# Discriminative sparse representation learning using multiclass hinge loss

Ryota Kamiya and Yoshikazu Washizawa The University of Electro-Communications, Tokyo, Japan. E-mail: k1731051@edu.cc.uec.ac.jp

Abstract-Sparse representation methods have been researched widely in recent years. Sparse representation classification methods, such as sparse representation classifier (SRC) and labelconsistent K-SVD (singular value decomposition) learn classification parameters, dictionary, and sparse representation simultaneously, so that they find an optimal sparse representation to discriminate categories. However, these classifiers use least square error (LSE) strategy to design the classifiers. LSE of the empirical risk is not optimal for classifiers because even if a training sample correctly classified, it may increase the empirical cost. We, therefore, introduce the hinge loss to design the classifier. The hinge loss is employed in support vector machines, and it shows better performance than LSE based methods. We provide an optimization algorithm to minimize the proposed criterion that is the linear combination of the hinge loss and sparse representation error. Experimental results show that the proposed method exhibited conventional sparse representation classification methods.

#### I. INTRODUCTION

As the internet of things (IoT) grows, the redundancy of data increases. Exploration of latent information is important, and the sparse representation is one of the promising tools in signal processing and machine learning. Sparse representation compactly expresses data in a vector having a few non-zero elements using an overcomplete dictionary, and it has been widely applied in many fields such as image classification and recognition, feature learning, image noise removal, and sensing [1], [3]. Compressed sensing is a method to acquire and reconstruct signals in compressed form using sparse representation [3]. The dictionary learning methods such as K-SVD (singular value decomposition) have also been proposed to obtain the overcomplete dictionary using given dataset, and applied to image compression and denoising [6].

For classification of redundant input data, sparse representation classification methods have been proposed. The sparse representation classifier (SRC) classifies input vectors using reconstruction residual of sparse representation for each class [2]. The label consistent K-SVD (LC-KSVD) is a classification method using K-SVD. LC-KSVD simultaneously minimizes reconstruction error of the sparse representation, and empirical classification error using class labels of training data [11]. LC-KSVD has been extended to the joint embedding and dictionary learning (JEDL) and the locality preserving KSVD (LP-KSVD) [4], [5]. However, these classifiers use the squared error between class labels and the output of classifiers, that does not evaluate misidentification directly. The squared error may be increased even if a training data is correctly classified. In this paper, we therefore, introduce the multiclass hinge loss for sparse representation classifier. The hinge loss is employed in support vector machines (SVMs) and shows better performance than least square error for classification problems since the hinge loss is always zero if a training data is correctly classified with appropriate margin. We define a new optimization criterion to minimize both reconstruction error of the sparse representation and misidentification error using the hinge loss. We provide an optimization algorithm to minimize our criterion.

### II. RELATED WORKS

We briefly introduce sparse expression methods of sparse representation classifier (SRC), K-SVD and LC-KSVD1 and LC-KSVD2.

### A. Sparse representation classifier (SRC)

SRC is a classifier using the sparse representation [2]. First, SRC obtains dictionaries  $D = [d_1, \ldots, d_K] \in \mathbb{R}^{n \times K}$  for each class using labeled training data, where *n* is the number of input dimension, and *K* is the number of dictionary vectors (atoms). For unknown input *x*, its sparse representation *s* is obtained by,

$$\min \|\boldsymbol{D}\boldsymbol{s} - \boldsymbol{x}\|^2 + \lambda_1 \|\boldsymbol{s}\|_1, \tag{1}$$

where  $\lambda_1$  is a parameter for trading-off the reconstruction error and the sparsity. Then, the residual  $r_i(\boldsymbol{x}) = \|\boldsymbol{x} - \boldsymbol{D}\delta_i(\boldsymbol{s})\|_2$  is obtained for each class i = 1, 2, ..., C, where  $\delta_i(\boldsymbol{s})$  is set to zero except for non-zero elements of  $\boldsymbol{s}$  associated with class  $i \in 1, 2, ..., C$ .  $\boldsymbol{x}$  is assigned to the class having the minimum residual  $r_i(\boldsymbol{x}), i = 1, ..., C$ .

#### B. K-SVD

K-SVD is a method to learn an overcomplete dictionary for sparse representation by minimizing the reconstruction error [6]. Let  $\boldsymbol{X} = [\boldsymbol{x}_1, \dots, \boldsymbol{x}_N] \in \mathbb{R}^{n \times N}$  be a set of *n*-dimensional training vector, where N is the number of training data. Then the optimization problem of K-SVD is

$$\min_{\boldsymbol{D},\boldsymbol{S}} \|\boldsymbol{X} - \boldsymbol{D}\boldsymbol{S}\|_F^2, \quad \text{subject to } \|\boldsymbol{s}_i\|_0 \le T_1, \qquad (2)$$

where  $S = [s_1, ..., s_N] \in \mathbb{R}^{K \times N}$  is a matrix of sparse vectors,  $\|\cdot\|_0$  is a  $l_0$  norm which is the number of non-zero elements,  $T_1$  is the sparsity constraint factor.

K-SVD alternately updates S and D. The update of the sparse matrix S by fixing the dictionary D is done for

## C. LC-KSVD

The label consistent K-SVD (LC-KSVD) is a classification method using K-SVD [12]. LC-KSVD considers additional two terms to be minimized.

$$\min_{\boldsymbol{D},\boldsymbol{S},\boldsymbol{A},\boldsymbol{W}} \|\boldsymbol{X} - \boldsymbol{D}\boldsymbol{S}\|_{F}^{2} + \alpha \|\boldsymbol{Q} - \boldsymbol{A}\boldsymbol{S}\|_{F}^{2} + \beta \|\boldsymbol{H} - \boldsymbol{W}\boldsymbol{S}\|_{F}^{2}$$
  
subject to  $\|\boldsymbol{s}_{i}\|_{0} \leq T_{1}, \quad i = 1, \dots, N.$ 
(3)

 $\boldsymbol{Q} \in \mathbb{R}^{K \times N}$  is the discriminative sparse code of the input signal for classification.

$$[\boldsymbol{Q}]_{i,j} = \begin{cases} 1 & \text{(the class label of } \boldsymbol{d}_i \text{ and } \boldsymbol{s}_j \text{ are the same}) \\ 0 & \text{(otherwise).} \end{cases}$$
(4)

 $A \in \mathbb{R}^{K \times K}$  is a linear transform, so that the sparse vectors having the same class label are linearly transformed to similar code vectors.  $\alpha > 0$  is a parameter to define the strength of the second term.

The third term is to design a classifier.  $W = [w_1, ..., w_K] \in \mathbb{R}^{C \times K}$  is the classifier to be trained, and  $H = [h_1, ..., h_N] \in \mathbb{R}^{C \times N}$  is a class label matrix of input signals, i.e., the *j*th element of  $h_i$  is one if  $h_i$  belongs the *j*th class, and otherwise zero.  $\beta$  is a parameter to define the strength of the third term. LC-KSVD1 is the case of  $\beta = 0$ , otherwise, it is denoted by LC-KSVD2. Jiang et al. [12] provided a K-SVD based algorithm to minimize the optimization problem (3).

For unknown input signal  $x \in \mathbb{R}^n$ , its sparse representation  $s \in \mathbb{R}^K$  is obtained. LC-KSVD2 estimates the class label as the largest element of Ws. In the LC-KSVD1, the class label is estimated in the same manner as SRC.

The squared error used in (3) is not optimal for classification. Let us consider a position that Q or H is one. If ASor WS has a value that is greater than one in the position, the total cost increases although it is not disadvantageous for classification. In a similar manner, if AS or WS has a value that is smaller than zero in the position that Q or H is zero, the total cost increases although it is not disadvantageous for classification. Therefore, the squared error criteria do not evaluate misidentification.

# III. SPARSE REPRESENTATION CLASSIFIER USING HINGE LOSS

We introduce the hinge loss to the sparse representation classifier. The hinge loss directly evaluates the misidentification, and is used in SVMs. Let  $y \in \{-1, +1\}$  be a binary class label, x be corresponding input vector, and f(x) be a decision function of a binary classifier, i.e. if f(x) > 0, x is



Fig. 1. hinge loss of binary classification



Fig. 2. multiclass hinge loss

assigned to the positive class, and otherwise, x is assigned to the negative class. The hinge loss is defined by

$$L(f) = \begin{cases} 1 - yf(\boldsymbol{x}) & (yf(\boldsymbol{x}) \le 1) \\ 0 & (yf(\boldsymbol{x}) \ge 1). \end{cases}$$
(5)

If yf(x) > 0, the sample x is correctly classified by f. The hinge loss has the margin one, and linearly give the penalty of  $yf(x) \le 1$  (Fig. 1).

The hinge loss is extended for multi-class problems. Let  $f_i(x)$  is the decision function for the *i*th class, i.e., x is assigned to  $\arg \max_{i=1,...,C} f_i(x)$ . The hinge loss for multi-class is

$$L(f) = \max\{0, 1 - (f_c(\boldsymbol{x}) - v)\},\tag{6}$$

$$w = \max_{i \in \{1, \dots, C\} \setminus \{c\}} f_i(\boldsymbol{x}) \tag{7}$$

where c is the class label of x. v is the maximum value of decision functions of the other classes. If the difference between the decision function value of the true class label c,  $f_c(x)$  and v is greater than one, L(f) linearly gives the penalty (Fig. 2).

We define the linear multiclass decision function for sparse vector  $\boldsymbol{s}, f_i(\boldsymbol{s}) = \boldsymbol{w}_i^\top \boldsymbol{s} + b_i$ , and let  $\boldsymbol{W} = [\boldsymbol{w}_1, \dots, \boldsymbol{w}_C] \in \mathbb{R}^{K \times C}, \ \boldsymbol{b} = [b_1, \dots, b_C] \in \mathbb{R}^C$ . Then the optimization problem of the proposed method is

$$\min_{\boldsymbol{D},\boldsymbol{S},\boldsymbol{W},\boldsymbol{b}} J = \sum_{i=1}^{N} \|\boldsymbol{x}_{i} - \boldsymbol{D}\boldsymbol{s}_{i}\|^{2} + \lambda \|\boldsymbol{s}_{i}\|_{1} + \mu L(\boldsymbol{W},\boldsymbol{b},\boldsymbol{s}_{i}), \quad (8)$$

where  $L(W, b, s_i)$  is the multiclass hinge loss function.

We alternately update  $(\boldsymbol{W}, \boldsymbol{b}), \boldsymbol{s}_i, i = 1, \dots, N$ , and  $\boldsymbol{D}$ .

a) Update of W and b: Let us optimize J for W and b while D and  $s_n$  n = 1, ..., N are fixed. Then the first and second terms are constant for W and b. Then we consider the subproblem,

$$\min_{\boldsymbol{W},\boldsymbol{b}} \sum_{i=1}^{N} L(\boldsymbol{W}, \boldsymbol{b}, \boldsymbol{s}_i).$$
(9)

The multiclass hinge loss (6) can be rewritten by using slack variable  $\xi$ ,

$$L(f) = \min \xi$$
  
subject to  $\xi \ge 0$   
 $\xi \ge 1 - (f_c(\boldsymbol{x}) - f_i(\boldsymbol{x})), \ i \in \{1, \dots, C\} \setminus \{c\},$   
(10)

Then the subproblem (13) is reduced to the following linear programming,

$$\min_{\boldsymbol{W},\boldsymbol{b},\boldsymbol{\xi}} \sum_{i=1}^{N} \xi_{i}$$
subject to  $\xi_{i} \geq 0, \ i = 1, \dots, N$ 

$$\xi_{i} \geq 1 - (\boldsymbol{w}_{c} - \boldsymbol{w}_{k})^{\top} \boldsymbol{s}_{i} + (b_{c} - b_{k}),$$

$$i = 1, \dots, N, \ k \in \{1, \dots, C\} \setminus \{c\}.$$
(11)

The solution can be obtained by a solver such as GLPK. Alternatively, this is a special case of multiclass SVM [13], and its implementations such as liblinear [14] obtain the solution.

b) Update of  $s_n$ : The sparse vector  $s_i$  can be optimized for each *i*. We, here omit the subscript *i*, and let  $s = s_+ - s_-$ ,  $(s_+ \ge 0, s_- \ge 0)$  and  $z = [s_+^\top | s_-^\top]^\top \in \mathbb{R}^{2K}$   $(z \ge 0)$ , where " $\ge$ " is for each element of vectors. Then the second term, the  $l_1$  regularization, of *J* can be rewritten to

$$\|\boldsymbol{s}\|_1 = \min_{\boldsymbol{z}} \mathbf{1}_{2K}^{\top} \boldsymbol{z}, \text{ subject to } \boldsymbol{z} \ge 0, \ [\boldsymbol{I}_K| - \boldsymbol{I}_K] \boldsymbol{z} = \boldsymbol{s},$$

where  $\mathbf{1}_{2K}$  is a (2K)-dimensional vector whose elements are one, and  $\mathbf{I}_K$  is the identity matrix of size K.

Let  $\tilde{w}_i = [w_i^{\top}| - w_i^{\top}]$ . Then the third term of J is transformed by using Eq. (10),

$$L(\boldsymbol{W}, \boldsymbol{b}, \boldsymbol{s}) = \min \xi$$
  
subject to  $\xi \ge 0$ ,  
 $\xi \ge 1 - (\tilde{\boldsymbol{w}}_c - \tilde{\boldsymbol{w}}_i)^\top \boldsymbol{z} + (b_c - b_i)$ , (12)  
 $i \in \{1, \dots, C\} \setminus \{c\}$ ,

where c is the class label of x.

By letting  $\tilde{z} = [z^{\top}|\xi]^{\top} \in \mathbb{R}^{2K+1}$ , and  $\tilde{D} = [D| - D]$ , the optimization problem (8) is reduced to the following constrained quadratic programming,

$$\min_{\tilde{\boldsymbol{z}}} \tilde{\boldsymbol{z}}^{\top} \begin{bmatrix} \tilde{\boldsymbol{D}}^{\top} \tilde{\boldsymbol{D}} & \boldsymbol{0}_{2K} \\ \boldsymbol{0}_{2K}^{\top} & \boldsymbol{0} \end{bmatrix} \tilde{\boldsymbol{z}} + [(\lambda \mathbf{1}_{2K}^{\top}) - 2\boldsymbol{x}^{\top} \tilde{\boldsymbol{D}}] \mu] \tilde{\boldsymbol{z}}$$
subject to  $\tilde{\boldsymbol{z}} \ge 0$ 

$$[(\tilde{\boldsymbol{w}}_{c}^{\top} - \tilde{\boldsymbol{w}}_{i}^{\top})] 1] \tilde{\boldsymbol{z}} \ge 1 - (b_{c} - b_{i}) i \in \{1, \dots, C\} \setminus \{c\}.$$

$$(13)$$

c) Update of D: Since the second and third terms of J are constant for D, D is updated by the least square,

$$\boldsymbol{D}^{\text{new}} = \underset{\boldsymbol{D}}{\arg\min} \|\boldsymbol{X} - \boldsymbol{D}\boldsymbol{S}\|_{F}^{2}$$
(14)

$$= \boldsymbol{X}\boldsymbol{S}^{\top}(\boldsymbol{S}\boldsymbol{S}^{\top})^{-1}.$$
 (15)

D is initiated by randomly selecting vectors from X, then D and S are initiated by the standard K-SVD [6]. We summarize our algorithm in Algorithm 1.

Algorithm 1	Optimization	algorithm	of the	proposed	method
Require: X					

Ensure: D, S, W, b

1: Initialize  $oldsymbol{D}^{(0)}$  by selecting randomly from  $oldsymbol{X}$ 

2: Compute  $D^{(0)}$  and  $S^{(0)}$  by using K-SVD

3: repeat

- 4: Update  $W^{(t)}$ ,  $b^{(t)}$  by (11);
- 5: for i = 1 to N do
- 6: Obtain  $\tilde{z}$  by Eq. (13)
- 7: Update  $oldsymbol{s}_i = [oldsymbol{I}_K| oldsymbol{I}_K]oldsymbol{z}$

8: end for

9: Compute  $D^{(t)}$  by (15);

10: **until** J converges

### IV. EXPERIMENT

We used the Isolet spoken letter recognition database. It has 150 subjects speaking the name of each letter of the alphabet twice. The described features contain spectral coefficients, contour features, sonorant features, pre-sonorant features, and postsonorant features. The number of input dimension n is 617, and the number of classes C is 26. The speakers are grouped into sets of 30 speakers each and are named Isolet 1 to Isolet 5. In the experiment, Isolet 1 to Isolet 5 are used to evaluate the performance of sparse codes expression classifiers. We conducted five fold cross validation for Isolet 1 to Isolet 5, i.e., each dataset is divided into five subsets, and four subsets are used for training, and remaining one subset is used for testing. We evaluated averaged classification accuracy for five folds.  $\alpha$  and  $\beta$  refer to the parameters the reference paper. The sparsity  $T_1$  and dictionary vector K used grid search. For the experiment the parameter setting was as follows;  $\alpha = 3.0$ ,  $\beta = 4.0$ , the sparsity  $T_1 = 10$ , and the number of dictionary vectors K = 25.

We show the experimental results in Table. I. The classification accuracy is almost the same for SRC, LCKSVD1 and LCKSVD2.

Table. II represents the number of nonzero elements in a sparse vector. We used OMP for SRC, LCKSVD1, and LCKSVD2, that is, the number of nonzero elements is directly determined by the hyperparameter. For this reason the number of nonzero elements in sparse vectors are the same for SRC and LC-KSVDs. The our method improves the identification rate more than the related method. Moreover the number of nonzero elements of the proposed method is large.

TABLE I CLASSIFICATION ACCURACY

CERISSIFICATION ACCORACT							
	Method						
	SRC	LCKSVD1					
	Mean ±STD (%)	Mean ±STD (%)					
Isolet1	91.73±1.44	89.16±2.91					
Isolet2	$89.87 \pm 1.19$	$86.85 \pm 1.46$					
Isolet3	85.32±2.72	79.61±2.82					
Isolet4	$80.10 \pm 3.57$	$76.63 \pm 3.16$					
Isolet5	88.77 ±1.65	83.96±2.16					
Mean	$87.16 \pm 2.11$	$86.60 \pm 2.50$					
	LC-KSVD2	Our Model					
	Mean ±STD (%)	Mean ±STD (%)					
Isolet1	$90.57 \pm 2.10$	$95.83 \pm 0.32$					
Isolet2	$87.56 \pm 1.45$	91.45±1.17					
Isolet3	$81.73 \pm 1.63$	$87.65 \pm 2.40$					
Isolet4	79.59 ±1.14	$88.38 \pm 0.55$					
Isolet5	$86.01 \pm 0.85$	86.71±2.36					
Mean	84 85+1 43	$90.04 \pm 1.36$					

 TABLE II

 NUMBER OF NON-ZERO IN SPARSE REPRESENTATION

SRC	10
LC-KSVD1	10
LC-KSVD2	10
Our method	11

# V. CONCLUSION AND FUTURE WORK

We have proposed a sparse representation classification method using multiclass hinge loss. The hinge loss directly evaluate the empirical misidentification risk, and optimization is reduced to a linear programming. We provided an algorithm to minimize the criterion which is the linear combination of sparse representation error and the hinge loss. The algorithm alternately updates weights for classifier W and b, the sparse representation S, and the dictionary D, thus the sparse representation of input vector x is not only approximation of x, but also discriminative representation.

Future works include the investigation of the relation between sparsity level, dictionary size, and classification accuracy, and a method using the  $l_0$  norm constraint.

#### REFERENCES

- [1] M. Elad, Sparse and Redundant Representations, Springer, 2010.
- [2] J. Wang, C. Y. Lu, M. Wang, P.P. Li, S.C. Yan, and X.G. Hu, "Robustface recognition via adaptive sparse representation," *IEEE Trans. Cybern.*, vol. 44, pp. 2368–2378, 2014.
- [3] D.L. Donoho, "Compressed sensing," *IEEE Trans. Inf. Theory* vol. 52, pp. 1289-1306, 2006.
  [4] Z. Zhang, F. Li, T.W.S. Chow, L. Zhang, and S. Yan, "Sparse codes auto-
- [4] Z. Zhang, F. Li, T.W.S. Chow, L. Zhang, and S. Yan, "Sparse codes autoextractor for classification: a joint embedding and dictionary learning framework for representation," *IEEE Trans. Signal Process.* vol. 64, pp. 3790-3805, 2016.
- [5] Y.-S. Lee, C.-Y. Wang, S. Mathulaprangsan, J.-H. Zhao, and J.-H. Zhao, "Locality-preserving K-SVD Based Joint Dictionary and Classifier Learning for Object Recognition," *Proceedings of the 2016 ACM on Multimedia Conference*, pp. 481-485, 2016.
- [6] M. Aharon, M. Elad, and A.M. Bruckstein, "K-SVD: An algorithm for designing of overcomplete dictionaries for sparse representation," *IEEE Trans. Signal Process.* vol. 54, no. 11, pp. 4311-4322, 2006.
- [7] S.S. Chen, D.L. Donoho, and M.A. Saunders, "Atomic decomposition by basis pursuit," *SIAM Journal on Scientific Computing*, vol. 20, no. 1, pp. 33-61, 1998.

- [8] Y. C. Pati, R. Rezaiifar and S. K. Perinkulam, "Orthogonal matching pursuit: recursive function approximation with applications to wavelet decomposition." *IEEE Signals, Systems and Computers*, 1993. 1993 Conference Record of The Twenty-Seventh Asilomar Conference on., 1993, pp.40-44.
- [9] D. R. Musicant, "MATLAB/CPLEX MEX-Files", Computer Sciences Department, University of Wisconsin, Madison, www.cs.wisc.edu/~musicant/data/cplex/., 2000
- [10] Z. Qiang, and B. Li, "Discriminative K-SVD for dictionary learning in face recognition." *Computer Vision and Pattern Recognition (CVPR)*, 2010 IEEE Conference on., IEEE, 2010, pp.2691-2698.
- [11] Z. Jiang, Z. Lin, and L. S. Davis, "Learning a discriminative dictionary for sparse coding via label consistent K-SVD," *Computer Vision and Pattern Recognition, 2011 IEEE Conference on.*, IEEE, 2011, pp.1697-1704.
- [12] Z. L. Jiang, Z. Lin, and L. S. Davis, "Label consistent K-SVD: Learning a discriminative dictionary for recognition," *IEEE transactions* on pattern analysis and machine intelligence, vol.35, 2013, pp.2651-2664.
- [13] K. Crammer and Y. Singer, "On the algorithmic implementation of multiclass kernel-based vector machines," *Journal of Machine Learning Research*, vol 2, pp. 265–292, 2001.
- [14] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin, "LIBLINEAR: A library for large linear classification," *Journal of Machine Learning Research*, vol 9, pp. 1871–1874, 2008.