

## Problem Statement

Let  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]^T = [\boldsymbol{\chi}_1, \dots, \boldsymbol{\chi}_n] \in \mathbb{R}^{m \times n}$  be a design matrix, where  $\mathbf{x}_i \in \mathbb{R}^n$  is the  $i$ -th object,  $\boldsymbol{\chi}_j \in \mathbb{R}^m$  is the  $j$ -th feature. Denote by  $\mathcal{J} = \{1, \dots, n\}$  the feature index set, and let  $\mathcal{A} \subseteq \mathcal{J}$  be a feature index subset. Let  $\mathbf{y} \in \{1, \dots, K\}^m$  be a target vector. Suppose a function  $\mathbf{f}$  approximates the probabilities of class label taking on each of the  $K$  possible values given an object  $\mathbf{x}$ , a feature index subset  $\mathcal{A}$ , and model parameters  $\mathbf{W} = [\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(K)}] \in \mathbb{R}^{n \times K}$

$$\mathbf{f}(\mathbf{x}, \mathcal{A}, \mathbf{W}) = \begin{bmatrix} p(y = 1 | \mathbf{x}, \mathcal{A}, \mathbf{W}) \\ \dots \\ p(y = K | \mathbf{x}, \mathcal{A}, \mathbf{W}) \end{bmatrix}.$$

Let  $\mathbf{a} \in \mathbb{B}^n = \{0, 1\}^n$  be an indicator vector such that  $a_j = 1$  if and only if  $j \in \mathcal{A}$ . The vector  $\mathbf{a}$  and the index set  $\mathcal{A}$  are related by

$$a_j^* = 1 \Leftrightarrow j \in \mathcal{A}^*, j \in \mathcal{J}. \quad (1)$$

Further in this paper we equate the binary vector  $\mathbf{a}$  and the feature index subset  $\mathcal{A}$ .

The data fitting problem is to find parameters  $\mathbf{W}^*$  such that

$$\mathbf{W}^* = \arg \min_{\mathbf{W} \in \mathbb{R}^{n \times K}} S(\mathbf{W}, \mathbf{a} | \mathbf{X}, \mathbf{y}, \mathbf{f}), \quad (2)$$

where  $S$  is the error function, which validates the quality of the parameters  $\mathbf{W}$  and the corresponding feature index subset  $\mathcal{A}$  given a design matrix  $\mathbf{X}$ , a target vector  $\mathbf{y}$ , and a hypothesis function  $\mathbf{h}$ .

The features are assumed to be noisy, irrelevant or multicollinear, which leads to additional error in estimating the optimum model parameters  $\mathbf{W}^*$  and increases the instability of them. Feature selection methods can be used to remove certain features from the design matrix  $\mathbf{X}$ . The feature selection procedure reduces the dimensionality of problem (2) and improves the stability of the optimum parameters  $\mathbf{W}^*$ . The feature selection problem is

$$\mathbf{a}^* = \arg \min_{\mathbf{a} \in \mathbb{B}^n} Q(\mathbf{a} | \mathbf{X}, \mathbf{y}), \quad (3)$$

where  $Q : \mathbb{B}^n \rightarrow \mathbb{R}$  is a quality criterion that determines the quality of a selected feature index subset  $\mathcal{A} \subseteq \mathcal{J}$ .

Problem (3) does not necessarily require estimation of the optimum parameters  $\mathbf{W}^*$  for data fitting procedure. It uses the relationships between the features  $\boldsymbol{\chi}_j, j \in \mathcal{J}$  and the target vector  $\mathbf{y}$ . The solution  $\mathbf{a}^*$  of problem (3) is used for determining the optimum parameters  $\mathbf{W}^*$  of problem (2).

One could combine the data fitting problem with the feature selection procedure

$$\mathbf{W}^*, \mathbf{a}^* = \arg \min_{\substack{\mathbf{W} \in \mathbb{R}^{n \times K} \\ \mathbf{a} \in \mathbb{B}^n}} S(\mathbf{W}, \mathbf{a} | \mathbf{X}, \mathbf{y}, \mathbf{f}). \quad (4)$$

The problem (4) is a mixed integer optimization problem which includes continuous variables  $\mathbf{W}$  and binary variables  $\mathbf{a}$ .

This study explores the softmax probability function

$$p(y = k|\mathbf{x}, \mathcal{A}, \mathbf{W}) = \frac{\exp\left(\mathbf{x}_{\mathcal{A}}^{\top} \mathbf{w}_{\mathcal{A}}^{(k)}\right)}{\sum_{j=1}^K \exp\left(\mathbf{x}_{\mathcal{A}}^{\top} \mathbf{w}_{\mathcal{A}}^{(j)}\right)},$$

where  $\mathbf{x}_{\mathcal{A}}, \mathbf{w}_{\mathcal{A}}^{(j)}$  is the reduced object and the parameter vector consisting of elements with indices in  $\mathcal{A}$ .

The likelihood function is the probability of the observed data given the parameters

$$L(\mathbf{W}) = p(\mathbf{y}|\mathbf{X}, \mathcal{A}, \mathbf{W}) = \prod_{i=1}^m \prod_{k=1}^K p(y = k|\mathbf{x}_i, \mathcal{A}, \mathbf{W})^{[y_i=k]}.$$

The error function is the negative logarithm of likelihood. Hence, minimizing the error function is equivalent to maximizing the likelihood

$$S(\mathbf{W}, \mathcal{A}|\mathbf{X}, \mathbf{y}, \mathbf{h}) = -\log L(\mathbf{W}) = -\sum_{i=1}^m \sum_{k=1}^K [y_i = k] \log p(y = k|\mathbf{x}_i, \mathcal{A}, \mathbf{W}). \quad (5)$$

This function is known as the cross-entropy error for multiclass classification problem. The error function (5) coincides with the logistic regression error function in the case of binary classification  $K = 2$ .

## 1 Quadratic Optimization Approach to the Multicollinearity Problem

In (Katrutsa stress), it shown that none of the feature selection methods considered (LARS, Lasso, Ridge, Stepwise and Genetic algorithm) solve problem (2) and give a model that is simultaneously stable, accurate and nonredundant. In contrast, in (Katrutsa QP) it was proposed quadratic programming approach to solving the multicollinearity problem for regression problem.

The main idea of the proposed approach is to minimize the number of similar features and maximize the number of relevant features.

To formalize this idea it was introduced the functions Sim and Rel:

$$\begin{aligned} \text{Sim: } \mathcal{J} \times \mathcal{J} &\rightarrow [0, 1], \\ \text{Rel: } \mathcal{J} &\rightarrow [0, 1]. \end{aligned} \quad (6)$$

These functions are problem-dependent, defined by the user before performing feature selection, and indicate how to measure feature similarity (Sim) and relevance to the target vector (Rel).

The criterion  $Q$  from problem (3) is represented as a quadratic function

$$Q(\mathbf{a}) = \mathbf{a}^{\top} \mathbf{Q} \mathbf{a} - \mathbf{b}^{\top} \mathbf{a}, \quad (7)$$

where  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  is a matrix of pairwise feature similarities, and  $\mathbf{b} \in \mathbb{R}^n$  is a vector of the relevances of features to the target vector.

To formalize this idea it was represented the criterion  $Q$  from problem (3) as a quadratic function

$$Q(\mathbf{a}) = \mathbf{a}^\top \mathbf{Q} \mathbf{a} - \mathbf{b}^\top \mathbf{a}, \quad (8)$$

where  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  is a matrix of pairwise feature similarities, and  $\mathbf{b} \in \mathbb{R}^n$  is a vector of the relevances of features to the target vector.

The matrix  $\mathbf{Q}$  is computed using Sim function :

$$\mathbf{Q} = [q_{ij}] = \text{Sim}(\boldsymbol{\chi}_i, \boldsymbol{\chi}_j), \quad (9)$$

and the vector  $\mathbf{b}$  is computed using Rel function:

$$\mathbf{b} = [b_i] = \text{Rel}(\boldsymbol{\chi}_i). \quad (10)$$

The optimum feature index set  $\mathcal{A}^*$  is defined by solution of the optimization problem

$$\mathbf{a}^* = \arg \min_{\mathbf{a} \in \mathbb{B}^n} \mathbf{a}^\top \mathbf{Q} \mathbf{a} - \mathbf{b}^\top \mathbf{a}, \quad (11)$$

One of the example of defining the Sim and Rel function proposed in (Katrutsa QP) is the following

$$q_{ij} = \text{Sim}(\boldsymbol{\chi}_i, \boldsymbol{\chi}_j) = |\hat{\rho}_{ij}|, \quad (12)$$

$$b_i = \text{Rel}(\boldsymbol{\chi}_i) = |\hat{\rho}_{iy}|, \quad (13)$$

where  $\hat{\rho}_{ij}$  are sample correlation coefficients between features  $\boldsymbol{\chi}_i$  and  $\boldsymbol{\chi}_j$ ,  $\hat{\rho}_{iy}$  are the sample correlation coefficient between feature  $\boldsymbol{\chi}_i$  and target vector  $\mathbf{y}$ . The sample correlation coefficient is defined as

$$\hat{\rho}_{ij} = \frac{(\boldsymbol{\chi}_i - \bar{\boldsymbol{\chi}}_i)^\top (\boldsymbol{\chi}_j - \bar{\boldsymbol{\chi}}_j)}{\|\boldsymbol{\chi}_i - \bar{\boldsymbol{\chi}}_i\|_2 \|\boldsymbol{\chi}_j - \bar{\boldsymbol{\chi}}_j\|_2}, \quad \bar{\boldsymbol{\chi}}_i = [\bar{\chi}_i, \dots, \bar{\chi}_i]^\top, \quad \bar{\boldsymbol{\chi}}_j = [\bar{\chi}_j, \dots, \bar{\chi}_j]^\top \quad (14)$$

where  $\bar{\chi}_i$  and  $\bar{\chi}_j$  are the means of  $\boldsymbol{\chi}_i$  and  $\boldsymbol{\chi}_j$  respectively.

We propose to extend this approach to the classification problems. The definition of the  $Q$  matrix is the same as in (12). In order to define relevancies to the target vector we build the model with only one feature  $\boldsymbol{\chi}_i$  and the target vector  $\mathbf{y}$ . Let denote by  $\hat{\mathbf{y}}_i$  the predictions of  $i$ th such model

$$\hat{\mathbf{y}}_i = \text{sign} \left( v_i^{(0)} + v_i^{(1)} \boldsymbol{\chi}_i \right), \quad i = 1, \dots, n.$$

We propose to use logistic regression to fit parameters  $\{(v_i^{(0)}, v_i^{(1)})\}_{i=1}^n$ . Then the elements of  $\mathbf{b} = [b_i]$  are defined as the absolute values of the sample correlation coefficient between the feature  $\boldsymbol{\chi}_i$  and the target vector  $\mathbf{y}$ :

$$b_i = \text{Rel}(\boldsymbol{\chi}_i) = |\hat{\rho}_{\hat{\mathbf{y}}_i \mathbf{y}}|. \quad (15)$$

## 1.1 Convex representation of the feature selection problem

The quadratic programming approach to the multicollinearity problem leads to problem (11), which is NP-hard because of the Boolean domain. Therefore, this problem is approximated with a convex optimization problem to solve it efficiently.

Assume that Sim gives a positive semidefinite matrix  $\mathbf{Q}$ . Then the quadratic form (8) is a convex function. To represent problem (11) in convex form, we have to replace the non-convex set  $\mathbb{B}^n$  with a convex set. A natural way to achieve this is to use the convex hull of  $\mathbb{B}^n$ :

$$\text{Conv}(\mathbb{B}^n) = [0, 1]^n.$$

Problem (11) is now approximated by the following *convex optimization problem*:

$$\begin{aligned} \mathbf{z}^* &= \arg \min_{\mathbf{z} \in [0, 1]^n} \mathbf{z}^\top \mathbf{Q} \mathbf{z} - \mathbf{b}^\top \mathbf{z} \\ \text{s.t. } &\|\mathbf{z}\|_1 \leq 1. \end{aligned} \tag{16}$$

This constraint is added to show that  $\mathbf{z}^*$  can be treated as a vector of non-normalized probabilities for every feature to be selected in the active set  $\mathcal{A}^*$ .

To return from a continuous vector  $\mathbf{z}^*$  to a Boolean vector  $\mathbf{a}^*$  and consequently to an active set  $\mathcal{A}^*$  (see equation (1)), we use the *significance threshold*  $\tau$ .

The value  $\tau$  is a significance threshold if  $z_j^* > \tau$  if and only if  $a_j^* = 1$  and  $j \in \mathcal{A}^*$ .

Tuning the value of  $\tau$  is problem-dependent and is based on the appropriate error rate, the number of features selected and the values of the evaluation criteria. To obtain the most appropriate significance threshold for a specific problem, we need to set a range of values for  $\tau$ .

## Mixed Integer Optimization Software

We use python-embedded [CVXPY](#) free software package for solving mixed integer optimization problem.

The considered types of problems which have open source solvers:

- LP — Linear Programming;
- SOCP — Second-Order Cone Programming;
- SDP — Semidefinite Programming;
- EXP — Problems with Exponential cone constraints;
- MIP — Mixed Integer Programming

Table 1 shows available solvers and types of the problems which they can handle. In this paper GUROBI solver is used.

Table 1: Solvers vs Problems

	LP	SOCP	SDP	EXP	MIP
CBC	X				X
GLPK	X				
GLPK-MI	X				X
Elemental	X	X			
ECOS	X	X		X	
ECOS-BB	X	X		X	X
GUROBI	X	X			X
MOSEK	X	X	X		
CVXOPT	X	X	X	X	
SCS	X	X	X	X	

## Computational experiment

This section provides experiments on the real datasets to show the performance of the proposed approach. In the computational experiment we compare the following methods: Quadratic Programming approach, Mixed Integer Programming, Lasso, Ridge, Elastic Net.