

# Implicit Learning of Simpler Output Kernels for Multi-label Prediction

Hanchen Xiong, Sandor Szedmak, Justus Piater  
UNIVERSITY OF INNSBRUCK, AUSTRIA



## OVERVIEW

### Research in Multi-label Prediction:

- exploiting and utilizing inter-label dependencies; 😊
- increasingly more sophisticated dependencies are used; ?
- “overfit” output structural dependencies when the desired ones are simpler 😞
- a regularization should be added on output structural dependencies ✓

### Our contributions: Joint SVM

- joint SVM  $\iff$  structural SVM with linearly decomposable score functions;
- in joint SVM, (non-)linear kernels can be assigned on outputs to capture inter-label dependencies;
- joint SVM shares the same computational complexity as a single SVM;
- when linear output kernels are used, a output-kernel regularization is implicitly added.
- yield promising results on 3 image annotation databases.

## FROM STRUCTURAL SVM TO JOINT SVM

Structural SVM is an extension of SVM for structured-outputs, in which, however, the margin to be maximized is defined as the score gap between the desired output and the first runner-up.

$$\arg \min_{\mathbf{W} \in \mathbb{R}^{\Psi}} \frac{1}{2} \|\mathbf{W}\|^2 + C \sum_{i=1}^m \max_{\mathbf{y}' \in \mathcal{Y}} \left\{ d(\mathbf{y}^{(i)}, \mathbf{y}') - \Delta_F(\mathbf{y}^{(i)}, \mathbf{y}') \right\} \quad (1)$$

where  $\Delta_F(\mathbf{y}^{(i)}, \mathbf{y}') = F(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}; \mathbf{W}) - F(\mathbf{x}^{(i)}, \mathbf{y}'; \mathbf{W})$ , the score function is  $F(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}; \mathbf{W}) = \langle \mathbf{W}, \phi(\mathbf{x}^{(i)}) \otimes \mathbf{y}^{(i)} \rangle$ , and Hamming distance  $d(\mathbf{y}^{(i)}, \mathbf{y}')$  is used on outputs. Because of linear decomposability, (1) can be rewritten as:

$$\begin{aligned} & \arg \min_{\mathbf{W} \in \mathbb{R}^{\mathcal{H}_\phi \times \mathbb{R}^T}} \frac{1}{2} \|\mathbf{W}\|_F^2 + C \sum_{i=1}^m \sum_{t=1}^T \max_{y'_t \in \{-1, +1\}} \left\{ d(y_t^{(i)}, y'_t) - \Delta_F(y_t^{(i)}, y'_t) \right\} \\ & \quad \downarrow \\ & \arg \min_{\mathbf{w}_1, \dots, \mathbf{w}_T \in \mathbb{R}^{\mathcal{H}_\phi}} \sum_{t=1}^T \left\{ \frac{1}{2} \|\mathbf{w}_t\|^2 + C \sum_{i=1}^m \max \left\{ 0, d(y_t^{(i)}, -y_t^{(i)}) - \Delta_F(y_t^{(i)}, -y_t^{(i)}) \right\} \right\} \quad (2) \\ & \quad \downarrow \\ & \arg \min_{\mathbf{w}_1, \dots, \mathbf{w}_T \in \mathbb{R}^{\mathcal{H}_\phi}} \sum_{t=1}^T \left\{ \frac{1}{2} \|\mathbf{w}_t\|^2 + 2C \sum_{i=1}^m \max \left\{ 0, 1 - y_t^{(i)} \mathbf{w}_t^\top \phi(\mathbf{x}^{(i)}) \right\} \right\} \end{aligned}$$

where  $\langle \cdot, \cdot \rangle_F$  denotes Frobenius product and  $\|\mathbf{W}\|_F$  is the Frobenius norm of matrix  $\mathbf{W}$ .

It can be seen in (2) that, with linearly decomposable score functions and output distances, structural SVM on multi-label learning is equivalent to learning  $T$  SVMs jointly: **Joint SVM**

$$\begin{aligned} & \arg \min_{\mathbf{W} \in \mathbb{R}^{\mathcal{H}_\psi \times \mathcal{H}_\phi}} \frac{1}{2} \|\mathbf{W}\|_F^2 + C \sum_{i=1}^m \bar{\xi}^{(i)} \\ & \quad \text{s.t.} \quad \langle \psi(\mathbf{y}^{(i)}), \mathbf{W} \phi(\mathbf{x}^{(i)}) \rangle \geq 1 - \bar{\xi}^{(i)}, \bar{\xi}^{(i)} \geq 0, i \in \{1, \dots, m\} \end{aligned} \quad (3)$$

where  $\mathbf{y}^{(i)} = [y_1^{(i)}, \dots, y_T^{(i)}]$ ,  $\mathbf{W} = [\frac{\mathbf{w}_1^\top}{T}; \dots; \frac{\mathbf{w}_T^\top}{T}]^\top$ . The corresponding dual form of Joint SVM is:

$$\begin{aligned} & \arg \min_{\alpha_1, \dots, \alpha_m} \sum_{i=1}^m \alpha_i - \sum_{i,j=1}^m \alpha_i \alpha_j K_\psi(\mathbf{y}^{(i)}, \mathbf{y}^{(j)}) K_\phi(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \\ & \quad \text{s.t.} \quad \forall i, 0 \leq \alpha_i \leq C \end{aligned} \quad (4)$$

therefore, the same computational complexity as a single regular SVM.

## IMPLICIT REGULARIZATION ON LINEAR OUTPUT KERNELS

When a linear output kernel is used to capture pairwise dependencies of labels, via which the output vectors can be linearly mapped as  $\psi(\mathbf{y}) = \mathbf{P}\mathbf{y}$ :  $K_\psi^{Lin}(\mathbf{y}^{(i)}, \mathbf{y}^{(j)}) = \mathbf{y}^{(i)\top} \mathbf{\Omega} \mathbf{y}^{(j)}$  where  $\mathbf{\Omega} = \mathbf{P}^\top \mathbf{P}$ . By denoting  $\mathbf{U} = \mathbf{P}^\top \mathbf{W}$ , we can have:

$$\begin{aligned} & \arg \min_{\mathbf{W} \in \mathbb{R}^{\mathcal{H}_\psi \times \mathcal{H}_\phi}} \frac{1}{2} \|\mathbf{W}\|_F^2 + C \sum_{i=1}^m \bar{\xi}^{(i)} \\ & \quad \text{s.t.} \quad \langle \mathbf{y}^{(i)}, \mathbf{U} \phi(\mathbf{x}^{(i)}) \rangle \geq 1 - \bar{\xi}^{(i)}, \bar{\xi}^{(i)} \geq 0, i \in \{1, \dots, m\} \end{aligned} \quad (5)$$

we use a compact regularization for both  $\mathbf{W}$  and  $\mathbf{\Omega}$ ,  $\frac{1}{2} \|\mathbf{W}^\top \mathbf{\Omega}^\top \mathbf{W}\|_F^2$ , resulting in:

$$\begin{aligned} & \arg \min_{\mathbf{U} \in \mathbb{R}^{\mathcal{H}_\psi \times \mathcal{H}_\phi}} \frac{1}{2} \|\mathbf{U}\|_F^2 + C \sum_{i=1}^m \bar{\xi}^{(i)} \\ & \quad \text{s.t.} \quad \langle \mathbf{y}^{(i)}, \mathbf{U} \phi(\mathbf{x}^{(i)}) \rangle \geq 1 - \bar{\xi}^{(i)}, \bar{\xi}^{(i)} \geq 0, i \in \{1, \dots, m\} \end{aligned} \quad (6)$$

Remarkably, a linear output kernel is implicitly learned, and absorbed in  $\mathbf{W}$ , also a regularization on the output kernel is also implicitly added.

## EXPERIMENTAL RESULTS AND COMPARISON

| Method      | Core5k      |             |             | Espgame     |             |             | iaprtc12    |             |             |
|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|             | P(%)        | R(%)        | F1(%)       | P(%)        | R(%)        | F1(%)       | P(%)        | R(%)        | F1(%)       |
| MBRM        | 24.0        | 25.0        | 24.0        | 18.0        | 19.0        | 18.0        | 24.0        | 23.0        | 23.0        |
| JEC         | 27.0        | 32.0        | 29.0        | 24.0        | 19.0        | 21.0        | 29.0        | 19.0        | 23.0        |
| TagProp     | 33.0        | 42.0        | 37.0        | 39.0        | 27.0        | <b>32.0</b> | 45.0        | <b>34.0</b> | <b>39.0</b> |
| FastTag     | 32.0        | <b>43.0</b> | 37.0        | <b>46.0</b> | 22.0        | 30.0        | <b>47.0</b> | 26.0        | 34.0        |
| JSVM        | <b>48.5</b> | 38.0        | <b>42.6</b> | 32.7        | <b>31.6</b> | <b>32.2</b> | 42.2        | 29.4        | 34.6        |
| JSVM+Pol(2) | 46.6        | 37.0        | 41.3        | 32.6        | 24.4        | 27.9        | 37.9        | 26.6        | 31.2        |
| JSVM+Pol(3) | 41.5        | 31.3        | 35.7        | 28.5        | 21.3        | 24.4        | 38.0        | 26.1        | 31.0        |

Table 1: Comparison between different versions of joint SVM and other related methods on three benchmark databases. P, R and F1 denote precision, recall and F1 measure respectively.

