

Multi-Dimensional Classification via Sparse Label Encoding

ICML

International Conference
On Machine Learning

Bin-Bin Jia

Min-Ling Zhang

School of Computer Science and Engineering,

MOE Key Laboratory of Computer Network & Information Integration

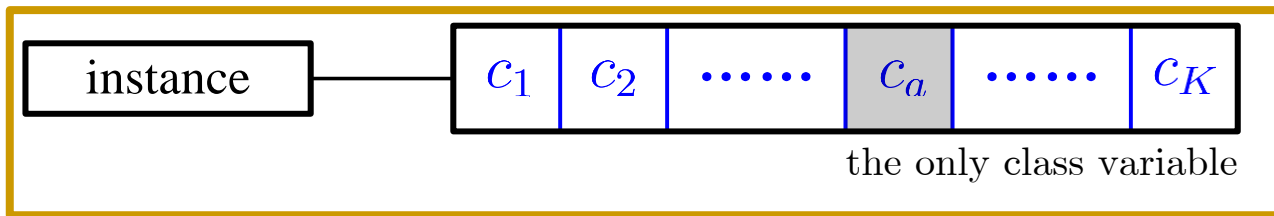
Southeast University, China



Virtual Event

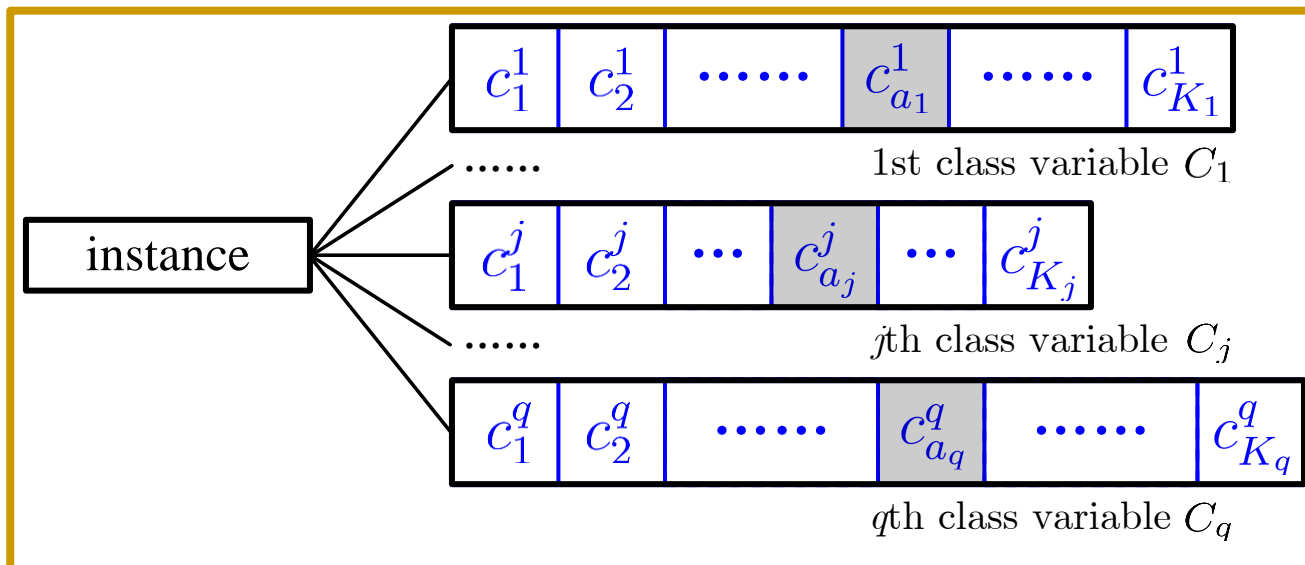
Multi-Dimensional Classification

Traditional Multi-class classification



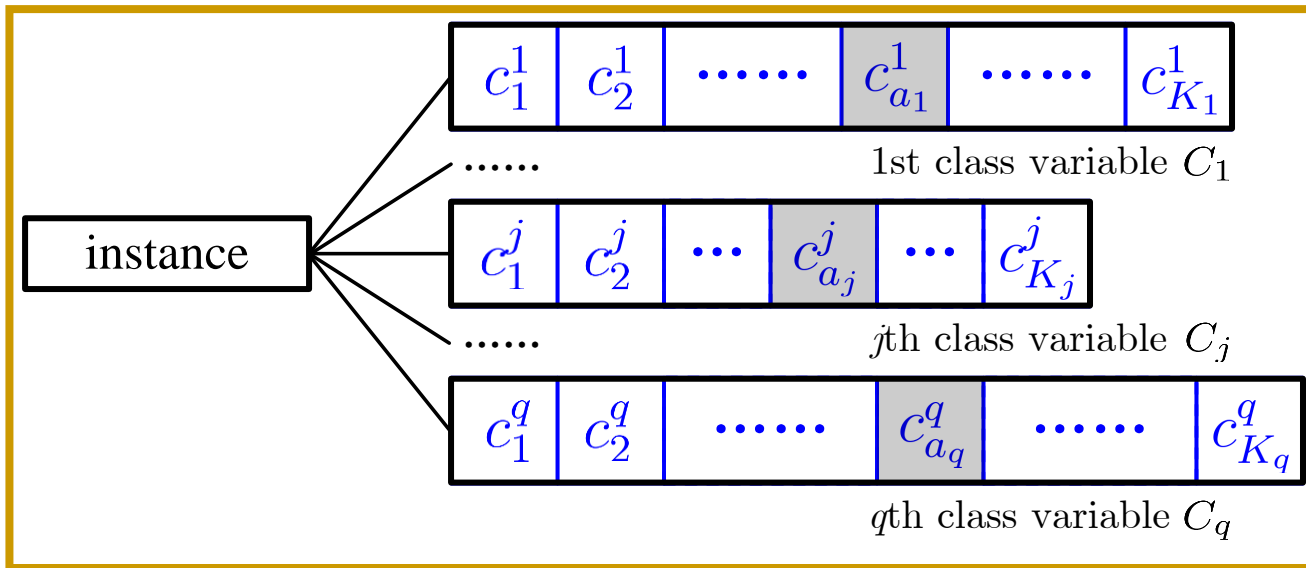
- Only one class variable

Multi-Dimensional Classification (MDC)



- Multiple class variables

The Problem



- Multiple class variables

Existing works: learn the predictive model **in the original output space** for MDC where the dependencies among class variables are considered.

Our work: make a first attempt to learn the predictive model **in its transformed label space** for MDC.

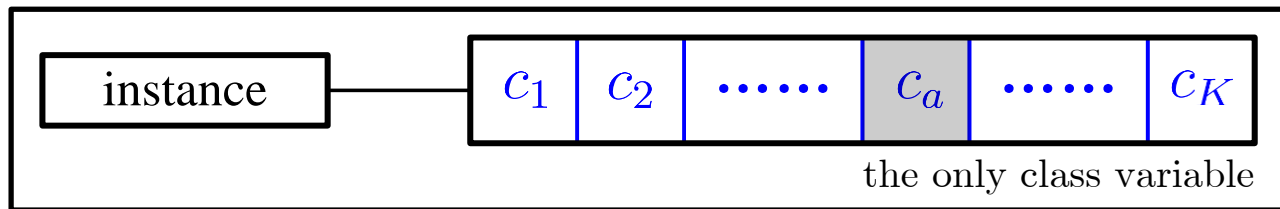
Outline

- Introduction
- The proposed SLEM Approach
- Experiments
 - Experimental setup
 - Experimental results
- Conclusion



Multi-Class Classification (MCC)

object



Input Space

represented by a **single instance** characterizing its properties

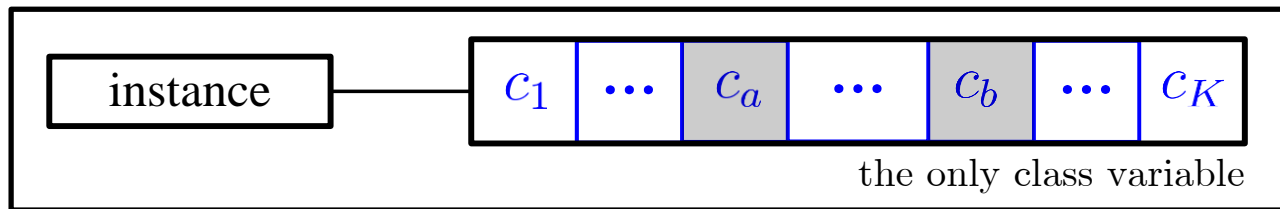
Output Space

associated with a **single class variable** characterizing its semantics

Only one label in the single class space is relevant.

Multi-Label Classification (MLC)

object



Input Space

represented by a single instance characterizing its properties

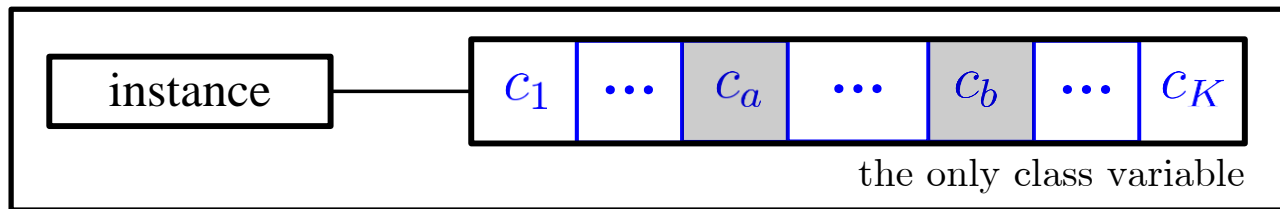
Output Space

associated with a single class variable characterizing its semantics

Multiple labels in the single class space are relevant.

Multi-Label Classification (MLC)

object



Input Space

represented by

Both MCC and MLC assume one single homogeneous class space.

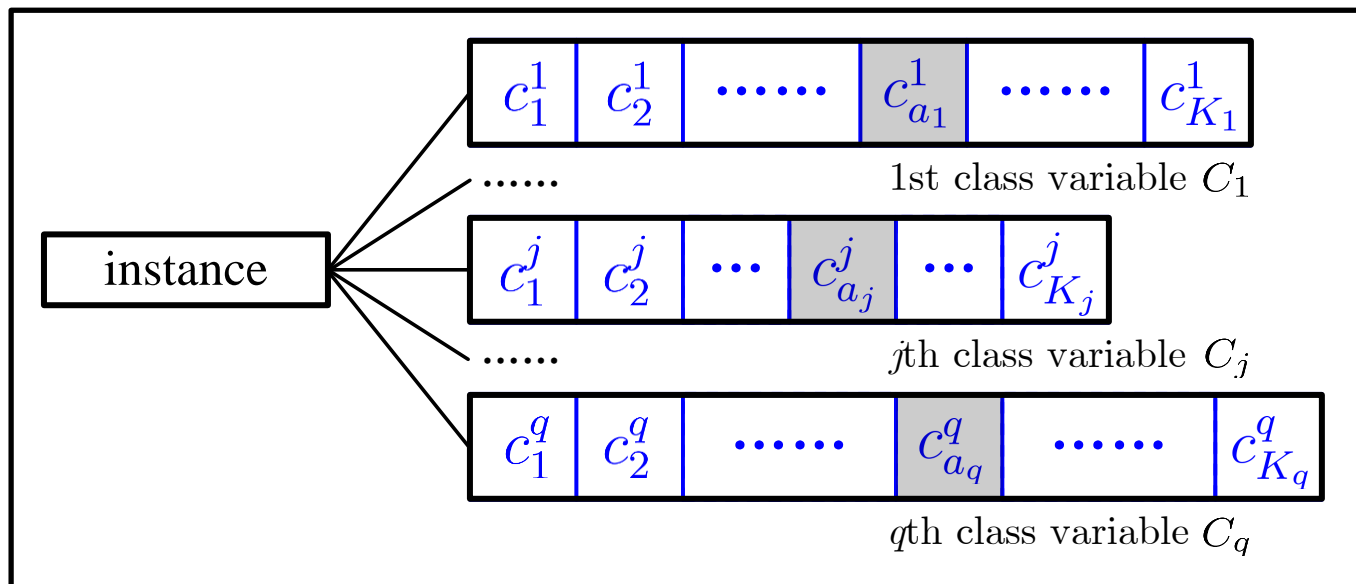
Output

associated with a single class variable characterizing its semantics

Multiple labels in the single class space are relevant.

Multi-Dimensional Classification

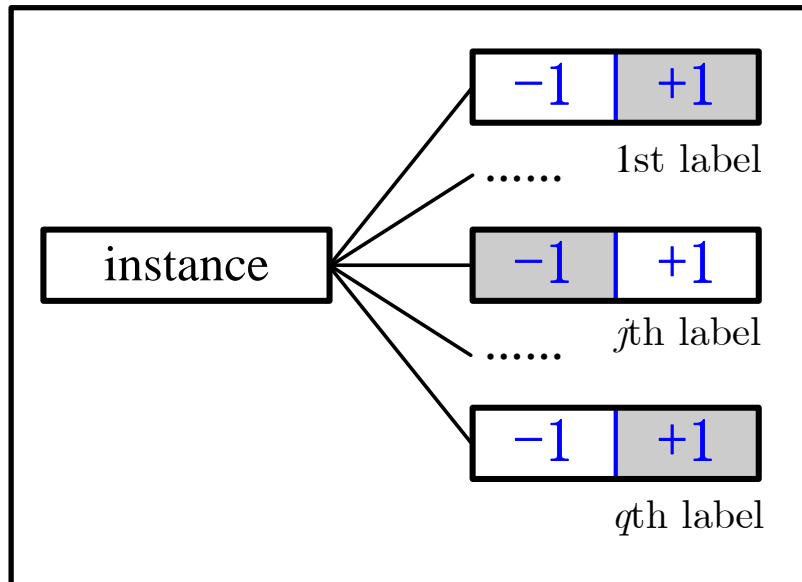
object



Multiple labels are relevant with each of them from one heterogeneous class space.

What's More...

Mathematical view of multi-label classification:



MLC's output space can be regarded as multiple **binary** class variables similar to MDC. But conceptually, they are all in the same class space.

Key difference between MDC and MLC

MDC usually assumes *heterogeneous* semantic spaces

MLC usually assumes *homogeneous* semantic space

MDC Examples



A piece of music



A news document

MDC Examples

Dim. 1: Genre     rock, popular,

Widely exist in real-world applications

- ❑ Text classification [Ortigosa-Hernandez et al., Neurocomputing12]
[Serafino et al., LNAI15] [Tu et al., ACM TIST17]
- ❑ Bioinformatics [Read et al., TKDE14] [Fernandez-Gonzalez et al., ICML
Workshop'15] [Bolt et al., IJAR17][Benjumbeda et al., IJAR18]
- ❑ Computer vision [Ma et al., Neurocomputing18]
- ❑ Resource allocation [Muktadir et al., IEICE TIS19]
- ❑ Other areas [Tekinerdogan, SoSE'19] [Verma et al. Sci Total Environ21]

Dim. 3: Zone

/inter-continental

A news document

Existing Approaches

Intuitive strategies:

- ❑ Binary Relevance (BR): training an independent multi-class classifier w.r.t. each class space
- ❑ Class Powerset (CP): training a single multi-class classifier by conducting powerset transformation

Other specifically designed approaches:

- ❑ Specifying chaining order over class variables [Zaragoza et al., IJCAI'11; Read et al., Pattern Recognition14]
- ❑ Partitioning class variables into groups [Read et al., TKDE14]
- ❑ Assuming DAG structure over class variables [Bolt & van der Gaag, IJAR17; Benjumbeda et al., IJAR18; Gil-Begue et al., Artif. Intell. Rev.21]



Existing Approaches

Intuitive strategies:

- ❑ Binary Relevance (BR): training an independent multi-class classifier w.r.t. each class space
- ❑ One-vs-All (OVA): training a multi-class classifier in the original output space !!!

Other specifically designed approaches:

- ❑ Specifying chaining order over class variables [Liu et al., IJCAI'11; Read et al., Pattern Recognition, 2014]
- ❑ Partitioning class variables into groups [Liu et al., IJCAI'14]
- ❑ Assuming DAG structure over class variables [Bolt & van der Gaag, IJAR17; Benjumbeda et al., IJAR18; Gil-Begue et al., Artif. Intell. Rev.21]

Why not in the transformed output space?

Outline

- Introduction
- **The proposed SLEM Approach**
- Experiments
 - Experimental setup
 - Experimental results
- Conclusion



Formal Definition of MDC

Settings

$\mathcal{X} = \mathbb{R}^d$: d -dimensional input (feature) space

$\mathcal{Y} = C_1 \times C_2 \times \cdots \times C_q$, where $C_j = \{c_1^j, c_2^j, \dots, c_{K_j}^j\}$

: output space which corresponds to the Cartesian product of q class spaces (dim.)

Inputs

$\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i) \mid 1 \leq i \leq m\}$: training data set, where

$$\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{id}]^\top \in \mathcal{X}$$

$$\mathbf{y}_i = [y_{i1}, y_{i2}, \dots, y_{iq}]^\top \in \mathcal{Y}$$

Outputs

f : multi-dimensional classifier $\mathcal{X} \rightarrow \mathcal{Y}$



Formal Definition of MDC

Settings

$\mathcal{X} = \mathbb{R}^d$: d -dimensional input (feature) space

$\mathcal{Y} = C_1 \times C_2 \times \dots \times C_n$ where $C_i = \{c_1^i, c_2^i, \dots, c_{k_i}^i\}$

General MDC approaches

Induce the MDC model *in the original output space*

Our SLEM approach

Induce the MDC model *in the transformed output space*

$$\mathbf{y}_i = [y_{i1}, y_{i2}, \dots, y_{iq}]' \in \mathcal{Y}$$

Outputs

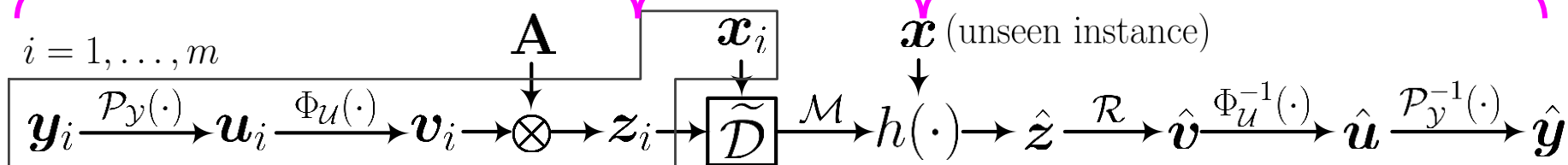
f : multi-dimensional classifier $\mathcal{X} \rightarrow \mathcal{Y}$

The SLEM Approach (1/8)

encoding

training

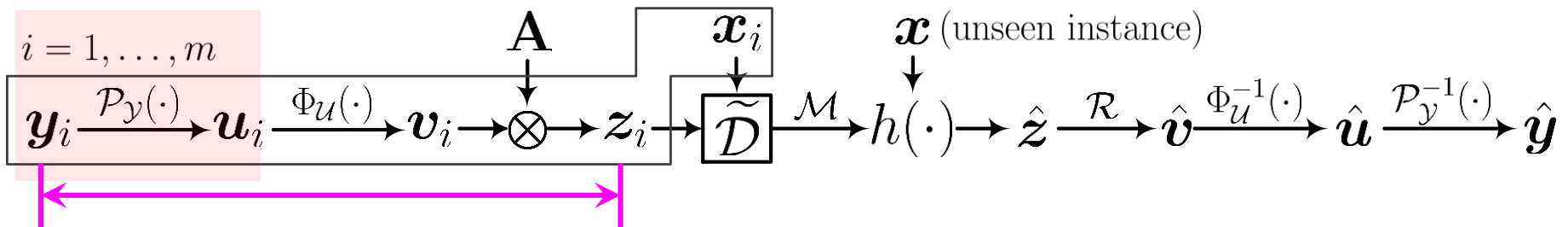
decoding



SLEM works in an *encoding-training-decoding* framework:

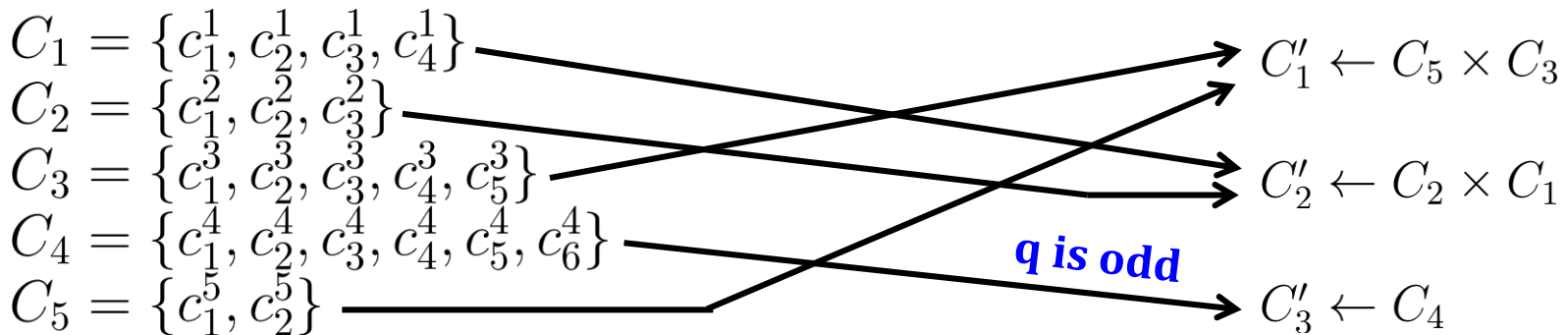
- ❑ In encoding phase, the categorical class vector \mathbf{y}_i is transformed into a real-valued label vector \mathbf{z}_i .
- ❑ In training phase, a multi-output regression model $h(\cdot)$ is induced in the encoded label space.
- ❑ In decoding phase, the predicted class vector $\hat{\mathbf{y}}$ for unseen instance \mathbf{x} is determined by conducting inverse operations in encoding phase based on $h(\mathbf{x})$.

The SLEM Approach (2/8)

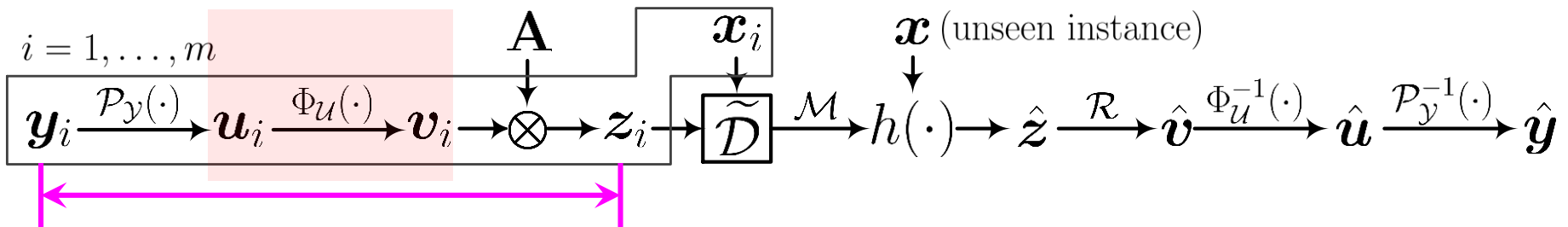


□ Pairwise grouping: $\mathcal{Y} = C_1 \times \dots \times C_q \Rightarrow \mathcal{U} = C'_1 \times \dots \times C'_{\lfloor \frac{q}{2} \rfloor}$

- Make the results of one-hot conversion in next step sparser.
- Group q class spaces into $\lfloor \frac{q}{2} \rfloor$ pairs (plus a singleton one if q is odd) according to the number of class labels in each class space.



The SLEM Approach (3/8)



□ **One-hot conversion:** Categorical space $\mathcal{U} \Rightarrow$ Binary $\{0,1\}$ space

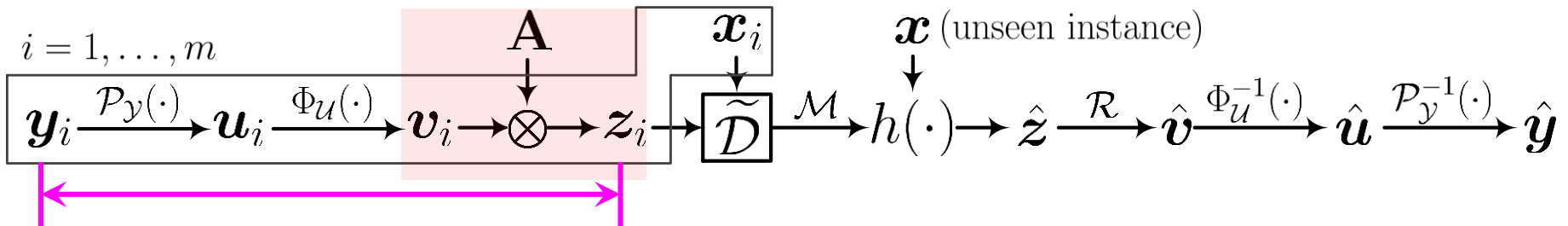
- Facilitate the following numeric computations.
- Convert each class label in categorical class vector u_i into its one-hot form and then concatenate them together.

$$C'_1 = \{a, b, c, d\}, C'_2 = \{\text{I, II, III, IV, V, VI}\}, C'_3 = \{\alpha, \beta, \gamma\} \quad (k = 3)$$

$$u_i = [b, \text{III}, \gamma]^\top \Rightarrow v_i = [0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1]^\top \quad (k\text{-sparse})$$

Local sparsity: there is one and only one '1' among each local group

The SLEM Approach (4/8)



□ Sparse linear encoding: Binary $\{0,1\}$ space \Rightarrow Real-valued space

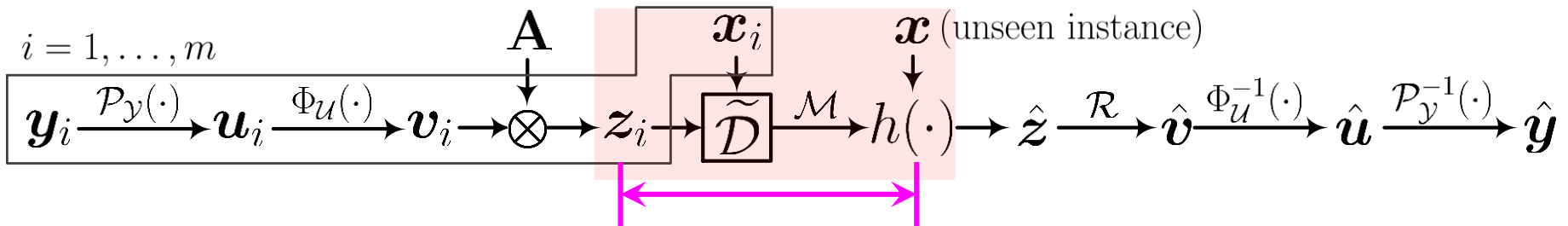
- Blend the heterogeneous class spaces into an integrated one.
- Encode binary label vector \mathbf{v}_i into real-valued vector \mathbf{z}_i with matrix \mathbf{A} which satisfies k -RIP (i.e., $\mathbf{z}_i = \mathbf{A}\mathbf{v}_i$):

Definition 1. For matrix \mathbf{A} , if there is a constant $\delta_k \in [0, 1)$ which satisfies

$$(1 - \delta_k) \|\mathbf{v}\|_2^2 \leq \|\mathbf{A}\mathbf{v}\|_2^2 \leq (1 + \delta_k) \|\mathbf{v}\|_2^2$$

where \mathbf{v} is any k -sparse vector, then \mathbf{A} is known as satisfying k -order Restricted Isometry Property (k -RIP) (e.g., Gaussian matrix and Bernoulli matrix).

The SLEM Approach (5/8)



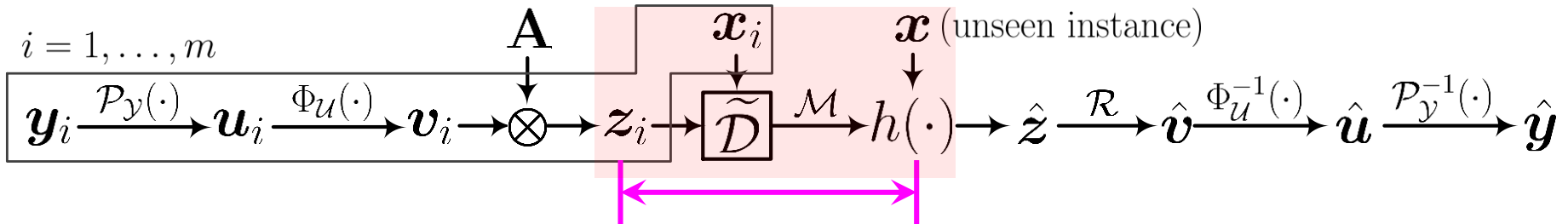
□ **Training:** Learn multi-output regression model $h(\cdot)$ with algorithm \mathcal{M}

- Train predictive model over $\tilde{\mathcal{D}} = \{(\mathbf{x}_i, \mathbf{z}_i) \mid 1 \leq i \leq m\}$.
- Learn multi-output regression model $h(\mathbf{x}) = \mathbf{W}^\top \mathbf{x} + \mathbf{b}$ via optimizing the following formulation:

$$\min_{\mathbf{W}, \mathbf{b}, \hat{\mathbf{V}}} \frac{1}{2} \|\mathbf{W}\|_F^2 + \lambda \sum_{i=1}^m [\|h(\mathbf{x}_i) - \mathbf{z}_i\|_2^2 + \gamma_1 (\|h(\mathbf{x}_i) - \mathbf{A}\hat{\mathbf{v}}_i\|_2^2 + \gamma_2 \|\hat{\mathbf{v}}_i - \mathbf{v}_i\|_1)]$$

Here, $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{s'}] \in \mathbb{R}^{d \times s'}$ and $\mathbf{b} = [b_1, b_2, \dots, b_{s'}]^\top$ are the model parameters of h to be determined, $\hat{\mathbf{V}} = [\hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_m]^\top \in \mathbb{R}^{m \times s}$ with $\hat{\mathbf{v}}_i$ corresponding to the recovered sparse vector for \mathbf{v}_i based on its prediction $h(\mathbf{x}_i)$, and λ , γ_1 and γ_2 are three trade-off parameters.

The SLEM Approach (5/8)



□ **Training:** Learn multi-output regression model $h(\cdot)$ with algorithm \mathcal{M}

- Train predictive model over $\tilde{\mathcal{D}} = \{(\mathbf{x}_i, \mathbf{z}_i) \mid 1 \leq i \leq m\}$.

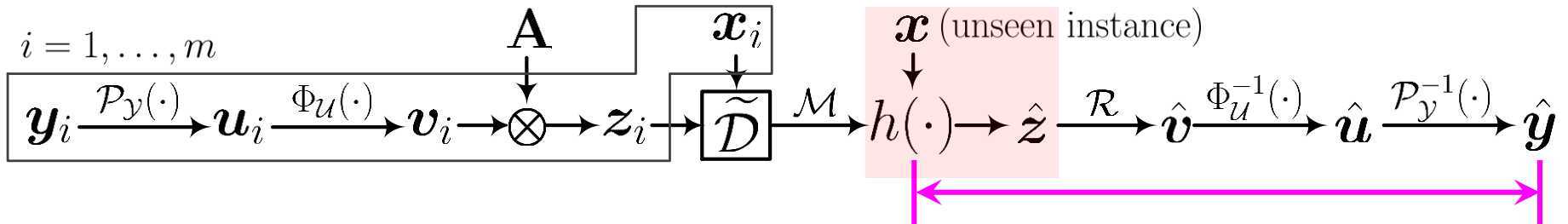
- **regularizer** by **Square loss** on **Facilitate sparse reconstruction**

minimize the following formulation:

$$\min_{\mathbf{W}, \mathbf{b}, \hat{\mathbf{V}}} \frac{1}{2} \|\mathbf{W}\|_F^2 + \lambda \sum_{i=1}^m \left[\|h(\mathbf{x}_i) - \mathbf{z}_i\|_2^2 + \gamma_1 (\|h(\mathbf{x}_i) - \mathbf{A}\hat{\mathbf{v}}_i\|_2^2 + \gamma_2 \|\hat{\mathbf{v}}_i - \mathbf{v}_i\|_1) \right]$$

Here, $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{s'}] \in \mathbb{R}^{d \times s'}$ and $\mathbf{b} = [b_1, b_2, \dots, b_{s'}]^\top$ are the model parameters of h to be determined, $\hat{\mathbf{V}} = [\hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_m]^\top \in \mathbb{R}^{m \times s}$ with $\hat{\mathbf{v}}_i$ corresponding to the recovered sparse vector for \mathbf{v}_i based on its prediction $h(\mathbf{x}_i)$, and λ , γ_1 and γ_2 are three trade-off parameters.

The SLEM Approach (6/8)



□ **Testing:** Obtain the real-valued prediction $\hat{\mathbf{z}} = h(\mathbf{x})$ for unseen instance \mathbf{x}

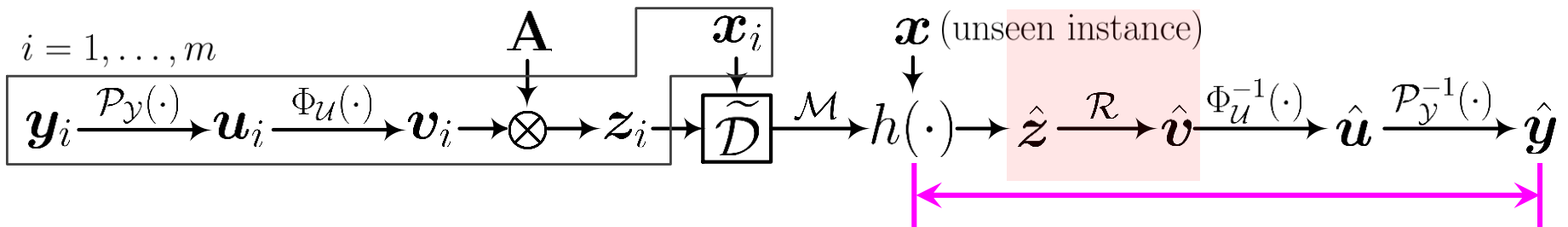
Suppose that \mathbf{y} is the ground-truth class vector of \mathbf{x} , then $\hat{\mathbf{z}} = h(\mathbf{x})$ should correspond to the prediction of \mathbf{z} :

$$\mathbf{y} \xrightarrow{\mathcal{P}_y(\cdot)} \mathbf{u} \xrightarrow{\Phi_u(\cdot)} \mathbf{v} \xrightarrow{\mathbf{A}} \mathbf{z}$$

The decoding phase corresponds to the inverse operations of encoding phase to obtain \mathbf{x} 's predicted class vector $\hat{\mathbf{y}}$ based on $\hat{\mathbf{z}}$.

Note that \mathbf{v} is k -sparse, and there is one and only one '1' in each local group (i.e., local sparsity).

The SLEM Approach (7/8)



□ Inverse of sparse label encoding

- Complete sparse reconstruction and keep local sparsity.
- Adapt the well-known orthogonal matching pursuit (OMP) algorithm to consider the local sparsity property.

The OMP algorithm recovers a k -sparse vector by choosing the column of \mathbf{A} that is most strongly correlated with the residual at each iteration and repeating this procedure k times.

Our key adaptation is to set the related columns to zero which belong to the same local group with the current selected column of \mathbf{A} .

The S

Algorithm 2 LOMP: $v = \mathcal{R}(z, \mathbf{A}, k, \mathcal{I})$

Input: The encoding matrix $\mathbf{A} \in \mathbb{R}^{s' \times s}$, the real-valued vector $z \in \mathbb{R}^{s'}$, sparsity level k , local sparsity information $\mathcal{I} : s_1, s_2, \dots, s_k$;

Output: The recovered k -sparse vector v ;

1: Initialize v as zero vector with length $s = \sum_{j=1}^k s_j$;

2: Initialize $r_0 = z, J = \emptyset, \mathbf{B} = \mathbf{A}$;

3: **for** $i = 1$ to k **do**

4: $j_* = \arg \max_j |\langle r_{i-1}, \mathbf{B}_{:j} \rangle|$;

5: $v(j_*) = 1$;

6: $J = J \cup \{j_*\}$;

7: $r_i = z - \mathbf{A}_{:J}(\mathbf{A}_{:J}^\top \mathbf{A}_{:J})^{-1} \mathbf{A}_{:J}^\top z$;

8: **for** $\kappa = 1$ to k **do**

9: $t_f = \sum_{t=1}^{\kappa} s_t$;

10: **if** $j_* \leq t_f$ **then**

11: $t_b = t_f - s_{\kappa}$;

12: $T = \{t_b + 1, t_b + 2, \dots, t_f\}$;

13: $\mathbf{B}_{:T} = 0$;

14: **break**;

15: **end if**

16: **end for**

17: **end for**

18: **Return** v .

$i = 1, \dots, m$

$$y_i \xrightarrow{\mathcal{P}_y(\cdot)} u_i$$

$$\hat{u} \xrightarrow{\mathcal{P}_y^{-1}(\cdot)} \hat{y}$$

□ Inverse

- Complet

- Adapt the algorithm to c

The OMP column of each iterat

Our key ad to the sam

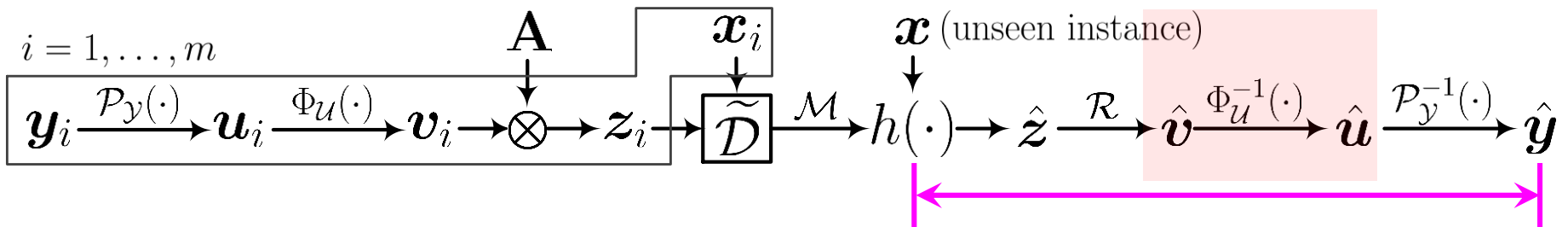


P) algo-

using the residual at

h belong of \mathbf{A} .

The SLEM Approach (8/8)

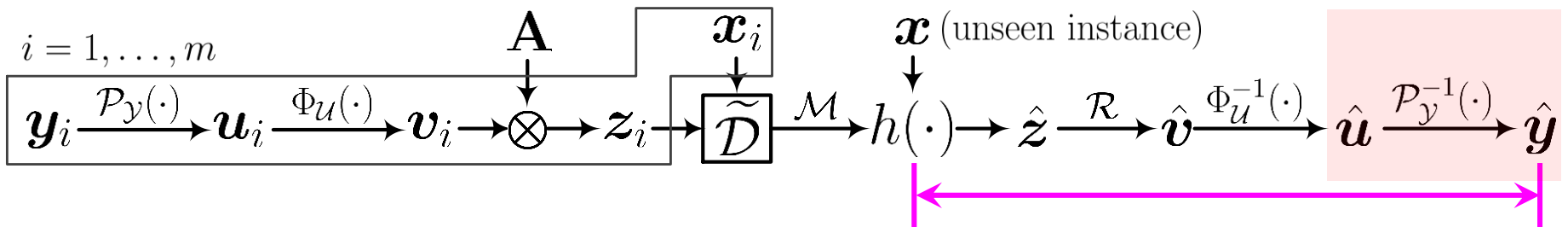


□ Inverse of one-hot conversion

$$C'_1 = \{a, b, c, d\}, C'_2 = \{\text{I, II, III, IV, V, VI}\}, C'_3 = \{\alpha, \beta, \gamma\} \quad (k = 3)$$

$$v_i = [0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1]^\top \Rightarrow u_i = [b, \text{III}, \gamma]^\top \quad (k\text{-sparse})$$

The SLEM Approach (8/8)



□ Inverse of one-hot conversion

$$C'_1 = \{a, b, c, d\}, C'_2 = \{I, II, III, IV, V, VI\}, C'_3 = \{\alpha, \beta, \gamma\} \quad (k = 3)$$

$$v_i = [0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1]^T \Rightarrow u_i = [b, III, \gamma]^T \quad (k\text{-sparse})$$

□ Inverse of pairwise grouping

$$\begin{aligned}
 C'_1 \leftarrow C_5 \times C_3 & \rightarrow C_1 = \{c_1^1, c_2^1, c_3^1, c_4^1\} \\
 C'_2 \leftarrow C_2 \times C_1 & \rightarrow C_2 = \{c_1^2, c_2^2, c_3^2\} \\
 C'_3 \leftarrow C_4 & \rightarrow C_3 = \{c_1^3, c_2^3, c_3^3, c_4^3, c_5^3\} \\
 & \rightarrow C_4 = \{c_1^4, c_2^4, c_3^4, c_4^4, c_5^4, c_6^4\} \\
 & \rightarrow C_5 = \{c_1^5, c_2^5\}
 \end{aligned}$$

Outline

- Introduction
- The proposed SLEM Approach
- **Experiments**
 - **Experimental setup**
 - **Experimental results**
- Conclusion



Experimental Setup

Experimental data sets

Characteristics of the experimental MDC data sets.

Data Set	#Exam.	#Dim.	#Labels/Dim.	#Features [†]
Jura	359	2	4,5	9 n
Oes10	403	16	3	298 n
Voice	3136	2	4,2	19 n
Scm20d	8966	16	4	61 n
Rf1	8987	8	4,4,3,4,4,3,4,3	64 n
Scm1d	9803	16	4	280 n
CoIL2000	9822	5	6,10,10,4,2	81 x
Flickr	12198	5	3,4,3,4,4	1536 n
Disfa	13095	12	5,5,6,3,4,4,5,4,4,4,6,4	136 n
Fera	14052	5	6	136 n
Adult	18419	4	7,7,5,2	5 n ,5 x

[†] n , x denote numeric and nominal features respectively.



Experimental Setup

Evaluation Metrics

Testing set: $\mathcal{S} = \{(\mathbf{x}_i, \mathbf{y}_i) \mid 1 \leq i \leq p\}$, where $\mathbf{y}_i = [y_{i1}, y_{i2}, \dots, y_{iq}]^\top$

Predicted class vector: $\hat{\mathbf{y}}_i = f(\mathbf{x}_i) = [\hat{y}_{i1}, \hat{y}_{i2}, \dots, \hat{y}_{iq}]^\top$

For each MDC test example $(\mathbf{x}_i, \mathbf{y}_i) : r^{(i)} = \sum_{j=1}^q \llbracket y_{ij} = \hat{y}_{ij} \rrbracket$

Hamming Score:
$$\text{HS}_{\mathcal{S}}(f) = \frac{1}{p} \sum_{i=1}^p \frac{1}{q} \cdot r^{(i)}$$

Exact Match:
$$\text{EM}_{\mathcal{S}}(f) = \frac{1}{p} \sum_{i=1}^p \llbracket r^{(i)} = q \rrbracket$$

Sub-Exact Match:
$$\text{SEM}_{\mathcal{S}}(f) = \frac{1}{p} \sum_{i=1}^p \llbracket r^{(i)} \geq q - 1 \rrbracket$$

Experimental Setup

Compared Algorithms

BR: Learn q independent multi-class classifier, one per dimension

CP: Learn a single multi-class classifier via powerset transformation

BCC: Learn q chain-structured multi-class classifiers, one per dimension

ESC: Group the class variables into groups

gMML: Learn a regressor for each class label as well as a Mahalanobis distance metric to train all regressor in a joint manner

Experimental Protocol

Ten-fold cross-validation + Pairwise t -test

Experimental Results

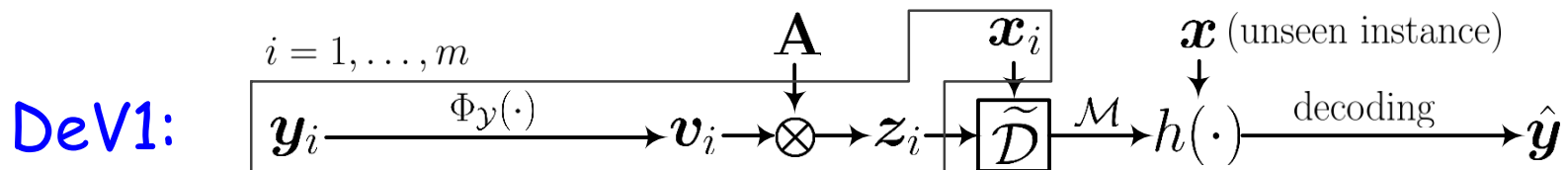
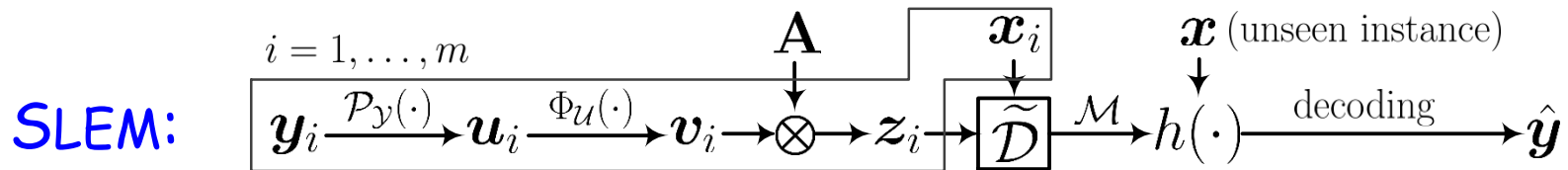
Win/tie/loss counts of pairwise t -test (at 0.05 significance level) between SLEM and each MDC approach.

Evaluation metric	SLEM against				
	BR	CP	BCC	ESC	gMML
HS	9/1/1	7/0/0	10/1/0	8/0/0	8/0/3
EM	10/1/0	5/2/0	10/1/0	7/1/0	9/1/1
SEM	7/2/2	5/1/1	8/2/1	6/1/1	5/4/2
In Total	26/4/3	17/3/1	28/4/1	21/2/1	22/5/6

Among all the 144 configurations, SLEM achieves superior or at least comparable performance against the five compared approaches in 132 cases.



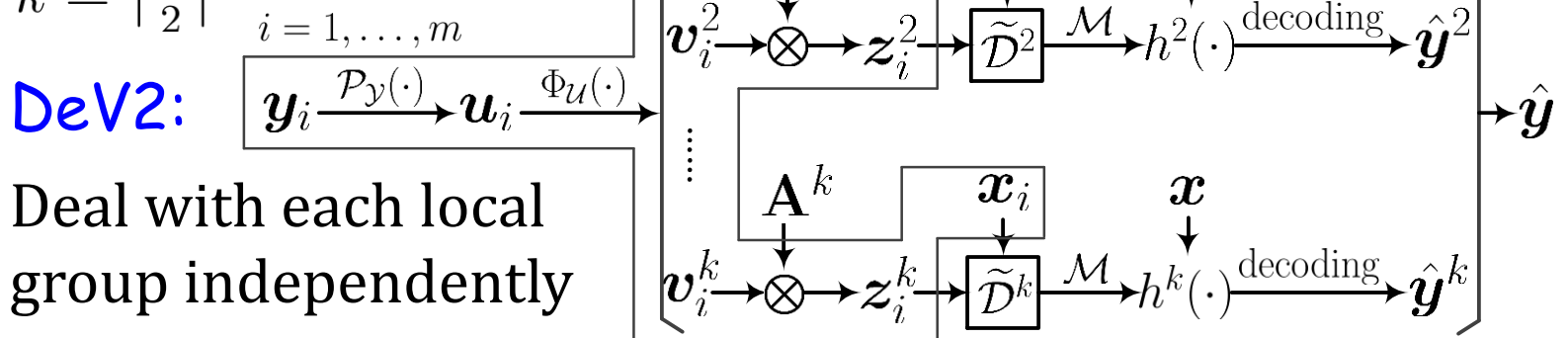
Further Analysis (1/2)



Omit pairwise grouping

$$\mathbf{v}_i = [\mathbf{v}_i^1, \mathbf{v}_i^2, \dots, \mathbf{v}_i^k]$$

$$k = \lceil \frac{q}{2} \rceil$$



Deal with each local group independently

Further Analysis (2/2)

Wilcoxon signed-ranks test for SLEM against its two degenerated versions in terms of each evaluation metric (significance level $\alpha = 0.05$; p -values shown in the brackets).

SLEM versus	Evaluation metric		
	HS	EM	SEM
DeV1	win [9.77e-04]	win [9.77e-04]	win [9.77e-04]
DeV2	win [3.91e-03]	win [3.91e-03]	win [3.91e-03]

- The superiority of SLEM against DeV1 shows the benefits of the pairwise grouping operation.
- The superiority of SLEM against DeV2 shows the benefits of encoding the heterogeneous class spaces into an integrated one.



Outline

- Introduction
- The proposed SLEM Approach
- Experiments
 - Experimental setup
 - Experimental results
- **Conclusion**



Conclusion

- Different from most existing MDC approaches, we propose a first attempt towards learning predictive models in the transformed label space instead of the original one.
- We design a novel MDC approach named SLEM which works in an *encoding-training-decoding* framework by utilizing the sparse property of the transformed label space.
- Experimental results clearly validate the superiority of SLEM against state-of-the-art MDC approaches.



Thanks !

<http://palm.seu.edu.cn/zhangml/files/SLEM.rar>

Email: jiabb@seu.edu.cn

