In LDA, the within-class scatter matrix is given by

$$\mathbf{S}_w = \begin{pmatrix} 0.5067 & 0.2367 & 0.2700 \\ 0.2367 & 0.1917 & 0.1800 \\ 0.2700 & 0.1800 & 0.3050 \end{pmatrix}, \tag{73}$$

and the between-class scatter matrix is given by

$$\mathbf{S}_b = \begin{pmatrix} 22.4305 & 10.1333 & -1.1886 \\ 10.1333 & 4.6483 & -0.5800 \\ -1.1886 & -0.5800 & 0.0893 \end{pmatrix}. \tag{74}$$

Through solving the eigenvalue problem (49), the Fisher's criteria of LDA are given by

$$\begin{aligned} J_1 &= 104.1481 \\ J_2 &= 1.2864. \end{aligned} \tag{75}$$

From (72) and (75), it can be verified that the relationship between the squared canonical correlation coefficients and Fisher's criterion of LDA is as described by (52). To verify the equality between the squared canonical correlation coefficients and SOCCs in (29), the right hand side of (29) is given by

$$\begin{aligned} R_1^2((\mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_2), \mathbf{Y}) + R_2^2((\mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_2), \mathbf{Y}) &= 0.9905 + 0.5626 \\ &= 1.5531, \end{aligned} \tag{76}$$

and the left hand side of (29) is given by the sum of the maxima in each iteration, that is $0.9779 + 0.4644 + 0.1108 = 1.5531$.

6.2. Application to synthetic data for binomial classification

In this case study, the proposed feature selection method for a binomial classification is investigated. The $N \times n$ feature matrix is sampled from the multivariate normal distribution, which is given by $\mathbf{X} \sim \mathcal{M}(\boldsymbol{\mu}, \boldsymbol{\Sigma}_N)$, where the mean values in the $n \times 1$ vector $\boldsymbol{\mu}$ are sampled from the normal distribution with mean 0 and standard deviation 0.1. The $n \times n$ covariance matrix $\boldsymbol{\Sigma}_N$ is sampled from the Wishart distribution, which is given by $\boldsymbol{\Sigma}_N \sim \mathcal{W}(\boldsymbol{\Sigma}_W, N)/N$, where $\boldsymbol{\Sigma}_W$ is a $n \times n$ diagonal matrix whose main diagonal is uniformly distributed on the interval (0, 1). Let the number of the instances is 600, i.e. $N = 600$, and the number of the candidate features are 100, i.e. $n = 100$. The 5th, 10th, and 15th features are used to construct the dummy response vector \mathbf{y} , which is sampled from the Bernoulli distribution (i.e. 1 trial binomial distribution) given by $\mathbf{y} \sim \mathcal{B}(\boldsymbol{\pi})$, where the probability vector $\boldsymbol{\pi} = (\pi_1, \dots, \pi_N)^\top$ is generated by the binomial logistic regression model, that is

$$\pi_i = \frac{1}{1 + \exp(-(-2x_{i,5} - 3x_{i,10} + 4x_{i,15}))}, \quad i = 1, \dots, N. \quad (77)$$

Given \mathbf{X} and \mathbf{y} , the aim of the feature selection study is to find the 3 correct feature indices (i.e. 5, 10, and 15).

The proposed OLS based feature selection method is compared with the mutual information based filter methods and embedded methods. For the mutual information based features selection methods, the continuous features are discretised into 4 categories by the mean values and the mean values \pm the standard deviation. For the OLS based feature selection, the continuous features are treated in two ways. One (denoted by OLS) implements the algorithm in subsection 5.2 and use continuous features directly. Another one (denoted by OLSd) implements the algorithm in subsection 5.3, where the continuous features are discretised into 4 categories by the mean values and the mean values \pm the standard deviation, and then encoded into matrices by $c - 1$ dummy encoding.

The simulation study is repeated 100 times to check how many times the feature selection methods choose the 3 correct features, and the results are given by Table 5. In this comparison, two OLS based feature selection methods choose the right features 95 times and 88 times, respectively, in the 100 tests, which are higher than what can be achieved by the mutual information based feature selection methods and the embedded methods.

6.3. Application to synthetic data for multinomial classification

In this case study, the feature selection for a 3-class multinomial classification is investigated. The $N \times n$ feature matrix is generated in the same way as in the last subsection. The number of the instances is 900,

Table 5. A comparison with mutual information based feature selection methods in binomial classification.

Method	Times of selecting correct features	Method	Times of selecting correct features	Method	Times of selecting correct features
OLS	95	MIM	73	ICAP	74
OLSd	88	JMI	79	DISR	79
mRMRd	76	CMIM	74	LASSO	85
mRMRq	77	CIFE	76	Net	87

i.e. $N = 900$, and the number of the candidate features are 100, i.e. $n = 100$. We use the 5th, 10th, and 15th features to construct the $N \times 3$ response matrix \mathbf{Y}' , which is c -label dummy encoded. \mathbf{Y}' is sampled from the categorical distribution (i.e. 1 trial multinomial distribution) given by $\mathbf{Y}' \sim C(\mathbf{\Pi})$, where the $N \times 3$ probability matrix $\mathbf{\Pi} = (\pi_{i,j})$ is composed of the probability vector for each class, that is $\mathbf{\Pi} = (\boldsymbol{\pi}_1, \boldsymbol{\pi}_2, \boldsymbol{\pi}_3)$. The probability vectors are generated by the multinomial logistic regression model [31, p. 270]. First, the probability ratios are given by

$$\begin{aligned} \frac{\pi_{i,1}}{\pi_{i,3}} &= \exp(-x_{i,5} - x_{i,10} + x_{i,15}) \\ \frac{\pi_{i,2}}{\pi_{i,3}} &= \exp(x_{i,5} - x_{i,10} - x_{i,15}), \quad i = 1, \dots, N. \end{aligned} \quad (78)$$

Second, the probability of π_3 is given by

$$\pi_{i,3} = \frac{1}{1 + \frac{\pi_{i,1}}{\pi_{i,3}} + \frac{\pi_{i,2}}{\pi_{i,3}}}, \quad i = 1, \dots, N. \quad (79)$$

Finally, $\boldsymbol{\pi}_1$ and $\boldsymbol{\pi}_2$ can be computed by substituting (79) into (78). To make the response matrix become $c - 1$ -label dummy encoded, the first column of \mathbf{Y}' is removed to form \mathbf{Y} , which becomes a $N \times 2$ matrix containing only 0 and 1. Given \mathbf{X} and \mathbf{Y} , the aim of the feature selection simulation is to find the 3 correct feature indices (i.e. 5, 10, and 15).

The task is repeated 100 times, and the number of times when a correct feature selection is achieved is shown in Table 6. Two OLS based methods still give the best results, especially OLS which uses the continuous features.

Table 6. A comparison with mutual information based feature selection methods in multinomial classification.

Method	Times of selecting correct features	Method	Times of selecting correct features	Method	Times of selecting correct features
OLS	92	MIM	82	ICAP	82
OLSd	84	JMI	80	DISR	80
mRMRd	83	CMIM	82	LASSO	79
mRMRq	84	CIFE	67	Net	79

6.4. Application to the datasets of NIPS feature selection challenge

Two datasets from the NIPS feature selection challenge³ are used for the feature selection methods evaluation. The detail of the datasets are illustrated in Table 7. Dexter dataset is from Reuters text categorisation task and Gisette dataset is from a handwriting recognition task. Both of the datasets have 2 classes. The features of the datasets are composed of real features and artificial features (called probes). As the probes do not carry information of the class labels, the desirable feature selection methods should avoid selecting them. The datasets are divided into training, validation, and test data. The labels of the test data are withheld by the data providers, and the performance on the test data are obtained by uploading the results to the challenge website.

Table 7. Summary of the NIPS feature selection challenge datasets.

Name	Feature (Real/Probe)	Train/Validation/Test
Dexter	20000 (9947/10053)	300/300/2000
Gisette	5000 (2500/2500)	6000/1000/6500

The feature values in both datasets are quantised to 1000 levels, and the features are treated as continuous. For the mutual information based methods, the 1000 levels are discretised into 10 equal width bins. For OLSd, the discretised features are encoded into matrices by $c - 1$ dummy encoding. For OLS, the continuous features are used directly.

The experiment is implemented in the following steps. First, the optimal features are selected by different methods using the training data. Then, given the selected optimal features, a linear Support Vector Machine

³<https://competitions.codalab.org/competitions/3931>

(SVM) is trained with the training data. Finally, the prediction results are generated by the SVM model on the training, validation, and test data, respectively. The classification performance is evaluated by the Area Under the Curve (AUC) of the Receiver Operating Characteristic (ROC) curve. Each method selects 20 optimal features. The AUC results on the training and validation data are shown in Fig. 1 and Fig. 2 for the Dexter and Gisette datasets, respectively. Generally, OLS which uses the continuous features gives the best classification performance. In Dexter dataset, OLS shows strikingly better results than other methods. The results on test data are given in Table 8. Although OLS method selects 1 probe in Dexter dataset, the rest of 19 real features (especially the first 8 features according to Fig. 1b) selected by OLS are more informative for classification than 20 real features selected by other methods, showing that OLS method can achieve the best AUC results.

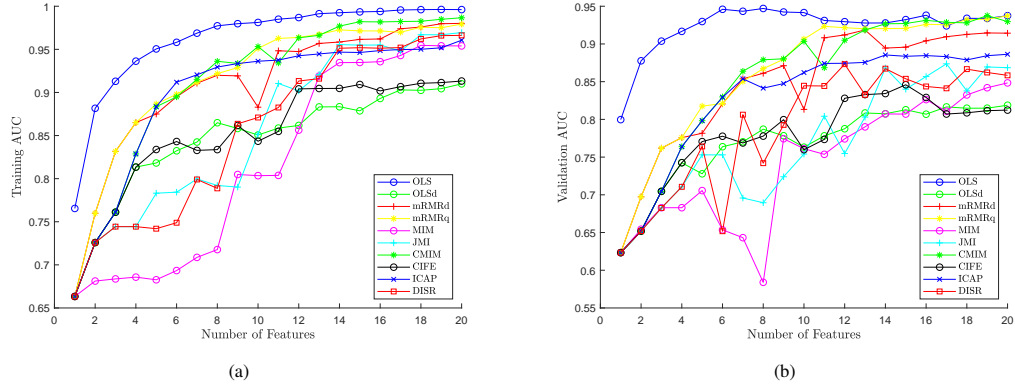


Fig. 1. AUC results of the feature selection methods on (a) training and (b) validation Dexter dataset.

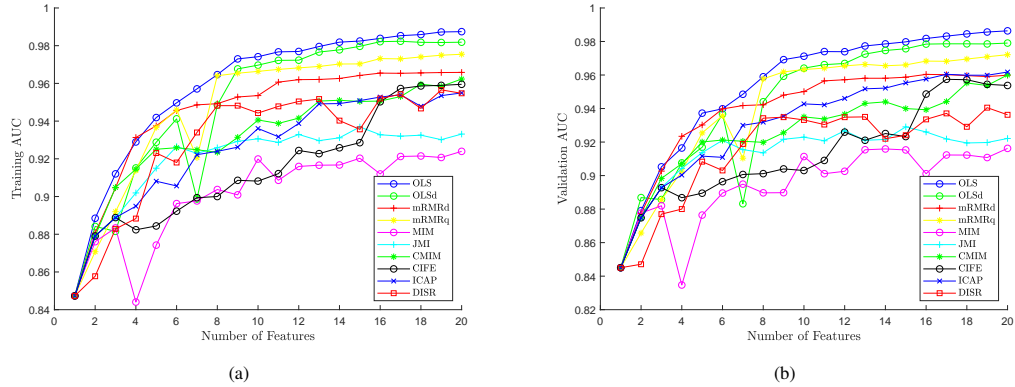


Fig. 2. AUC results of the feature selection methods on (a) training and (b) validation Gisette dataset.

Table 8. Results on the NIPS feature selection challenge test data.

		OLS	OLSd	mRMRd	mRMRq	MIM	JMI	CMIM	CIFE	ICAP	DISR
Dexter	AUC	0.9551	0.8413	0.9246	0.9355	0.8774	0.8917	0.9444	0.8367	0.8848	0.8908
	Probe	1	4	0	1	0	0	1	14	2	0
Gisette	AUC	0.9873	0.9824	0.9662	0.9776	0.9324	0.9352	0.9667	0.9605	0.9580	0.9490
	Probe	0	0	0	0	0	0	0	0	0	0

6.5. Application to the MNIST dataset of handwritten digit

The MNIST dataset ⁴ has a training set of 60000 instances and a test set of 10000 instances, which are categorised into 10 classes (i.e. 0 to 9 digits). The data is balanced and each class has around 6000 instances for training set and 1000 for test set. The instance is a 28×28 pixel box containing grey levels ranging from 0 to 255. The 784 pixels are features for the classification. The aim of the feature selection is to choose 50 optimal features for a LDA classifier. The performance of the feature selection methods is evaluated by the classification accuracy (ACC).

For the mutual information based methods and the OLSd method, the continuous grey level features are discretised into 0 and 1 using the method given in [32]. The OLS method uses the continuous feature directly. The hyperparameters of the embedded methods are optimised by the grid search. The feature selection and LDA classifier training are carried out in the training datasets, and tested in the test dataset. The results of the ACC for the training dataset and the test dataset are shown in Fig. 3. The performance of the feature selection methods is close at the beginning, where only a few features are selected. The ACC of training and test datasets starts divided into two groups after more than 14 features are selected. The OLS, OLSd, mRMRq, CMIM, ICAP, and CIFE methods are in the first tier and the rest of the methods are in the second tier. When 50 optimal features are selected, the OLS method is in the first place and the OLSd method is in the second place for both the training and test sets.

The 50 pixels selected by each feature selection method are compared in Fig. 4. All the feature selection criteria can indicate the pixels in the centre area is relevant to the digit classification, while the pixels in the corners are irrelevant. However, it is found that some mutual information based methods and the embedded methods tend to select the blocks of neighbouring pixels, while the proposed methods tend to select the pixels more spread over the entire picture. As the neighbouring pixels normally contain the redundant information,

⁴<http://yann.lecun.com/exdb/mnist/>

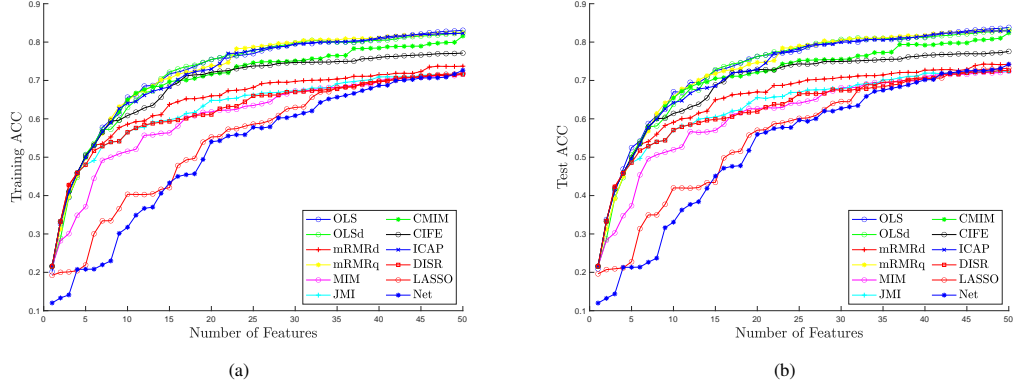


Fig. 3. ACC results of the feature selection methods on (a) training and (b) test MNIST dataset.

the criteria without redundancy control, e.g. MIM [12], will improperly select the relevant but redundant pixels, which are as shown in Fig. 4. Therefore, the results imply that, in this case, the proposed OLS methods provides best redundancy control.

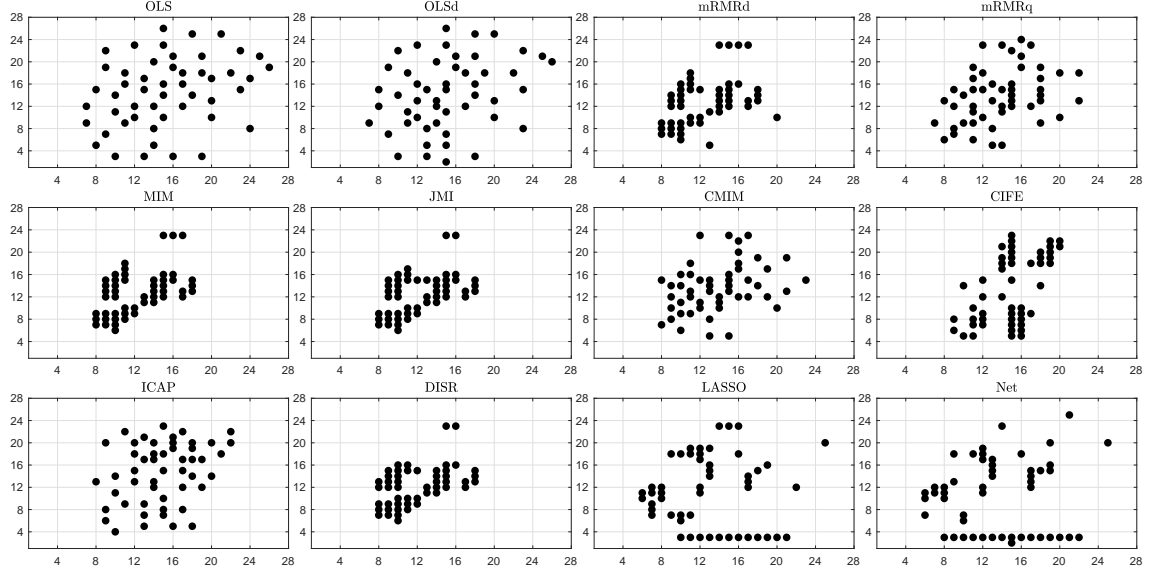


Fig. 4. The comparison of the selected pixels when using different feature ranking criteria.

6.6. Application to the UCI datasets

The four real world datasets are from the UCI machine learning repository ⁵ [33]. The detailed information of the datasets are shown in Table 9.

Table 9. Summary of the UCI datasets.

Name	Mfeat	Breast	CNAE	Lymph
Data Type	Continuous	Continuous	Discrete	Categorical/Discrete
No. of instances	2000	569	1080	142
No. of Features	649	30	856	18
No. of Classes	10	2	9	2

The multi-feature digit (Mfeat) dataset [34] consists of handwritten numerals (0 to 9) extracted from a collection of Dutch utility maps. Each class has 200 instances, and there are 2000 instances in total. The instances are originally scanned to produce the images of 8 bits greyscale. 649 continuous features are extracted from the raw images, including 76 Fourier coefficients, 216 profile correlations, 64 Karhunen-Loève coefficients, 240 pixel averages in 2×3 windows, 47 Zernike moments, and 6 morphological features.

The Wisconsin diagnostic breast cancer (Breast) dataset [35] has 569 instances are classified into malignant (212) or benign (357) breast cancer. The original data are from 10 measurements, including the cell nucleus' radius, texture, perimeter, etc. The 30 features are constructed by taking the mean, standard error, and largest (mean of the three largest values) of these measurements.

CNAE-9 (CNAE) dataset [36] contains 1080 documents of business descriptions of Brazilian companies, which are categorised into 9 economic activities. The dataset is balanced, and each activity has 120 documents. The frequencies of the 856 words in the document are used as features. The feature matrix are highly sparse and 99.22% of the matrix is filled with zeros.

Lymphography (Lymph) dataset [37] contains 148 instances in 4 classes, i.e. normal find (2), metastases (81), malign lymph (61), and fibrosis (4). To make the dataset balanced, only the instances belong to metastases and malign lymph are used. The 18 medical diagnostic attributes form categorical and discrete features.

For the continuous features, the mutual information based feature selection methods use the discretised features which are categorised into 3 categories by the mean values \pm the standard deviations. The OLSd

⁵<https://archive.ics.uci.edu/ml/datasets.php>

method uses the dummy features which encode the discretised features into matrices via $c - 1$ dummy encoding, while the OLS method uses the continuous features directly. For the categorical/discrete features, the mutual information based methods can directly use them, and OLSd use $c - 1$ dummy encoding features. The OLS method uses the discrete feature directly, and adopts the categorical features which are encoded into vectors by the ordinal encoding. The hyperparameters of the LASSO and the elastic net are optimised by the grid search. The 10-fold cross validation for the ACC of the LDA classifier is applied. The mean and standard deviation are extracted from the training data for the continuous feature discretisation. The feature selection and the LDA classifier training are implemented on the 10 training datasets. The feature selection performance is evaluated by the average ACC on the training and validation datasets, which are shown in Fig. 5 to Fig. 8.

For the Mfeat dataset, except that the mRMRq method gives a significantly worse ACC when the 3rd to 8th features are selected, the other methods gives similar results. When 20 features are selected, the OLS method achieves the highest ACC in the training set and the second highest ACC in the validation set, while the OLSd method achieves the second highest ACC in the training set and the highest ACC in the validation set.

For the Breast dataset, the ACC of the feature selection methods varies in the small range between 0.91 and 0.97. The OLS gives the best ACC in the training set from the 2nd to 20th feature, and the validation set from the 2nd to 8th features.

For the CNAE dataset, except that the CIFE method gives a significantly worse ACC, the other mutual information based methods give similar results. The OLS gives the best ACC in the training set from the 19th to 50th features. When 50 features are selected, the top 5 methods in the test set are CMIM, mRMRq, DISR, OLS, OLSd.

For the Lymph dataset, except that OLS gives a slightly better ACC in the training set when the first 4 features are selected and the 8th to 12th features are selected, no method is evidently stronger or weaker than other methods.

It can be seen that no feature selection method ranks first in all four datasets. However, the proposed OLS and OLSd are generally in the top five methods, and the methods do not show a significant weakness in any of the datasets. In addition, as no discretisation is required, the proposed OLS method is more convenient than the mutual information based methods when the features are continuous. The embedded methods generally give worse ACC comparing with the filter methods in the first few selected features. The reason can be that the filter methods select the most important features first, while the embedded methods

select the feature together without ranking them and the features are merely sorted by the absolute values of the features' fitted coefficients.

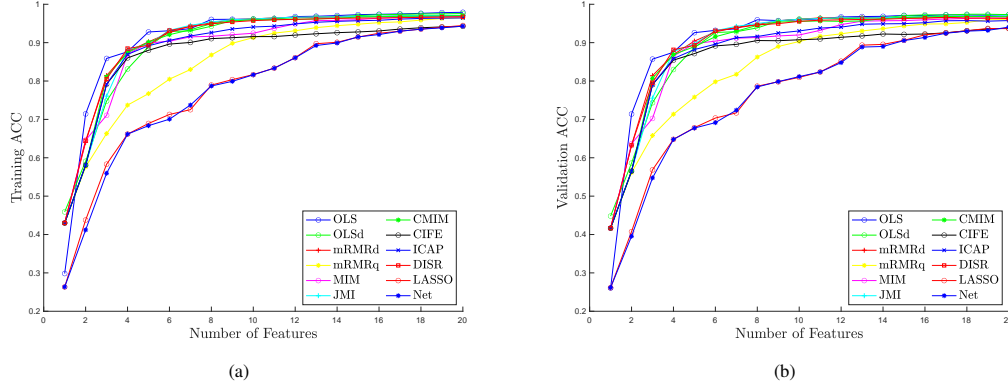


Fig. 5. The average ACC results of the 10-fold cross validation for the feature selection methods on (a) training and (b) validation Mfeat dataset.

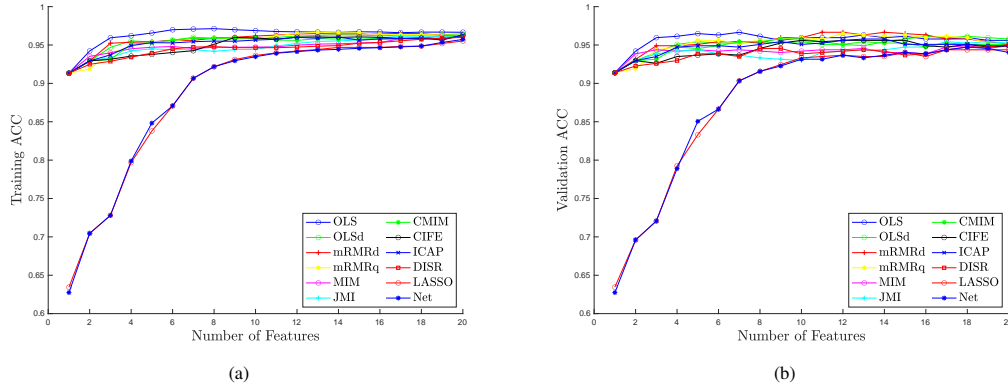


Fig. 6. The average ACC results of the 10-fold cross validation for the feature selection methods on (a) training and (b) validation Breast dataset.

6.7. Comparison with the most recent methods

The two most recent methods are compared with the proposed OLS method. The ranking criteria of the first method is OBCC [26], which has been briefly described in the introduction. The second is the Shapley value based feature selection method [38, 39]. The Shapley value is firstly introduced by Lloyd S. Shapley as a concept of coalitional game theory [40]. Recently, it has been developed as a unified approach

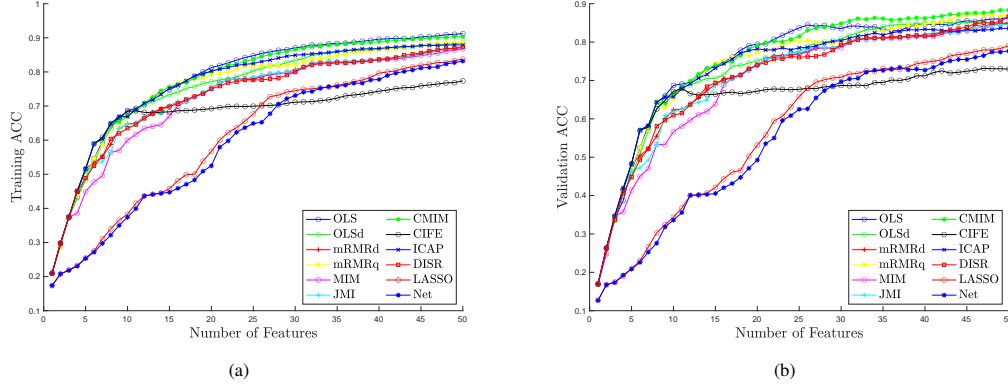


Fig. 7. The average ACC results of the 10-fold cross validation for the feature selection methods on (a) training and (b) validation CNAE dataset.

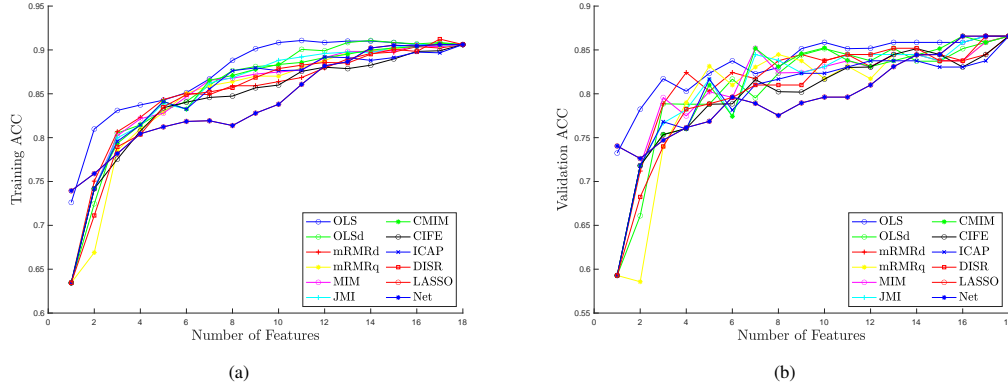


Fig. 8. The average ACC results of the 10-fold cross validation for the feature selection methods on (a) training and (b) validation Lymph dataset.

to the feature importance analysis for machine learning models [38], which can also be used for feature selection [39]. One reason of the popularity of the Shapley value is it has some desirable properties, such as additivity and uniqueness [38]. In this paper, the ranking criterion of the Shapley value based features selection method, denoted as SHAP, is the sum of absolute Shapley values for all the training data and classes with a LDA model, which we use `MATLAB` function `shapley` to compute. As the computation speed of SHAP is slow for the large dataset, only the small datasets Breast and Lymph in the last subsection and two new small datasets of Flag and Parkinsons are used for this comparison.

The Flag dataset [41] contains details of 194 nations and their flags, which are classified into 8 religions including Catholic, other Christian, Muslim, Buddhist, Hindu, Ethnic, Marxist, and others. The information

forms 28 categorical and discrete features including population of the nation, number of different colours in the flag, number of vertical bars in the flag, etc, where the categorical features are ordinal encoded.

The Parkinsons dataset [42] is composed of 195 biomedical voice measurements from people, where 147 measurements are from people with Parkinson’s disease and 48 measurements are from healthy people. The 22 continuous features including average, maximum and minimum vocal fundamental frequencies, etc.

All three methods, i.e. OLS, OBCC and SHAP, can be used in both categorical and numerical features. The 10-fold cross validation for the ACC of the LDA classifier is applied. The feature selection and the LDA classifier training are implemented on the 10 training datasets. The feature selection performance is evaluated by the average ACC on the training and validation datasets, which are shown in Fig. 9 to Fig. 12.

Theoretically, the SOCC and OBCC are identical when selecting the first feature. Therefore, the two methods always select the same first feature and give the same ACC results in Fig. 9 to Fig. 12. For the Flag dataset, as the OBCC adopts the point-biserial correlation coefficient, which can only be used in binomial classification issues, we extend it to the more general Pearson correlation coefficient for this multinomial classification case. The performance of the OBCC method in the training datasets is significantly worse than other two methods, which is shown in Fig. 11a. For the Parkinsons dataset, the first 6 features selected by the SHAP method is not helpful to the LDA classifiers. The proposed OLS method achieves the best ACC results in the Breast dataset and comparable results in the rest datasets. After the comprehensive comparison, the proposed OLS method does not show the obvious deficiency in any datasets, which implies the great robustness of this method across the different datasets.

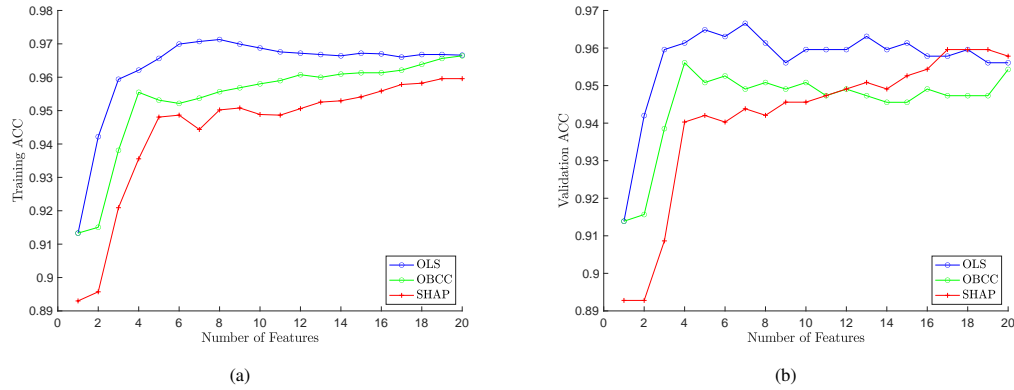


Fig. 9. The average ACC results of the 10-fold cross validation for the OLS, OBCC and SHAP methods on (a) training and (b) validation Breast dataset.

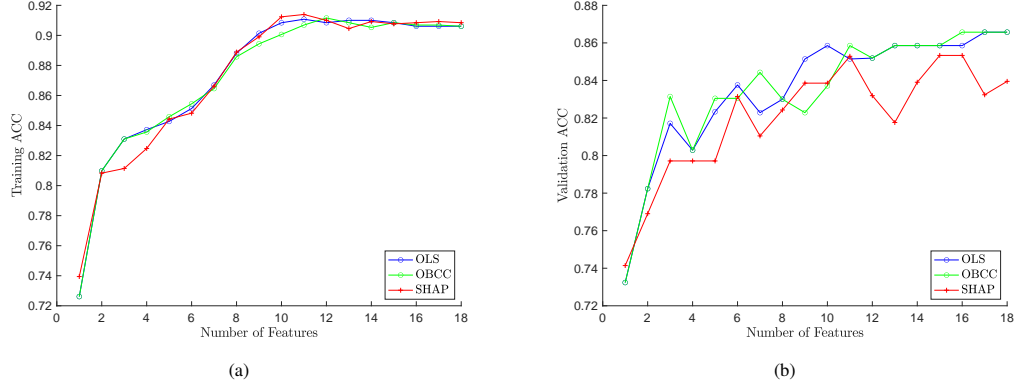


Fig. 10. The average ACC results of the 10-fold cross validation for the OLS, OBCC and SHAP methods on (a) training and (b) validation Lymph dataset.

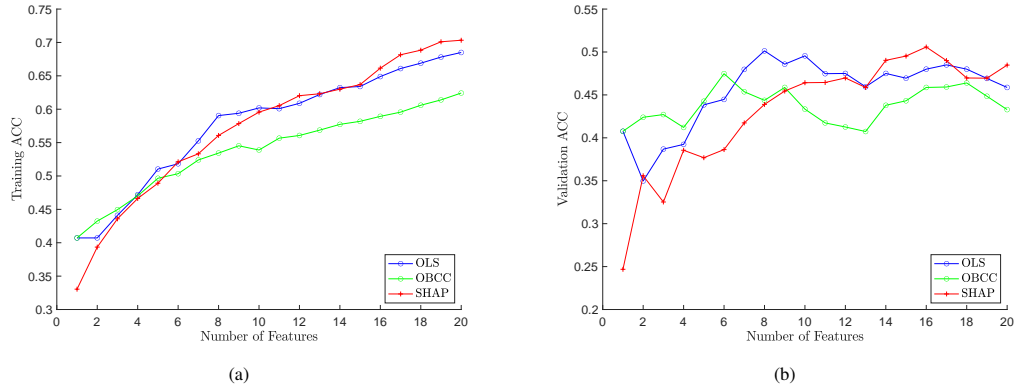


Fig. 11. The average ACC results of the 10-fold cross validation for the OLS, OBCC and SHAP methods on (a) training and (b) validation Flag dataset.

7. Conclusions

This paper proposes a novel OLS based feature selection method for classification. The method is based on the newly proposed concept of SOCCs which, for the first time, reveals an important relationship of the OLS based solution to a least-squares problem with the multiple correlation coefficient and the canonical correlation coefficient. Utilising the relationships, the OLS based feature selection method is developed where either the multiple correlation coefficient (for binomial classification) or the canonical correlation coefficient (for multinomial classification) is used as the feature ranking criterion. The relationship between CCA and LDA is analysed to demonstrate the statistical implication of the canonical correlation coefficient

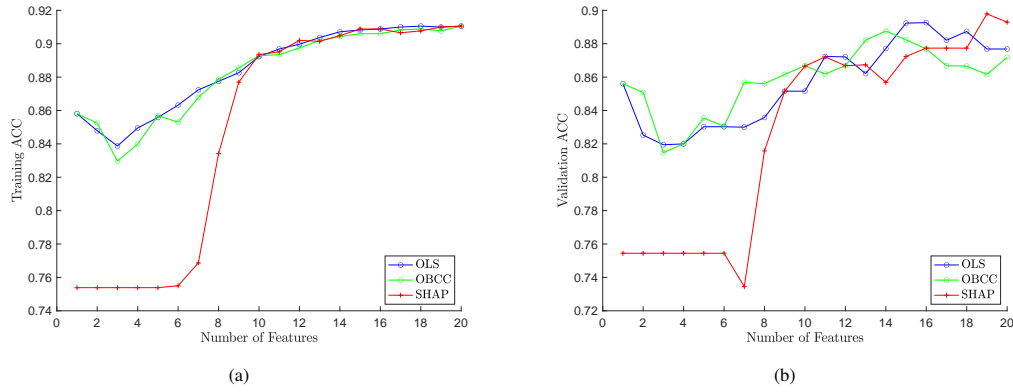


Fig. 12. The average ACC results of the 10-fold cross validation for the OLS, OBCC and SHAP methods on (a) training and (b) validation Parkinsons dataset.

hence the proposed SOCCs in LDA based classification problem. The speed advantage of the OLS based feature selection method in greedy search has been analysed. In empirical studies, a simple example has been used to illustrate the procedure of the OLS based feature selection method, and to demonstrate the relationship of the SOCCs with canonical correlation coefficient and Fisher's criterion. The synthetic and real world datasets have been used to compare the mutual information based methods and the embedded methods with the new OLS based method, showing that the OLS based method is always in the top 5 among 12 candidate methods. In addition, as the mutual information estimation normally requires the discretisation on the continuous features, the OLS based methods, which can deal with both numerical (including continuous and discrete) and categorical features, is more convenient to use than the mutual information based methods in addressing continuous feature based classification problems.

Acknowledgements

The authors would like to acknowledge that this work was supported by the UK Engineering and Physical Science Research Council Grant EP/R018480/1.

References

- [1] I. Guyon, A. Elisseeff, An introduction to variable and feature selection, *Journal of Machine Learning Research* 3 (2003) 1157–1182.

- [2] L. M. Abualigah, A. T. Khader, M. A. Al-Betar, O. A. Alomari, Text feature selection with a robust weight scheme and dynamic dimension reduction to text document clustering, *Expert Systems with Applications* 84 (2017) 24–36.
- [3] L. M. Q. Abualigah, A. T. Khader, E. S. Hanandeh, A new feature selection method to improve the document clustering using particle swarm optimization algorithm, *Journal of Computational Science* 25 (2018) 456–466.
- [4] H. Peng, F. Long, C. Ding, Feature selection based on mutual information: criteria of max-dependency, max-relevance, and min-redundancy, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 27 (8) (2005) 1226–1238.
- [5] L. M. Q. Abualigah, Feature selection and enhanced krill herd algorithm for text document clustering, Springer, 2019.
- [6] L. M. Q. Abualigah, Multi-verse optimizer algorithm: a comprehensive survey of its results, variants, and applications, *Neural Computing and Applications* (2020) 1–21.
- [7] R. Tibshirani, Regression shrinkage and selection via the lasso, *Journal of the Royal Statistical Society: Series B (Methodological)* 58 (1) (1996) 267–288.
- [8] L. Breiman, J. Friedman, C. Stone, R. Olshen, *Classification and Regression Trees*, Taylor & Francis, Boca Raton, FL, USA, 1984.
- [9] C. Ding, H. Peng, Minimum redundancy feature selection from microarray gene expression data, *Journal of Bioinformatics and Computational Biology* 3 (2) (2005) 185–205.
- [10] S. Chen, S. A. Billings, W. Luo, Orthogonal least squares methods and their application to non-linear system identification, *International Journal of Control* 50 (5) (1989) 1873–1896.
- [11] W. W. Cooley, P. R. Lohnes, *Multivariate Data Analysis*, Wiley, New York, NY, USA, 1971.
- [12] G. Brown, A. Pocock, M.-J. Zhao, M. Luján, Conditional likelihood maximisation: a unifying framework for information theoretic feature selection, *Journal of machine learning research* 13 (2012) 27–66.
- [13] Y. Li, T. Li, H. Liu, Recent advances in feature selection and its applications, *Knowledge and Information Systems* 53 (3) (2017) 551–577.

- [14] Z. Zeng, H. Zhang, R. Zhang, C. Yin, A novel feature selection method considering feature interaction, *Pattern Recognition* 48 (8) (2015) 2656–2666.
- [15] 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) 1–45.
- [16] S. Sharmin, M. Shoyaib, A. A. Ali, M. A. H. Khan, O. Chae, Simultaneous feature selection and discretization based on mutual information, *Pattern Recognition* 91 (2019) 162–174.
- [17] S. Kashef, H. Nezamabadi-pour, A label-specific multi-label feature selection algorithm based on the pareto dominance concept, *Pattern Recognition* 88 (2019) 654–667.
- [18] P. Zhang, G. Liu, W. Gao, Distinguishing two types of labels for multi-label feature selection, *Pattern Recognition* 95 (2019) 72–82.
- [19] L. Hu, Y. Li, W. Gao, P. Zhang, J. Hu, Multi-label feature selection with shared common mode, *Pattern Recognition* (2020) 107344.
- [20] H. A. Dau, A. Bagnall, K. Kamgar, C.-C. M. Yeh, Y. Zhu, S. Gharghabi, C. A. Ratanamahatana, E. Keogh, The ucr time series archive, *IEEE/CAA Journal of Automatica Sinica* 6 (6) (2019) 1293–1305.
- [21] J. Ircio, A. Lojo, U. Mori, J. A. Lozano, Mutual information based feature subset selection in multi-variate time series classification, *Pattern Recognition* 108 (2020) 107525.
- [22] M. J. Korenberg, Identifying nonlinear difference equation and functional expansion representations: the fast orthogonal algorithm, *Annals of biomedical engineering* 16 (1) (1988) 123–142.
- [23] M. J. Korenberg, A robust orthogonal algorithm for system identification and time-series analysis, *Biological cybernetics* 60 (4) (1989) 267–276.
- [24] S. Chen, C. F. Cowan, P. M. Grant, Orthogonal least squares learning algorithm for radial basis function networks, *IEEE Transactions on neural networks* 2 (2) (1991) 302–309.
- [25] H.-L. Wei, S. A. Billings, Feature subset selection and ranking for data dimensionality reduction, *IEEE transactions on pattern analysis and machine intelligence* 29 (1) (2006) 162–166.

- [26] J. R. A. Soares, H.-L. Wei, S. A. Billings, A novel logistic-narx model as a classifier for dynamic binary classification, *Neural Computing and Applications* 31 (1) (2019) 11–25.
- [27] J. Cohen, P. Cohen, *Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences*, Wiley, Hillsdale, MI, USA, 1975.
- [28] T. H. Cormen, C. E. Leiserson, R. L. Rivest, C. Stein, *Introduction to Algorithms*, 3rd Edition, MIT Press, Cambridge, MA, USA, 2009.
- [29] T. Sun, S. Chen, Class label versus sample label-based cca, *Applied Mathematics and Computation* 185 (1) (2007) 272–283.
- [30] R. A. Fisher, The use of multiple measurements in taxonomic problems, *Annals of Eugenics* 7 (2) (1936) 179–188.
- [31] D. W. Hosmer Jr, S. Lemeshow, R. X. Sturdivant, *Applied Logistic Regression*, 3rd Edition, Wiley, Hoboken, NJ, USA, 2013.
- [32] R. Salakhutdinov, I. Murray, On the quantitative analysis of deep belief networks, in: *Proceedings of the 25th International Conference on Machine Learning*, 2008, pp. 872–879.
- [33] D. Dua, C. Graff, *UCI machine learning repository* (2019).
URL <http://archive.ics.uci.edu/ml>
- [34] M. van Breukelen, R. P. Duin, D. M. Tax, J. Den Hartog, Handwritten digit recognition by combined classifiers, *Kybernetika* 34 (4) (1998) 381–386.
- [35] W. N. Street, W. H. Wolberg, O. L. Mangasarian, Nuclear feature extraction for breast tumor diagnosis, in: *Biomedical Image Processing and Biomedical Visualization*, Vol. 1905, International Society for Optics and Photonics, 1993, pp. 861–870.
- [36] P. M. Ciarelli, E. Oliveira, Agglomeration and elimination of terms for dimensionality reduction, in: *2009 Ninth International Conference on Intelligent Systems Design and Applications*, IEEE, 2009, pp. 547–552.
- [37] R. S. Michalski, I. Mozetic, J. Hong, N. Lavrac, The multi-purpose incremental learning system AQ15 and its testing application to three medical domains, in: *AAAI*, 1986, pp. 1041–1045.

- [38] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, *Advances in Neural Information Processing Systems* 30 (2017) 4765–4774.
- [39] D. Fryer, I. Strümke, H. Nguyen, Shapley values for feature selection: The good, the bad, and the axioms (2021). [arXiv:2102.10936](#).
- [40] L. S. Shapley, A value for n-person games, *Contributions to the Theory of Games* 2 (28) (1953) 307–317.
- [41] G. H. John, R. Kohavi, K. Pfleger, Irrelevant features and the subset selection problem, in: *Machine Learning Proceedings 1994*, Elsevier, 1994, pp. 121–129.
- [42] M. A. Little, P. E. McSharry, E. J. Hunter, J. Spielman, L. O. Ramig, Suitability of dysphonia measurements for telemonitoring of Parkinson’s disease, *IEEE Transactions on Biomedical Engineering* 56 (4) (2009) 1015–1022.

Appendix A. The proof of the relationship between SOCCs and multiple correlation coefficient

In the ordinary least-squares problem, the linear regression model with N instances is given by

$$\mathbf{y} = (\mathbf{1}, \mathbf{X}) \begin{pmatrix} \beta_0 \\ \boldsymbol{\beta} \end{pmatrix} + \mathbf{e}. \quad (\text{A.1})$$

It can be shown that the least-squares estimation of β_0 and $\boldsymbol{\beta}$, denoted as $\hat{\beta}_0$ and $\hat{\boldsymbol{\beta}}$, satisfies the equation (see Appendix C for proof)

$$\begin{aligned} \hat{\beta}_0 &= \bar{y} - \bar{\mathbf{x}}^\top \hat{\boldsymbol{\beta}} \\ &= \bar{y} - (\bar{x}_1, \dots, \bar{x}_n) \hat{\boldsymbol{\beta}}, \end{aligned} \quad (\text{A.2})$$

where \bar{y} is the sample mean of \mathbf{y} , and \bar{x}_i is the sample mean of \mathbf{x}_i . Substituting (A.2) into (A.1), the linear model (A.1) is simplified to

$$\mathbf{y}_C = \mathbf{X}_C \hat{\boldsymbol{\beta}} + \hat{\mathbf{e}}, \quad (\text{A.3})$$

where $\hat{\mathbf{e}}$ is an estimation for the error term \mathbf{e} , \mathbf{y}_C is the centred response variable given by

$$\mathbf{y}_C = \begin{pmatrix} y_1 - \bar{y} \\ \vdots \\ y_N - \bar{y} \end{pmatrix}, \quad (\text{A.4})$$

and \mathbf{X}_C is the matrix of the centred independent variables given by

$$\begin{aligned}\mathbf{X}_C &= (\mathbf{x}_{C1}, \dots, \mathbf{x}_{Cn}) \\ &= \begin{pmatrix} x_{1,1} - \bar{x}_1 & \dots & x_{1,n} - \bar{x}_n \\ \vdots & \ddots & \vdots \\ x_{N,1} - \bar{x}_1 & \dots & x_{N,n} - \bar{x}_n \end{pmatrix}.\end{aligned}\tag{A.5}$$

Equation (A.3) implies that $\hat{\boldsymbol{\beta}}$ satisfies the normal equation (see Appendix C for proof)

$$(\mathbf{X}_C^\top \mathbf{X}_C) \hat{\boldsymbol{\beta}} = \mathbf{X}_C^\top \mathbf{y}_C,\tag{A.6}$$

transforming the least-squares problem with the intercept into the least-squares problem without the intercept.

When \mathbf{X}_C has full column rank, the unnormalised reduced QR decomposition is performed on \mathbf{X}_C as

$$\mathbf{X}_C = \mathbf{W}_C \mathbf{A},\tag{A.7}$$

where \mathbf{A} is a $n \times n$ invertible upper triangular matrix and \mathbf{W}_C is a $N \times n$ matrix with the orthogonal columns $\mathbf{w}_{C1}, \dots, \mathbf{w}_{Cn}$. As $\mathbf{W}_C = \mathbf{X}_C \mathbf{A}^{-1}$, it can be seen that \mathbf{w}_{Ci} , which is the linear transformation of $\mathbf{x}_{C1}, \dots, \mathbf{x}_{Cn}$, has zero sample mean. Substituting (A.7) into (A.3) yields

$$\mathbf{y}_C = \mathbf{W}_C \hat{\mathbf{g}} + \hat{\mathbf{e}},\tag{A.8}$$

where $\hat{\mathbf{g}} = \mathbf{A} \hat{\boldsymbol{\beta}} = (\hat{g}_1, \dots, \hat{g}_n)^\top$. The parameter vector $\hat{\mathbf{g}}$ obviously satisfies the normal equation

$$\mathbf{W}_C^\top \mathbf{W}_C \hat{\mathbf{g}} = \mathbf{W}_C^\top \mathbf{y}_C.\tag{A.9}$$

Thus, the ordinary least-squares problem (A.6) about \mathbf{X}_C and \mathbf{y}_C is transformed into the OLS problem (A.9) about \mathbf{W}_C and \mathbf{y}_C .

The residual sum of squares for OLS is given by

$$\begin{aligned}\hat{\mathbf{e}}^\top \hat{\mathbf{e}} &= (\mathbf{y}_C - \mathbf{W}_C \hat{\mathbf{g}})^\top (\mathbf{y}_C - \mathbf{W}_C \hat{\mathbf{g}}) \\ &= \mathbf{y}_C^\top \mathbf{y}_C - 2 \hat{\mathbf{g}}^\top \mathbf{W}_C^\top \mathbf{y}_C + \hat{\mathbf{g}}^\top \mathbf{W}_C^\top \mathbf{W}_C \hat{\mathbf{g}}.\end{aligned}\tag{A.10}$$

Because of (A.9), this equation becomes

$$\hat{\mathbf{e}}^\top \hat{\mathbf{e}} = \mathbf{y}_C^\top \mathbf{y}_C - \hat{\mathbf{g}}^\top \mathbf{W}_C^\top \mathbf{W}_C \hat{\mathbf{g}}.\tag{A.11}$$

As \mathbf{W}_C is orthogonal, the inner product $\mathbf{W}_C^\top \mathbf{W}_C$ is the diagonal matrix $\text{diag}(\mathbf{w}_{C1}^\top \mathbf{w}_{C1}, \dots, \mathbf{w}_{Cn}^\top \mathbf{w}_{Cn})$. Thus, (A.11) can be rewritten to

$$\hat{\mathbf{e}}^\top \hat{\mathbf{e}} = \mathbf{y}_C^\top \mathbf{y}_C - \sum_{i=0}^n \hat{g}_i^2 \mathbf{w}_{Ci}^\top \mathbf{w}_{Ci}. \quad (\text{A.12})$$

Both sides of (A.12) are divided by $\mathbf{y}_C^\top \mathbf{y}_C$, that is

$$\frac{\hat{\mathbf{e}}^\top \hat{\mathbf{e}}}{\mathbf{y}_C^\top \mathbf{y}_C} = 1 - \sum_{i=1}^n \frac{\hat{g}_i^2 \mathbf{w}_{Ci}^\top \mathbf{w}_{Ci}}{\mathbf{y}_C^\top \mathbf{y}_C}. \quad (\text{A.13})$$

Due to the orthogonality of \mathbf{W}_C , the computation of the parameter vector $\hat{\mathbf{g}}$ can be simplified as

$$\hat{g}_i = \frac{\mathbf{w}_{Ci}^\top \mathbf{y}_C}{\mathbf{w}_{Ci}^\top \mathbf{w}_{Ci}}. \quad (\text{A.14})$$

Substituting (A.14) into (A.13),

$$\begin{aligned} \frac{\hat{\mathbf{e}}^\top \hat{\mathbf{e}}}{\mathbf{y}_C^\top \mathbf{y}_C} &= 1 - \sum_{i=1}^n \frac{\mathbf{y}_C^\top \mathbf{w}_{Ci} \mathbf{w}_{Ci}^\top \mathbf{y}_C}{\mathbf{w}_{Ci}^\top \mathbf{w}_{Ci} \mathbf{y}_C^\top \mathbf{y}_C} \\ &= 1 - \sum_{i=1}^n h_i, \end{aligned} \quad (\text{A.15})$$

where h_i is the SOCC.

The definition of the multiple correlation coefficient is based on the centred linear regression model (A.3). The squared multiple correlation coefficient (or called coefficient of determination) has the following relationship with the total sum of squares SST and the residual sum of squares SSR of the model (A.3)

$$R^2(\mathbf{X}, \mathbf{y}) = 1 - \frac{SSR(\mathbf{X}_C, \mathbf{y}_C)}{SST(\mathbf{X}_C, \mathbf{y}_C)}, \quad (\text{A.16})$$

where

$$\begin{aligned} SST(\mathbf{X}_C, \mathbf{y}_C) &= \mathbf{y}_C^\top \mathbf{y}_C \\ SSR(\mathbf{X}_C, \mathbf{y}_C) &= (\mathbf{y}_C - \hat{\mathbf{y}}_C)^\top (\mathbf{y}_C - \hat{\mathbf{y}}_C) = \hat{\mathbf{e}}^\top \hat{\mathbf{e}}. \end{aligned} \quad (\text{A.17})$$

Comparing (A.15) and (A.16), it is found

$$R^2(\mathbf{X}, \mathbf{y}) = \sum_{i=1}^n h_i. \quad (\text{A.18})$$

Appendix B. The proof of the relationship between SOCCs and canonical correlation coefficients

According to (26a), the sum of the squared canonical correlation coefficients, i.e. the sum of the eigenvalues, are given by

$$\sum_{k=1}^{n \wedge m} R_k^2(\mathbf{X}, \mathbf{Y}) = \text{tr}(\mathbf{R}_{\mathbf{X}, \mathbf{X}}^{-1} \mathbf{R}_{\mathbf{X}, \mathbf{Y}} \mathbf{R}_{\mathbf{Y}, \mathbf{Y}}^{-1} \mathbf{R}_{\mathbf{Y}, \mathbf{X}}), \quad (\text{B.1})$$

where the operator tr denotes the trace of the matrix. If the columns of \mathbf{X} are zero sample mean and orthogonal, the correlation matrix of \mathbf{X} is identity matrix, so $\mathbf{R}_{\mathbf{X},\mathbf{X}}^{-1} = \mathbf{I}$. It is known that the multiple correlation between \mathbf{Y} and each \mathbf{x}_i can be evaluated by [11, p. 174]

$$\begin{pmatrix} R^2(\mathbf{x}_1, \mathbf{Y}) \\ \vdots \\ R^2(\mathbf{x}_n, \mathbf{Y}) \end{pmatrix} = \text{diag}(\mathbf{R}_{\mathbf{X},\mathbf{Y}}\mathbf{R}_{\mathbf{Y},\mathbf{Y}}^{-1}\mathbf{R}_{\mathbf{Y},\mathbf{X}}), \quad (\text{B.2})$$

where the operator diag obtains the main diagonal of the matrix. Therefore, according to (B.1) and (B.2), the following equation holds when the columns of \mathbf{X} are centred and orthogonal.

$$\sum_{k=1}^{n \wedge m} R_k^2(\mathbf{X}, \mathbf{Y}) = \sum_{i=1}^n R^2(\mathbf{x}_i, \mathbf{Y}). \quad (\text{B.3})$$

Through the unnormalised reduced QR decomposition,

$$\begin{aligned} \mathbf{X}_C &= \mathbf{W}_C \mathbf{A} \\ \mathbf{Y}_C &= \mathbf{V}_C \mathbf{B} \end{aligned} \quad (\text{B.4})$$

where \mathbf{W}_C is a $N \times n$ matrix with the centred orthogonal columns given by

$$\mathbf{W}_C = (\mathbf{w}_{C1}, \dots, \mathbf{w}_{Cn}), \quad (\text{B.5})$$

\mathbf{V}_C is a $N \times m$ matrix with the centred orthogonal columns given by

$$\mathbf{V}_C = (\mathbf{v}_{C1}, \dots, \mathbf{v}_{Cm}), \quad (\text{B.6})$$

\mathbf{A} is a $n \times n$ invertible upper triangular matrix, and \mathbf{B} is a $m \times m$ invertible upper triangular matrix. It is noticed that the transformation from \mathbf{X} (or \mathbf{Y}) to \mathbf{W}_C (or \mathbf{V}_C) is affine. As the canonical correlation coefficient is invariant under affine transformations,

$$R_k(\mathbf{X}, \mathbf{Y}) = R_k(\mathbf{W}_C, \mathbf{V}_C) \quad k = 1, \dots, n \wedge m. \quad (\text{B.7})$$

As the columns of \mathbf{W}_C are centred and orthogonal, the following equation holds according to (B.3) and (B.7).

$$\sum_{k=1}^{n \wedge m} R_k^2(\mathbf{X}, \mathbf{Y}) = \sum_{k=1}^{n \wedge m} R_k^2(\mathbf{W}_C, \mathbf{V}_C) = \sum_{i=1}^n R^2(\mathbf{w}_{Ci}, \mathbf{V}_C). \quad (\text{B.8})$$

Define squared orthogonal correlation matrix as

$$\mathbf{H} = \begin{pmatrix} h_{1,1} & \dots & h_{1,m} \\ \vdots & \ddots & \vdots \\ h_{n,1} & \dots & h_{n,m} \end{pmatrix}, \quad (\text{B.9})$$

where $h_{i,j}$ is the SOCC given by

$$h_{i,j} = \frac{\mathbf{v}_{Cj}^\top \mathbf{w}_{Ci} \mathbf{w}_{Ci}^\top \mathbf{v}_{Cj}}{\mathbf{w}_{Ci}^\top \mathbf{w}_{Ci} \mathbf{v}_{Cj}^\top \mathbf{v}_{Cj}}. \quad (\text{B.10})$$

Due to (20), the multiple correlation coefficient between \mathbf{V}_C and each \mathbf{w}_C can be evaluated by

$$\begin{aligned} R^2(\mathbf{w}_{C1}, \mathbf{V}_C) &= \sum_{j=1}^m h_{1,j} \\ &\vdots \\ R^2(\mathbf{w}_{Cn}, \mathbf{V}_C) &= \sum_{j=1}^m h_{n,j}. \end{aligned} \quad (\text{B.11})$$

Substituting (B.11) into (B.8), it is found that the sum of the squared canonical correlation coefficients between \mathbf{X} and \mathbf{Y} is equal to the sum of all entries of the squared orthogonal correlation matrix \mathbf{H} , that is

$$\sum_{k=1}^{n \wedge m} R_k^2(\mathbf{X}, \mathbf{Y}) = \sum_{i=1}^n \sum_{j=1}^m h_{i,j}. \quad (\text{B.12})$$

Appendix C. Preliminary knowledge of linear regression

The ordinary least-squares estimation of a linear model with intercept is to find the optimal parameters $\hat{\beta}_0 \in \mathbb{R}$ and $\hat{\beta} \in \mathbb{R}^n$ to minimise the squared residual given by

$$L = (\mathbf{y} - \hat{\beta}_0 \mathbf{1} - \mathbf{X} \hat{\beta})^\top (\mathbf{y} - \hat{\beta}_0 \mathbf{1} - \mathbf{X} \hat{\beta}), \quad (\text{C.1})$$

where $\mathbf{y} \in \mathbb{R}^N$ is the response vector, $\mathbf{X} \in \mathbb{R}^{N \times n}$ is the design matrix, and $\mathbf{1}$ is the $N \times 1$ vector of ones. The squared residual is minimised when the derivative of L with respect to the parameters is zero, that is

$$\begin{aligned} \frac{\partial L}{\partial \hat{\beta}_1} &= \frac{\partial (\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}_1 \hat{\beta}_1 - \hat{\beta}_1^\top \mathbf{X}_1^\top \mathbf{y} + \hat{\beta}_1^\top \mathbf{X}_1^\top \mathbf{X}_1 \hat{\beta}_1)}{\partial \hat{\beta}_1} \\ &= -2\mathbf{X}_1^\top \mathbf{y} + 2\mathbf{X}_1^\top \mathbf{X}_1 \hat{\beta}_1 \\ &= 0, \end{aligned} \quad (\text{C.2})$$

where $\hat{\beta}_1 = (\hat{\beta}_0, \hat{\beta})^\top$ and $\mathbf{X}_1 = (\mathbf{1}, \mathbf{X})$. According to (C.2), it is known the optimal parameters satisfy the equation given by

$$\begin{aligned}\mathbf{X}_1^\top \mathbf{y} &= \mathbf{X}_1^\top \mathbf{X}_1 \hat{\beta}_1 \\ \begin{pmatrix} \mathbf{1}^\top \\ \mathbf{X}^\top \end{pmatrix} \mathbf{y} &= \begin{pmatrix} \mathbf{1}^\top \\ \mathbf{X}^\top \end{pmatrix} (\mathbf{1}, \mathbf{X}) \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta} \end{pmatrix} \\ \begin{pmatrix} N\bar{y} \\ \mathbf{X}^\top \mathbf{y} \end{pmatrix} &= \begin{pmatrix} N & N\bar{\mathbf{x}}^\top \\ N\bar{\mathbf{x}} & \mathbf{X}^\top \mathbf{X} \end{pmatrix} \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta} \end{pmatrix},\end{aligned}\tag{C.3}$$

where $\bar{y} \in \mathbb{R}$ is the sample mean of \mathbf{y} and $\bar{\mathbf{x}} \in \mathbb{R}^n$ is the column vector composed of the sample mean of each column of \mathbf{X} . Then, the two important equations can be found in (C.3). The first one is

$$\begin{aligned}N\bar{y} &= N\hat{\beta}_0 + N\bar{\mathbf{x}}^\top \hat{\beta} \\ \bar{y} &= \hat{\beta}_0 + \bar{\mathbf{x}}^\top \hat{\beta}.\end{aligned}\tag{C.4}$$

The second one is

$$\begin{aligned}\mathbf{X}^\top \mathbf{y} &= N\bar{\mathbf{x}}\hat{\beta}_0 + \mathbf{X}^\top \mathbf{X} \hat{\beta} \\ (\mathbf{X}_C + \mathbf{1}\bar{\mathbf{x}}^\top)^\top (\mathbf{y}_C + \bar{y}\mathbf{1}) &= N\bar{\mathbf{x}}\hat{\beta}_0 + (\mathbf{X}_C + \mathbf{1}\bar{\mathbf{x}}^\top)^\top (\mathbf{X}_C + \mathbf{1}\bar{\mathbf{x}}^\top) \hat{\beta} \\ \mathbf{X}_C^\top \mathbf{y}_C + \bar{y}\mathbf{X}_C^\top \mathbf{1} + \bar{\mathbf{x}}\mathbf{1}^\top \mathbf{y}_C + \bar{y}\bar{\mathbf{x}}\mathbf{1}^\top \mathbf{1} &= N\bar{\mathbf{x}}\hat{\beta}_0 + (\mathbf{X}_C^\top \mathbf{X}_C + \mathbf{X}_C^\top \mathbf{1}\bar{\mathbf{x}}^\top + \bar{\mathbf{x}}\mathbf{1}^\top \mathbf{X}_C + \bar{\mathbf{x}}\mathbf{1}^\top \mathbf{1}\bar{\mathbf{x}}^\top) \hat{\beta} \\ \mathbf{X}_C^\top \mathbf{y}_C + \bar{y}\mathbf{X}_C^\top \mathbf{1} + \bar{\mathbf{x}}\mathbf{1}^\top \mathbf{y}_C + N\bar{y}\bar{\mathbf{x}} &= N\bar{\mathbf{x}}\hat{\beta}_0 + (\mathbf{X}_C^\top \mathbf{X}_C + \mathbf{X}_C^\top \mathbf{1}\bar{\mathbf{x}}^\top + \bar{\mathbf{x}}\mathbf{1}^\top \mathbf{X}_C + N\bar{\mathbf{x}}\bar{\mathbf{x}}^\top) \hat{\beta},\end{aligned}\tag{C.5}$$

where \mathbf{X}_C is the column centred matrix of \mathbf{X} by its sample mean $\bar{\mathbf{x}}$ and \mathbf{y}_C is the centred vector of \mathbf{y} by its sample mean \bar{y} . As \mathbf{X}_C and \mathbf{y}_C have zero sample means, it is known that $\mathbf{X}_C^\top \mathbf{1}$ and $\mathbf{1}^\top \mathbf{y}_C$ are zeros. Due to this and (C.4), the equation (C.5) can be rewritten as

$$\begin{aligned}\mathbf{X}_C^\top \mathbf{y}_C + N\bar{y}\bar{\mathbf{x}} &= N\bar{\mathbf{x}}\hat{\beta}_0 + (\mathbf{X}_C^\top \mathbf{X}_C + N\bar{\mathbf{x}}\bar{\mathbf{x}}^\top) \hat{\beta} \\ \mathbf{X}_C^\top \mathbf{y}_C + N\bar{y}\bar{\mathbf{x}} &= N\bar{\mathbf{x}}(\bar{y} - \bar{\mathbf{x}}^\top \hat{\beta}) + (\mathbf{X}_C^\top \mathbf{X}_C + N\bar{\mathbf{x}}\bar{\mathbf{x}}^\top) \hat{\beta} \\ \mathbf{X}_C^\top \mathbf{y}_C &= \mathbf{X}_C^\top \mathbf{X}_C \hat{\beta}\end{aligned}\tag{C.6}$$