

A greedy pursuit approach to classification using multi-task multivariate sparse representations

Charles Jeon^{*}, Umamahesh Srinivas[†], Vishal Monga[†]

^{*} Department of Electrical and Systems Engineering, University of Pennsylvania, USA

[†] Department of Electrical Engineering, Pennsylvania State University, USA

July 2012

Abstract

In this report, we propose an extension of the well-known Simultaneous Orthogonal Matching Pursuit (SOMP) algorithm to solve a multi-task multivariate classification problem using sparse representations. This corresponds to the case where an event is described by multiple representations, and separate training dictionaries are designed for each such representation.

I. NOTATION

Let $\mathbf{y}_i \in \mathbb{R}^m, i = 1, \dots, T$ be T different representations¹ of the same physical event, which is to be classified into one of C different classes. Let $\mathbf{Y} := [\mathbf{y}_1 \dots \mathbf{y}_T] \in \mathbb{R}^{m \times T}$. Assuming n training samples/events in total, we design T dictionaries $\mathbf{D}_i \in \mathbb{R}^{m \times n}, i = 1, \dots, T$, corresponding to the T representations. We define a new composite dictionary $\mathbf{D} := [\mathbf{D}_1 \dots \mathbf{D}_T] \in \mathbb{R}^{m \times nT}$. Further, each dictionary \mathbf{D}_i is represented as the concatenation of the sub-dictionaries from all classes corresponding to the i -th representation of the event:

$$\mathbf{D}_i := [\mathbf{D}_i^1 \mathbf{D}_i^2 \dots \mathbf{D}_i^C], \quad (1)$$

where \mathbf{D}_i^j represents the collection of training samples for representation i that belong to the j -th class. So, we have:

$$\mathbf{D} := [\mathbf{D}_1 \dots \mathbf{D}_T] = [\mathbf{D}_1^1 \mathbf{D}_1^2 \dots \mathbf{D}_1^C \dots \mathbf{D}_T^1 \mathbf{D}_T^2 \dots \mathbf{D}_T^C]. \quad (2)$$

An important assumption in designing \mathbf{D} is that the k -th column from each of the dictionaries $\mathbf{D}_i, i = 1, \dots, T$, taken together offer multiple representations of the k -th training sample/event.

II. MULTI-TASK MULTIVARIATE SPARSE REPRESENTATIONS

A test event \mathbf{Y} can now be represented as a linear combination of training samples as follows:

$$\mathbf{Y} = [\mathbf{y}_1 \dots \mathbf{y}_T] = \mathbf{D}\mathbf{S} = [\mathbf{D}_1^1 \mathbf{D}_1^2 \dots \mathbf{D}_1^C \dots \mathbf{D}_T^1 \mathbf{D}_T^2 \dots \mathbf{D}_T^C] [\boldsymbol{\alpha}_1 \dots \boldsymbol{\alpha}_T], \quad (3)$$

where the coefficient vectors $\boldsymbol{\alpha}_i \in \mathbb{R}^{nT}, i = 1, \dots, T$, and $\mathbf{S} = [\boldsymbol{\alpha}_1 \dots \boldsymbol{\alpha}_T] \in \mathbb{R}^{nT \times T}$.

We examine the structure of the coefficient matrix \mathbf{S} and make some crucial observations.

- It is reasonable to assume that the i -th representation of the test event (i.e. \mathbf{y}_i) can be approximately represented by the linear span of the training samples belonging to the i -th representation alone (i.e. only those training samples in \mathbf{D}_i). So the columns of \mathbf{S} have the following structure: each vector $\boldsymbol{\alpha}_i$ has non-zero coefficients only in the locations corresponding to the columns of \mathbf{D}_i and has zeros elsewhere. As a result, \mathbf{S} exhibits block-diagonal structure.
- Each representation \mathbf{y}_i of the test event is a *sparse* linear combination of the training samples in \mathbf{D}_i . Suppose the event belongs to class $c \in \{1, \dots, C\}$; then only those coefficients in $\boldsymbol{\alpha}_i$ that correspond to \mathbf{D}_i^c are expected to be non-zero.
- Furthermore, the non-zero weights of training samples in the linear combination exhibit one-to-one correspondence across representations. If the j -th training sample from the c -th class in \mathbf{D}_1 has a non-zero contribution to \mathbf{y}_1 , then for all $i = 2, \dots, T$, \mathbf{y}_i has non-zero contributions from the j -th training sample of the c -th class in \mathbf{D}_i .

This suggests a joint sparsity model similar to the model introduced in [1]. However, the multi-task nature of the problem with different dictionaries \mathbf{D}_i does not permit us to apply the SOMP algorithm from [1] directly. Since \mathbf{S} obeys column correspondence, we introduce a new matrix $\mathbf{S}' \in \mathbb{R}^{n \times T}$ as the transformation of \mathbf{S} with the zero coefficients removed,

$$\mathbf{S}' = \begin{bmatrix} \boldsymbol{\alpha}_1^1 & \dots & \boldsymbol{\alpha}_i^1 & \dots & \boldsymbol{\alpha}_T^1 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \boldsymbol{\alpha}_1^C & \dots & \boldsymbol{\alpha}_i^C & \dots & \boldsymbol{\alpha}_T^C \end{bmatrix},$$

where $\boldsymbol{\alpha}_i^j$ refers to the sub-vector extracted from $\boldsymbol{\alpha}_i$ that corresponds to coefficients from the j -th class. Note that, in the i -th column of \mathbf{S}' , only the coefficients corresponding to \mathbf{D}_i are retained (for $i = 1, \dots, T$).

¹The term *multi-task* is used to refer to these multiple representations in some application domains such as heterogeneous sensor fusion.

We can now apply row-sparsity constraints similar to the approach in [1]. Our modified optimization problem becomes:

$$\hat{\mathbf{S}}' = \arg \min_{\mathbf{S}'} \|\mathbf{S}'\|_{\text{row-0}} \quad \text{subject to} \quad \|\mathbf{Y} - \mathbf{D}\mathbf{S}'\|_F \leq \epsilon, \quad (4)$$

for some tolerance $\epsilon > 0$. We minimize the number of non-zero rows, while the constraint guarantees a good approximation.

The matrix \mathbf{S} can be transformed into \mathbf{S}' by introducing matrices $\mathbf{H} \in \mathbb{R}^{nT \times T}$ and $\mathbf{J} \in \mathbb{R}^{n \times nT}$,

$$\mathbf{H} = \text{diag}[\mathbf{1} \ \mathbf{1} \ \dots \ \mathbf{1}], \mathbf{J} = [\mathbf{I}_n \ \mathbf{I}_n \ \dots \ \mathbf{I}_n],$$

where $\mathbf{1} \in \mathbb{R}^n$ is the vector of all ones, and \mathbf{I}_n denotes the n -dimensional identity matrix. Finally, we obtain $\mathbf{S}' = \mathbf{J}(\mathbf{H} \circ \mathbf{S})$, where \circ denotes the Hadamard product, $(\mathbf{H} \circ \mathbf{S})_{ij} \triangleq h_{ij}s_{ij}$ for all i, j .

III. EXTENSION OF SOMP FOR MULTI-TASK MULTIVARIATE SPARSE REPRESENTATIONS

Eq. (4) represents a hard optimization problem due to presence of the non-invertible transformation from \mathbf{S} to \mathbf{S}' . We bypass this difficulty by proposing a modified version of the SOMP algorithm for the multi-task multivariate case.

Recall that the original SOMP algorithm gives K distinct atoms (assuming K iterations) from a dictionary \mathbf{D} that best represent the data matrix \mathbf{Y} . In every iteration k , SOMP measures the residual for each atom in \mathbf{D} and creates an orthogonal projection with maximal correlation. Extending this to the multi-task setting, for every representation $i, i = 1, \dots, T$, we can identify the index set that gives the highest correlation with the residual at the k -th iteration as follows:

$$\lambda_{i,k} = \arg \max_{j=1, \dots, n} \sum_{q=1}^T w_q \|\mathbf{R}_{k-1}^t \mathbf{d}_{q,j}\|_p, p \geq 1,$$

where w_q denotes the weight (confidence) assigned to the q -th representation, $\mathbf{d}_{q,j}$ represents the j -th column of $\mathbf{D}_q, q = 1, \dots, T$, and the superscript $(\cdot)^t$ indicates the matrix transcript operator. After finding $\lambda_{i,k}$, we modify the index set to:

$$\Lambda_{i,k} = \Lambda_{i,k-1} \cup \{\lambda_{i,k}\}, i = 1, \dots, T.$$

Thus, by finding the index set for the T distinct representations, we can create an orthogonal projection with each of the atoms in their corresponding representations. The algorithm is summarized below in Algorithm 2.

Algorithm 1 SOMP for multi-task multivariate sparse representation-based classification

Input: Dictionary \mathbf{D} as defined in Section I, signal matrix \mathbf{Y} , number of iterations K

Initialization: residual $\mathbf{R}_0 = \mathbf{Y}$, index set $\Lambda_0 = \phi$, iteration counter $k = 1$

while $k \leq K$ **do**

(1) Find the index of the atom that best approximates all residuals:

$$\lambda_{i,k} = \arg \max_{j=1, \dots, n} \sum_{q=1}^T w_q \|\mathbf{R}_{k-1}^t \mathbf{d}_{q,j}\|_p, p \geq 1$$

(2) Update the index set $\Lambda_{i,k} = \Lambda_{i,k-1} \cup \{\lambda_{i,k}\}, i = 1, \dots, T$

(3) Compute the orthogonal projector $\mathbf{p}_{i,k} = \left(\mathbf{D}_{\Lambda_{i,k}}^t \mathbf{D}_{\Lambda_{i,k}} \right)^{-1} \mathbf{D}_{\Lambda_{i,k}}^t \mathbf{y}_i$, for $i = 1, \dots, T$, where $\mathbf{D}_{\Lambda_{i,k}} \in \mathbb{R}^{n \times k}$ consists of the k atoms in \mathbf{D}_i indexed in $\Lambda_{i,k}$

(4) Update the Residual Matrix $\mathbf{R}_k = \mathbf{Y} - [\mathbf{D}_{\Lambda_{1,k}} \mathbf{p}_{1,k} \ \dots \ \mathbf{D}_{\Lambda_{T,k}} \mathbf{p}_{T,k}]$

(5) Increment k : $k \leftarrow k + 1$

end while

Output: Index set $\Lambda_i = \Lambda_{i,K}, i = 1, \dots, T$; sparse representation $\hat{\mathbf{S}}'$ whose non-zero rows indexed for each representation by $\Lambda_i, i = 1, \dots, T$, are the K rows of the matrix $\left(\mathbf{D}_{\Lambda_{i,K}}^t \mathbf{D}_{\Lambda_{i,K}} \right)^{-1} \mathbf{D}_{\Lambda_{i,K}}^t \mathbf{Y}$.

REFERENCES

- [1] J. A. Tropp, A. C. Gilbert, and M. J. Strauss, "Algorithms for simultaneous sparse approximation. Part I: Greedy pursuit," *Signal Processing*, vol. 86, pp. 572–588, Apr. 2006.