# FEATURE-AWARE LABEL SPACE DIMENSION REDUCTION FOR MULTI-LABEL CLASSIFICATION Yao-Nan Chen (r99922008@csie.ntu.edu.tw) and Hsuan-Tien Lin (htlin@csie.ntu.edu.tw) Department of Computer Science and Information Engineering, National Taiwan University

### MULTI-LABEL CLASSIFICATION SETUP

Which tags  $\mathcal{Y}$  are associated with this picture  $\mathbf{x}$ ?



[ building, taipei 101, <del>day view</del>,  $\mathcal{Y} = \{$ night view, skyscraper, fireworks, <del>new</del> york, fireworks, car, face, taipei world financial center, university, etc.}

(CC BY-SA SElefant from Wikimedia Commons)

- Given: N examples  $\left\{ \left( \mathbf{x}_n \in \mathbb{R}^d, \mathcal{Y}_n \subseteq \{1, 2, \cdots, K\} \right) \right\}_{n=1}^N$
- Goal: classifier  $g(\mathbf{x})$  that closely predicts the label-set  $\mathcal{Y}$  associated with some unseen inputs x, presumably by exploiting hidden relations between labels, e.g.
  - taipei 101 & taipei world financial center highly correlated
  - skyscraper subset of building
  - day view & night view disjoint

### LABEL SPACE DIMENSION REDUCTION

 $\mathcal{Y} \subseteq \{1, 2, \cdots, K\}$  equivalent to  $\mathbf{y} \in \{0, 1\}^K$ 

- feature space dimension reduction: compress  $\mathbf{x}$  to remove irrelevant, redundant (possibly related), or noisy information, and achieve better efficiency & performance
  - principal component analysis (PCA): linearly project  $\mathbf{x}$  to  $\mathbf{w}_m^T \mathbf{x}$  with minimum projection error
  - canonical correlation analysis (CCA): linearly project  $\mathbf{x}$  to  $\mathbf{w}_m^T \mathbf{x}$  in order to maximize correlation with some  $\mathbf{v}_m^T \mathbf{y}$
- label space dimension reduction: analogously, but compress y instead

1. compress: transform  $\{(\mathbf{x}_n, \mathbf{y}_n)\}$  to  $\{(\mathbf{x}_n, \mathbf{t}_n)\}$  with  $\mathbf{t}_n = \operatorname{compress}(\mathbf{y}_n) \in \mathbb{R}^M \text{ and } M \ll K$ 

- 2. learn: train some  $\mathbf{r}(\mathbf{x})$  from  $\{(\mathbf{x}_n, \mathbf{t}_n)\}$
- 3. decompress:  $g(\mathbf{x}) = \operatorname{decompress}(\mathbf{r}(\mathbf{x}))$
- compressive sensing (Hsu et al., NIPS 2009): linearly project y to  $\mathbf{t}[m] = \mathbf{v}_m^T \mathbf{y}$ with random  $\mathbf{v}_m$ 's (for incoherence)
- principal label space transformation (PLST; Tai and Lin, NC 2012): linearly project y to  $\mathbf{t}[m] = \mathbf{v}_m^T \mathbf{y}$  with minimum projection error (sibling of PCA)

# FEATURE-AWARE LABEL SPACE DIMENSION REDUCTION

• feature space dimension reduction

unsupervised (not using  $\mathbf{y}$ )

PCA, locally linear embedding, etc.

• label space dimension reduction

feature-unaware (not using  $\mathbf{x}$ )

PLST, compressive sensing, etc.

—can we improve PLST by **feature-aware** label space dimension reduction?

# CONDITIONAL PRINCIPAL LABEL SPACE TRANSFORMATION

• idea 1: exploit dual role of CCA to be feature-aware

proposed OCCA :  $\min_{\mathbf{W},\mathbf{V}} \| \mathbf{X} \|$ 

- project to easiest-by-linear-regression directions

• idea 2: keep benefits of PLST for compression existing  $PLST : \min || \mathbf{Y}$ 

- project to most representative directions

 $\min_{\mathbf{W},\mathbf{V}} \|\mathbf{X}\mathbf{W}^T - \mathbf{Y}\mathbf{V}^T\|_F^2 + \|$ learning error compression error

- theoretical guarantee (Tai and Lin, NC 2012): when using linear regression as  $\mathbf{r}$ ,

hamming loss  $\leq$  learning error + compression error

- algorithmic simplicity: closed-form optimal V contains top eigenvectors of

| OCCA  | PLST                      | CPLST  |
|---|---------------------------|--|
| $\mathbf{Y}^T (\underbrace{\mathbf{X} \mathbf{X}^{\dagger}}_{\text{hot matrix}} - \mathbf{I}) \mathbf{Y}$ | $\mathbf{Y}^T \mathbf{Y}$ | $\begin{array}{ c c c c c c c c c c c c c c c c c c c$ |
| Hat HIATIX  |                           | hat matrix   |

(note:  $\mathbf{Z}$ , i.e. the mean-shifted  $\mathbf{Y}$ , is actually used for better projection)

supervised (using  $\mathbf{y}$ )

CCA, sliced inverse regression, etc.

—supervised generally better for learning from compress(x) to y

feature-aware (using  $\mathbf{x}$ )

???

project  $\mathbf{x}$  to  $\mathbf{w}_m^T \mathbf{x}$  in order to maximize correlation with some  $\mathbf{v}_m^T \mathbf{y}$  $\equiv$  project y to  $\mathbf{v}_m^T \mathbf{y}$  in order to maximize correlation with some  $\mathbf{w}_m^T \mathbf{x}$  $\approx$  project y to  $\mathbf{v}_m^T \mathbf{y}$  in order to minimize difference to some  $\mathbf{w}_m^T \mathbf{x}$ 

$$\mathbf{W}^T - \mathbf{W}^T ||_F^2$$
, s.t.  $\mathbf{V}\mathbf{V}^T = \mathbf{I}$ 

$$\| \mathbf{Y} - \mathbf{Y}\mathbf{V}^T\mathbf{V} \|_F^2$$
, s.t.  $\mathbf{V}\mathbf{V}^T = \mathbf{I}$ 

### • proposed algorithm: conditional principal label space transformation (CPLST)

$$\|\mathbf{Y} - \mathbf{Y}\mathbf{V}^T\mathbf{V}\|_{F_{\star}}^2$$
, s.t.  $\mathbf{V}\mathbf{V}^T = \mathbf{I}$ 

# – physical meaning: exploit conditional (feature-aware) correlations - kernelization: replace linear regression with kernel ridge regression as r

# EXPERIMENTAL RESULTS

### on yeast data set:



• CPLST: optimize learning+compression error, and hence **best hamming loss** on 8 benchmark data sets:

| algorithms | CPLST vs. PLST           | CPLST vs. PLST           | kernel-CPLST vs. PLST            |
|------------|--------------------------|--------------------------|----------------------------------|
|            | + linear regression      | + decision tree          | + kernel ridge regression        |
| M = 20% K  | <b>3 win</b> , 5 similar | <b>2 win</b> , 6 similar | <b>5 win</b> , 1 lose, 2 similar |

CPLST consistently better than or similar to PLST across data & algorithms

## SUMMARY

Conditional Principal Label Space Transformation, which

- readily-strong PLST



• PBR: baseline, with standard basis as  $\mathbf{v}_m$ 

• OCCA: optimize learning error, but worst in compression error

• PLST: optimize compression error, but worst in learning error

• projects to **conditional** principal directions by combining ideas behind **CCA** (featureaware) and PLST (optimal compression)

• can be **kernelized** for exploiting feature power

• achieves **better/similar** practical performance **consistently** when compared with the