

# Efficient Leave-One-Out Strategy for Supervised Feature Selection

Dingcheng Feng, Feng Chen\*, and Wenli Xu

**Abstract:** Feature selection is a key task in statistical pattern recognition. Most feature selection algorithms have been proposed based on specific objective functions which are usually intuitively reasonable but can sometimes be far from the more basic objectives of the feature selection. This paper describes how to select features such that the basic objectives, e.g., classification or clustering accuracies, can be optimized in a more direct way. The analysis requires that the contribution of each feature to the evaluation metrics can be quantitatively described by some score function. Motivated by the conditional independence structure in probabilistic distributions, the analysis uses a leave-one-out feature selection algorithm which provides an approximate solution. The leave-one-out algorithm improves the conventional greedy backward elimination algorithm by preserving more interactions among features in the selection process, so that the various feature selection objectives can be optimized in a unified way. Experiments on six real-world datasets with different feature evaluation metrics have shown that this algorithm outperforms popular feature selection algorithms in most situations.

**Key words:** leave-one-out; feature selection objectives; evaluation metrics

## 1 Introduction

Feature selection is a fundamental task for many statistical pattern recognition applications such as image processing, speech recognition, text mining, and bioinformatics<sup>[1-6]</sup>. Feature selection is usually a discrete process using advanced combinatorial mathematics, especially discrete optimization. Due to its importance in applications and the optimization difficulties, feature selection has attracted much attention with the development of pattern recognition and machine learning methods<sup>[1,7,8]</sup>.

Many new feature selection methods have been

recently proposed using the formalism of nonlinear dimensionality reduction, especially the manifold learning techniques, such as the Laplacian eigenmap<sup>[9]</sup> and locally linear embedding<sup>[10]</sup>. Laplacian score, which is closely related to the canonical Fisher score<sup>[11]</sup>, evaluates the importance of each feature by its ability to preserve the locality relations among samples<sup>[12]</sup>. SPEC (SPECtrum decomposition of graph Laplacian) extends the Laplacian score by providing a unified perspective to systematically select features based on the properties of graph Laplacians with a series of heuristic strategies for deriving the SPEC framework and illustrations of the connection between Laplacian score and ReliefF<sup>[13,14]</sup>. The Trace Ratio (TR) criterion generalizes both the Laplacian score and SPEC by evaluating a single feature to evaluating a feature subset with a “trace-ratio” objective function defined with respect to the feature selection matrix<sup>[15]</sup>. New objectives for feature selection have recently been proposed. For example, Multi-Cluster Feature Selection (MCFS) aims to select features which best preserve the multi-cluster structure within the data<sup>[16]</sup>,

---

• Dingcheng Feng, Feng Chen, and Wenli Xu are with the National Laboratory for Information Science and Technology, Tsinghua University, Beijing 100084, China. E-mail: fdc08@mails.tsinghua.edu.cn; {chenfeng, xuwl}@mail.tsinghua.edu.cn.

\* To whom correspondence should be addressed.

Manuscript received: 2012-10-15; revised: 2013-06-05; accepted: 2013-06-07

while the Similarity Preserving Feature Selection (SPFS) framework attempts to select features which best preserve the sample similarity matrix<sup>[6]</sup>.

Previous feature selection algorithms have been based on specific objective functions, while the “ultimate” feature evaluation metrics (e.g., classification or clustering accuracies) usually differ from these specific objectives. This study optimizes the “ultimate” metrics (which is a more fundamental objective). Specifically, this method selects feature subsets such that the fundamental objectives are optimized. Further, the model requires that the feature selection algorithm returns a quantitative measure of the contribution of each feature to the fundamental objectives by some score function. A leave-one-out feature selection algorithm is used to provide an approximate solution for these purposes. The leave-one-out algorithm is motivated by the conditional independence structure in probabilistic distributions, and can be viewed as an effort to overcome the limitation of the conventional “greedy backward elimination” method. The algorithm is compared with many popular feature selection algorithms for real-world datasets with different evaluation metrics with the results showing that the leave-one-out algorithm outperforms previous methods in most situations.

## 2 Feature Selection Objectives

A common “framework” for most feature selection algorithms can be summarized as the following two closely related procedures: (1) Formulate the feature selection to optimize some objective function which represents some reasonable intuitive idea explicitly. (2) Utilize an optimization strategy which leads to meaningful solutions of the designed objectives. The objective functions should relate to the feature selection motivation. However, the feature selection formulation can differ greatly, which can lead to different optimization objectives. The following illustrates the differences by describing several of the most popular feature selection objectives.

- SPEC<sup>[14]</sup> (The Laplacian score<sup>[12]</sup> is used as special case). The objective is to select a feature subset to form the best lower dimensional representation by

$$\min_{\mathbf{S}} \mathcal{J}(\mathbf{S}) = \text{tr}(\mathbf{Y}\mathcal{L}\mathbf{Y}^T) \quad (1)$$

where  $\mathcal{L} = \mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}}$  is the normalized Laplacian matrix,  $\mathbf{L} = \mathbf{D} - \mathbf{W}$  is the Laplacian matrix,  $W_{ij} =$

$\frac{1}{n_i}$  if  $y_i = y_j = l$  and 0 elsewhere (supervised

case) or  $W_{ij} = \exp\left(-\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2}{2\sigma^2}\right)$  (unsupervised

case) is the pairwise instance similarity matrix,  $\mathbf{D}$  is the degree matrix of  $\mathbf{W}$  defined by a diagonal matrix with  $D_{ii} = \sum_{j=1}^m W_{ij}$  ( $i = 1, \dots, n$ ),  $\mathbf{Y} = [\mathbf{y}_j^{(i)}]$  ( $i = 1, \dots, m$ ,  $j = 1, \dots, k$ ) is the submatrix of the data matrix defined by  $\mathbf{y}^{(i)} = \mathbf{S}\mathbf{x}^{(i)}$  (representing data points with respect to selected variables), and  $\mathbf{S} \in \{0, 1\}^{k \times n}$  ( $\mathbf{S}\mathbf{1}_{n \times 1} = \mathbf{1}_{k \times 1}$ ,  $\|\mathbf{1}_{k \times 1}^T \mathbf{S}\|_0 = k$ ) is the feature selection matrix.

- TR criterion<sup>[15]</sup>. The objective is to select a feature subset to simultaneously minimize the within-class affinity and maximize the between-class affinity by

$$\max_{\mathbf{S}} \mathcal{J}(\mathbf{S}) = \frac{\text{tr}(\mathbf{S}\mathbf{X}\mathbf{L}_b\mathbf{X}^T\mathbf{S}^T)}{\text{tr}(\mathbf{S}\mathbf{X}\mathbf{L}_w\mathbf{X}^T\mathbf{S}^T)} \quad (2)$$

where  $\mathbf{X} \in \mathbb{R}^{n \times m}$  is the data matrix with  $m$  samples and  $n$  features,  $\mathbf{S} \in \{0, 1\}^{k \times n}$  ( $\mathbf{S}\mathbf{1}_{n \times 1} = \mathbf{1}_{k \times 1}$ ,  $\|\mathbf{1}_{k \times 1}^T \mathbf{S}\|_0 = k$ ) is a feature selection matrix,  $\mathbf{L}_w = \mathbf{D}_w - \mathbf{A}_w$ ,  $\mathbf{L}_b = \mathbf{D}_b - \mathbf{A}_b$ ,  $\mathbf{A}_w$  and  $\mathbf{A}_b$  represent the within-class (or local) and between-class (or global) affinity relationships among samples,  $\mathbf{D}_w$  is the degree matrix of  $\mathbf{A}_w$  defined by a diagonal matrix with  $(\mathbf{D}_w)_{ii} = \sum_{j=1}^m (\mathbf{A}_w)_{ij}$  ( $i = 1, \dots, n$ ), and  $\mathbf{D}_b$  is the degree matrix of  $\mathbf{A}_b$  defined by a diagonal matrix with  $(\mathbf{D}_b)_{ii} = \sum_{j=1}^m (\mathbf{A}_b)_{ij}$  ( $i = 1, \dots, n$ ).

- SPFS<sup>[6]</sup>. The objective is to select a feature subset which best preserves the similarity among samples by

$$\min_{\mathbf{S}} \mathcal{J}(\mathbf{S}) = \|\mathbf{X}^T \mathbf{S}^T \mathbf{S} \mathbf{X} - \mathbf{W}\| \quad (3)$$

where  $\mathbf{X} \in \mathbb{R}^{n \times m}$  is the data matrix with  $m$  samples and  $n$  features,  $\mathbf{S} \in \{0, 1\}^{k \times n}$  ( $\mathbf{S}\mathbf{1}_{n \times 1} = \mathbf{1}_{k \times 1}$ ,  $\|\mathbf{1}_{k \times 1}^T \mathbf{S}\|_0 = k$ ) is a feature selection matrix, and  $\mathbf{W}$  is a predefined similarity matrix which can be constructed either by using label information (supervised case) or using distance matrices (unsupervised case).

- MCFS<sup>[16]</sup>. The objective is to select a feature subset which best preserves the multi-cluster structure among samples by

$$\max_{\mathbf{S}} \mathcal{J}(\mathbf{S}) = \sum_{j \in F(\mathbf{S})} \text{MCFS}(j) \quad (4)$$

where  $F(\mathbf{S})$  is the selected feature subset,  $\text{MCFS}(j)$  is the score for feature  $j$  defined by  $\text{MCFS}(j) = \max_k |a_{k,j}|$ ,  $\mathbf{a}_k \in \mathbb{R}^{n \times 1}$  is the optimal solution to the  $L_1$ -regularized regression problem

$$\min_{\mathbf{a}_k} \|\mathbf{y}^{(k)} - \mathbf{X}^T \mathbf{a}_k\| + \|\mathbf{a}_k\|_1 \quad (5)$$

where  $\mathbf{X} \in \mathbb{R}^{n \times m}$  is the data matrix with  $m$  samples and  $n$  features,  $\mathbf{y}^{(k)} \in \mathbb{R}^{m \times 1}$  is the solution to the generalized eigenvalue decomposition problem  $\mathbf{L}\mathbf{y} = \lambda\mathbf{D}\mathbf{y}$  ( $\mathbf{y}^{(k)}$  corresponds to the  $k$ -th smallest eigenvalue),  $\mathbf{L} = \mathbf{D} - \mathbf{W}$  is the Laplacian matrix,  $\mathbf{W}$  is the weight matrix defined on a graph (can have multiple choices), and  $\mathbf{D}$  is the degree matrix of  $\mathbf{W}$  defined by a diagonal matrix with  $D_{ii} = \sum_{j=1}^m W_{ij}$  ( $i = 1, \dots, n$ ).

SPEC, TR, SPFS, and MCFS are all appropriate feature selection formulations, which select feature subsets based on reasonable intuition. However, the direct connections among these different motivations are difficult to identify, and users have difficulty choosing which algorithm to apply in specific applications. Most feature selection evaluation criteria are the classification or clustering accuracies, which are the “ultimate” objective for the feature selection. The next section describes a more direct perspective to optimizing the more fundamental objectives of the feature selection.

### 3 Leave-One-Out Feature Selection Framework

Three popular evaluation metrics for feature selection are used here to evaluate a dataset  $\mathbf{X} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}] \in \mathbb{R}^{n \times m}$  with the corresponding labels  $\mathbf{y} = [y^{(1)}, \dots, y^{(m)}] \in \mathbb{R}^{1 \times m}$  where  $y^{(i)} \in \{1, \dots, C\}$  is the label of sample point  $\mathbf{x}^{(i)}$ . Specifically, the selected feature subset  $F(\mathcal{S}) \in \{1, \dots, n\}$  is evaluated using the following evaluation metric functions on  $\mathbf{X}(F, :) = [\mathbf{x}_F^{(1)}, \dots, \mathbf{x}_F^{(m)}] \in \mathbb{R}^{k \times m}$  (denoted by  $\mathbf{X}_F$ ) and  $\mathbf{y}$ , which can be computed and  $\mathbf{x}_F^{(i)} = \mathbf{x}^{(i)}(F)$  denotes the elements indexed by  $F$  in  $\mathbf{x}^{(i)}$ .

- Classification accuracy<sup>[6, 12, 14-16]</sup>:

$$\max \mathcal{J}_{\text{classify}} = \frac{\sum_{i=1}^{m_{\text{test}}} \delta(y^{(i)}, \hat{f}(\mathbf{x}_F^{(i)}))}{m_{\text{test}}} \quad (6)$$

where  $\hat{f}$  is a classifier trained with  $m_{\text{train}}$  data points randomly sampled from  $(\mathbf{X}_F, \mathbf{Y})$ ,  $m_{\text{test}}$  is the number of remaining samples (test samples), and  $\delta(x, y)$  is the delta function which equals 1 if  $x = y$  and 0 otherwise.

- Clustering accuracy<sup>[12, 16]</sup>:

$$\max \mathcal{J}_{\text{cluster}} = \frac{\sum_{i=1}^m \delta(y^{(i)}, \text{map}(c^{(i)}))}{m} \quad (7)$$

where  $y^{(i)}$  is the provided label,  $c^{(i)}$  is the cluster label of data point  $\mathbf{x}^{(i)}$ , and  $\text{map}(c^{(i)})$  maps each cluster label  $c^{(i)}$  to the equivalent label in data.

- Normalized mutual information<sup>[12, 16]</sup>:

$$\max \mathcal{J}_{mi} = \frac{\text{MI}(C, C')}{\max(H(C), H(C'))} \quad (8)$$

where

$$\text{MI}(C, C') = \sum_{c_i \in C, c'_j \in C'} p(c_i, c'_j) \log_2 \frac{p(c_i, c'_j)}{p(c_i)p(c'_j)}$$

and  $p(c_i)$  is the probability of a sample proportion in cluster  $c_i$ , and  $H(C)$  is the entropy of cluster  $C$ .

The feature selection objectives (Eqs. (6)-(8)) are more fundamental than the case in Eqs. (1)-(4). Previous studies have proposed some “intuitive” feature selection criteria (e.g., Eqs. (1)-(4)), and evaluate the selected feature subset by the “ultimate” evaluation criteria (e.g., Eqs. (6)-(8)). Though these two classes of objectives are closely connected as illustrated by the effectiveness of previous feature selection algorithms, good solutions of the intuitive criteria cannot guarantee good solutions of the ultimate in generic situations.

This study aims to provide a unified framework to find good feature subsets with respect to arbitrary evaluation objectives. Specifically, the model is based on

$$\max_F \mathcal{J}(F) \quad (9)$$

where  $\mathcal{J}$  can be arbitrary feature selection objectives. We are interested in cases when  $\mathcal{J}$  is Eqs. (6)-(8) which are more fundamental and direct feature selection criteria. Moreover, the model will give individual representation of the importance of each feature, i.e., give each feature a “score” which describes its contribution to the objective. Though the “fundamental objective” and “individual representation” are desirable, Eq. (9) cannot be solved in the general case.

Thus, the structure of conditional independence relations<sup>[17]</sup> is used to provide reasonable approximations for Eq. (9). The Markov blanket of a random variable  $y$  with respect to  $\{y, \mathbf{x}\}$  is a variable subset  $\text{MBL}(y) \subset \mathbf{x} = \{x_1, \dots, x_n\}$  when  $y$  is independent of all the variables in  $\mathbf{x}$ . The smallest Markov blanket is called the Markov boundary of  $y$ , denoted by  $\text{MB}(y)$ , which is usually unique with one sufficient condition that  $p(y, \mathbf{x})$  is strictly positive<sup>[17]</sup>. One property of a Markov blanket is that the intersection of two Markov blankets is still a Markov blanket. Note that the feature subset selection can also be viewed as identifying the Markov boundary of class label  $y$  from the feature variables  $\mathbf{x}$  such that  $p(y|\mathbf{x}) = p(y|\text{MB}(y))$ . Given  $\text{MB}(y)$ , where  $y$  is

independent of all  $\mathbf{x} \setminus \text{MB}(y)$ , then  $\text{MB}(y)$  contains the information for the conditional distribution of  $y$  and, therefore, is a good feature subset. According to the property of Markov blankets,  $y$  is independent of all  $\mathbf{x} \setminus A$  provided that  $A$  is a variable subset satisfying  $\text{MB}(y) \subset A \subset \mathbf{x}$ . However, when  $\text{MB}(y) \not\subset A$ , i.e.,  $\exists x_i$  such that  $x_i \in \text{MB}(y)$  but  $x_i \notin A$ ,  $y$  is not independent of  $\mathbf{x} \setminus A$  given  $A$  and, therefore,  $A$  cannot be viewed as a good feature subset.

The conditional independence suggests the intuition that removing an important feature may have a bigger impact on the conditional distribution of  $p(y|\mathbf{x})$  than adding a redundant feature. Note that  $\mathbf{x}$  is also a (trivial) Markov blanket of  $y$ . Therefore, the importance of feature  $x_i$  can be reflected in  $D(p(y|\mathbf{x}), p(y|\mathbf{x}_{-i}))$  where  $\mathbf{x}_{-i} = \mathbf{x} \setminus \{x_i\}$  and  $D$  is some distance metric of two distributions. Since general supervised learning is based on detecting the properties of the full conditional distribution  $p(y|\mathbf{x})$ , the individual features can reasonably be evaluated by their contributions to  $p(y|\mathbf{x})$ . Specifically, the score for feature  $x_i$  can be given as  $w(x_i) = \mathcal{J}(\mathbf{x}) - \mathcal{J}(\mathbf{x}_{-i})$ , where  $\mathcal{J}$  is an arbitrary evaluation criterion of a feature subset such as in Eqs. (6)-(8).

Algorithm 1 summarizes this feature selection idea which can use the previous evaluation metrics  $\mathcal{J}$  in Eqs. (6)-(8). The algorithm returns the quantitative importance for each feature to the metric. If a larger (smaller) value of  $\mathcal{J}$  corresponds to a good feature subset, then a larger (smaller) weight,  $w(x_i)$ , will correspond to a better (worse) feature. Algorithm 1 runs the first step of the “greedy backward elimination”, but stops at this step and then performs similar computations for the other features. Therefore, Algorithm 1 avoids the drawback of the “greedy backward elimination” which usually converges to a local optimal feature subset, since the effect of the interactions among features on the evaluation metric is neglected in the selection process. Algorithm 1 overcomes this limitation by evaluating the  $n - 1$  features together, which preserves possible interactions among features.

## 4 Experiments

This section shows that the leave-one-out algorithm can outperform many popular feature selection algorithms with real-world datasets (Table 1). Figure 1 compares the different feature selection algorithms

### Algorithm 1 The Leave-one-out Feature Selection Algorithm

**Input:** Feature matrix  $\mathbf{X} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}] \in \mathbb{R}^{n \times m}$ , label vector  $\mathbf{y} = [y^{(1)}, \dots, y^{(m)}] \in \mathbb{R}^{1 \times m}$ , a feature evaluation metric  $\mathcal{J}$  such as in Eqs. (6), (7), and (8).

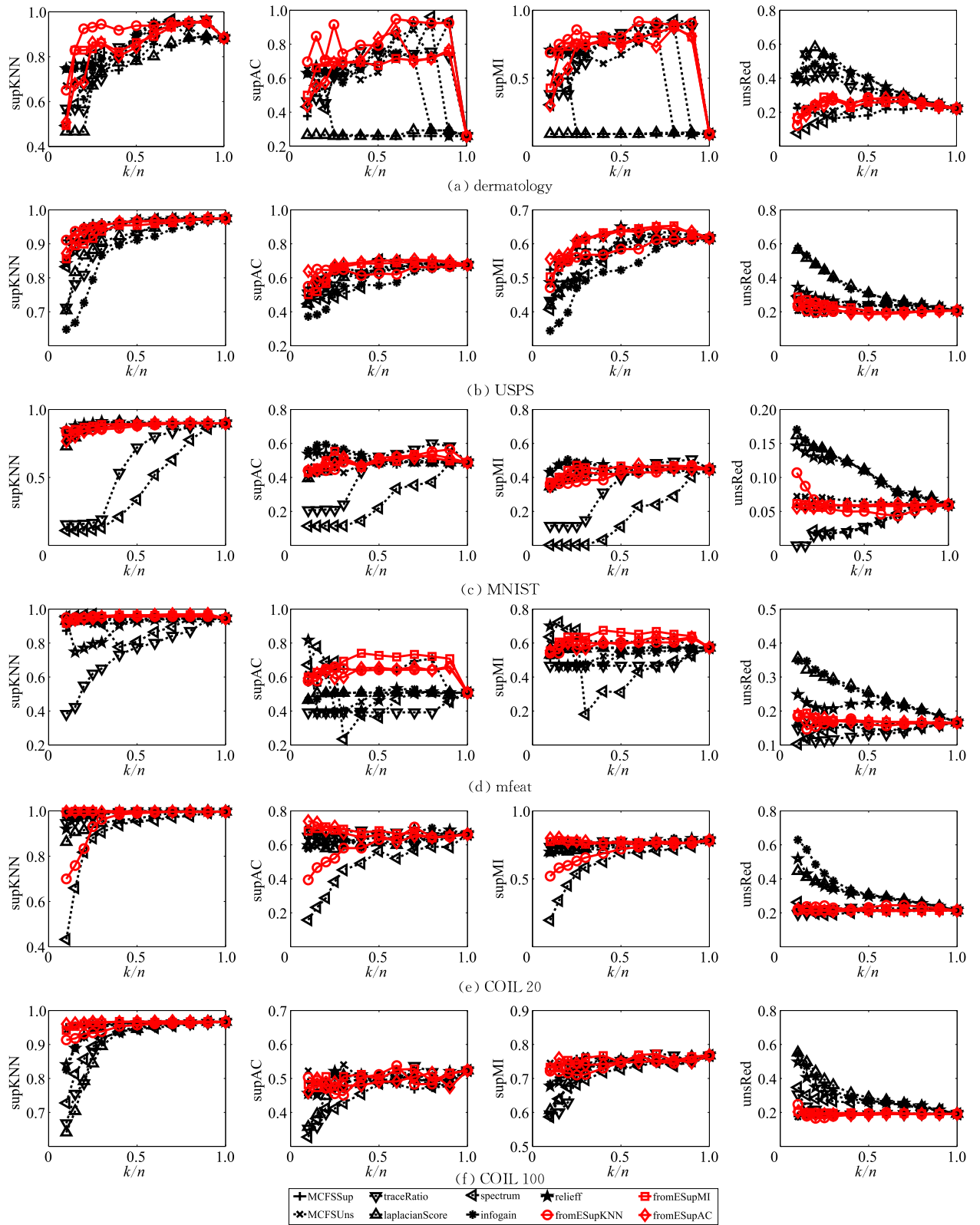
**Output:** A weight function for each feature.

- 1: Compute  $\mathcal{J}(\mathbf{X}, \mathbf{y})$ .
- 2: **for**  $i \in \{1, \dots, n\}$  **do**
- 3:    $F_{-i} = \{1, \dots, n\} \setminus \{i\}$ .
- 4:   Compute  $\mathcal{J}(\mathbf{X}(F_{-i}, :), \mathbf{y})$ .
- 5:   Compute  $w(i) = \mathcal{J}(\mathbf{X}, \mathbf{y}) - \mathcal{J}(\mathbf{X}(F_{-i}, :), \mathbf{y})$ .
- 6: **end for**
- 7: **return**  $w$ .

in terms of the classification accuracies, clustering accuracies, and normalized mutual information for real-world datasets. “MCFSSup” and “MCFSSuns” denote the supervised and unsupervised version of the MSFS algorithm, “infogain” (information gain) and “relief”<sup>[13]</sup> denote two conventional feature selection algorithms, “fromESupKNN”, “fromESupMI”, and “fromESupAC” denote the leave-one-out feature selection algorithm where the evaluation metric  $\mathcal{J}$  is the classification accuracy in Eq. (6) indicated as “supKNN” (using the  $k$ -nearest neighbor classifier with  $k = 1$ ), the clustering accuracy in Eq. (7) indicated as “supAC” (using the best result of the  $K$ -means clustering method for 10 times), and the normalized mutual information in Eq. (8) indicated as “supMI”. The classification accuracy, Eq. (6), in “fromESupKNN” was computed from test data where the datasets were partitioned as  $m_{\text{train}} \approx 0.5m$  and  $m_{\text{test}} \approx 0.5m$ . The clustering accuracy, Eq. (7), and the normalized mutual information, Eq. (8), in “fromESupMI” and “fromESupAC” were computed using all the data. The unsupervised “redundancy rate”<sup>[18]</sup> was computed as an evaluation metric and was defined as  $\mathcal{J}_{\text{red}} = \frac{1}{F(S)} \sum_{i,j \in F(S)} \rho_{i,j}$  where  $\rho_{i,j}$  is the correlation between the  $i$ -th and the  $j$ -th features and  $F(S)$  is the selected feature subset. Note

Table 1 Real-world datasets.

Dataset	Dimension	Size	Classes
dermatology	34	366	6
USPS	256	9298	10
MNIST	784	4000	10
mfeat	649	2000	10
COIL20	1024	1440	20
COIL100	1024	7200	100



**Fig. 1** Classification accuracies (column 1), clustering accuracies (column 2), mutual information (column 3), and redundancy rate (column 4) with different numbers ( $\frac{k}{n} = 0.1, 0.15, 0.20, 0.25, 0.30, 0.40, 0.50, 0.60, 0.70, 0.80, 0.90, 1$ ) of selected features in six real-world datasets (row 1-row 6).

that smaller values of  $\mathcal{J}_{\text{red}}$  indicate that the redundancy is smaller in  $F(S)$  and  $F(S)$  is thus a better feature subset.

Figure 1 shows that the leave-one-out feature selection has better performance in most situations, where the classification accuracy, Eq. (6), is computed from the test data with the datasets partitioned as  $m_{\text{train}} \approx 0.6m$  and  $m_{\text{test}} \approx 0.4m$  and the clustering accuracy, Eq. (7), the normalized mutual information, Eq. (8) and the redundancy rate are computed using all the data. The leave-one-out algorithm selects features which *directly* maximize these classification accuracies, clustering accuracies, and normalized mutual information while preserves the interactions among features in the selection process. Other feature selection algorithms use specific objective functions such as Eqs. (1)-(4) which are reasonable in some situations, but only indirectly solve Eq. (9). The gap between Eq. (9) and the specific feature selection objectives such as in Eqs. (1)-(4) can be quite large, which is why these algorithms are inefficient in some situations.

## 5 Conclusions

In the past several decades, feature selection has drawn much attention due to its fundamental importance in pattern recognition and the difficulty in optimizing the process. Different from previous algorithms which are built on some specific objective functions, this paper presents a leave-one-out feature selection algorithm that directly optimizes more basic feature selection objectives such as the classification and clustering accuracies. The leave-one-out algorithm borrows the concept from conditional independence structures that removing important features usually causes larger changes in the conditional distribution than removing redundant features, since the label variable is determined by a small set of variables. With this observation, the leave-one-out algorithm defines the score of each feature as the performance change with respect to the absence of the feature from the full feature set. The importance to the objectives is then reflected in the scores, which overcomes the limitation of the greedy backward elimination algorithm that loses interactions among features in the selection process. Experiments on real-world datasets with various kinds of feature evaluation metrics verify the effectiveness of the leave-one-out algorithm.

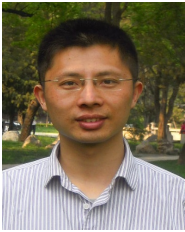
## Acknowledgements

We thank the editors and reviewers for their constructive suggestions and comments. This work was supported by the National Natural Science Foundation of China (Nos. 61071131 and 61271388), the Beijing Natural Science Foundation (No. 4122040), the Research Project of Tsinghua University (No. 2012Z01011), and the Doctoral Fund of the Ministry of Education of China (No. 20120002110036).

## References

- [1] I. Guyon and A. Elisseeff, An introduction to variable and feature selection, *J. Mach. Learn. Res.*, vol. 3, nos. 7-8, pp. 1157-1182, Oct. 2003.
- [2] H. Peng, F. Long, and C. Ding, Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy, *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 27, no. 8, pp. 1226-1238, Aug. 2005.
- [3] M. C. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006.
- [4] J. Li and M. Sun, Non-independent term selection for Chinese text categorization, *Tsinghua Science and Technology*, vol.14, no.1, pp. 113-120, Feb. 2009.
- [5] F. Nie, H. Huang, X. Cai, and C. Ding, Efficient and robust feature selection via joint  $\ell_2$ ,  $\ell_1$ -norms minimization, in *Advances in Neural Information Processing Systems 23*, Vancouver, BC, Canada, 2010, pp. 1813-1821.
- [6] Z. Zhao, L. Wang, H. Liu, and J. Ye, On similarity preserving feature selection, *IEEE Trans. Knowl. Data Eng.*, vol. 25, no. 3, pp. 619-632, Mar. 2013.
- [7] C. Hou, F. Nie, D. Yi, and Y. Wu, Feature selection via joint embedding learning and sparse regression, in *Proc. 22nd Int. Joint Conf. on Artificial Intelligence*, Barcelona, Spain, 2011, pp. 1324-1329.
- [8] S. Xiang, F. Nie, G. Meng, C. Pan, and C. Zhang, Discriminative least squares regression for multiclass classification and feature selection, *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 23, no. 11, pp. 1738-1754, Nov. 2012.
- [9] M. Belkin and M. Niyogi, Laplacian eigenmaps for dimensionality reduction and data representation, *Neural Comput.*, vol. 15, no. 6, pp. 1373-1396, Jun 2003.
- [10] S. Roweis and L. Saul, Nonlinear dimensionality reduction by locally linear embedding, *Science*, vol. 290, no. 5500, pp. 2323-2326, Dec. 2000.
- [11] R. Duda, P. Hart, and D. Stork, *Pattern Classification*, Wiley-Interscience, 2001.
- [12] X. He, D. Cai, and P. Niyogi, Laplacian score for feature selection, in *Advances in Neural Information Processing Systems 18*, Vancouver, BC, Canada, 2006, pp. 507-514.
- [13] I. Kononenko, Estimating attributes: Analysis and extensions of relief, in *Machine Learning: ECML-94*, Springer, 1994, pp. 171-182.
- [14] Z. Zhao and H. Liu, Spectral feature selection for supervised and unsupervised learning, in *Proc. 24th Int. Conf. on Machine Learning*, Corvallis, USA, 2007, pp. 1151-1157.

- [15] F. Nie, S. Xiang, Y. Jia, C. Zhang, and S. Yan, Trace ratio criterion for feature selection, in *Proc. 23rd AAAI Conf. on Artificial Intelligence*, Chicago, USA, 2008, pp. 671-676.
- [16] D. Cai, C. Zhang, and X. He, Unsupervised feature selection for multi-cluster data, in *Proc. 16th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, Washington DC, USA, 2010, pp. 333-342.
- [17] J. Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1988.
- [18] Z. Zhao, L. Wang, and H. Liu, Efficient spectral feature selection with minimum redundancy, in *Proc. 24th AAAI Conf. on Artificial Intelligence*, Atlanta, USA, 2010, pp. 673-678



**Dingcheng Feng** received his BS degree in automation from Harbin Institute of Technology, Harbin, China, in 2008. He is currently working toward his PhD degree in control science and engineering in the Department of Automation, Tsinghua University, Beijing, China. His research interests include structure learning, variational inference, causal inference, and feature selection.



**Feng Chen** received his BS and MS degrees in automation from Saint-Petersburg Poly-technic University, Saint-Petersburg, Russia, in 1994 and 1996, respectively, and his PhD degree in automation from Tsinghua University, Beijing, China, in 2000. He is currently a professor with the National Laboratory

for Information Science and Technology, Department of Automation, Tsinghua University. His research interests include computer vision, video processing, variational inference, and structure learning in graphical models.



**Wenli Xu** received his BS degree in electrical engineering and his MS degree in automatic control engineering from Tsinghua University, Beijing, China, in 1970 and 1980, respectively, and his PhD degree in electrical and computer engineering from the University of Colorado at Boulder, in 1990. He is currently a professor with the National Laboratory for Information Science and Technology, Department of Automation, Tsinghua University. His research interests include automatic control, intelligent robot, and computer vision.