



Aalto University
School of Science
and Technology

Structured output prediction for multilabel classification

Hongyu Su

Helsinki Institute for Information Technology HIIT
Department of Computer Science, Aalto University

August 7, 2015

Multilabel classification

- ▶ *Multilabel classification* is an important research field in machine learning.
- ▶ Input variable $\mathbf{x} \in \mathcal{X}$ is in d dimensional input space $\mathcal{X} = \mathbb{R}^d$.
- ▶ Output variable $\mathbf{y} = (y_1, \dots, y_l) \in \mathcal{Y}$ is a binary vector consist of l binary variables $y_j \in \{+1, -1\}$.
- ▶ \mathbf{y} is called a multilabel, y_j is called a microlabel.
- ▶ Output space is composed by a Cartesian product of l sets

$$\mathcal{Y} = \mathcal{Y}_1 \times \dots \times \mathcal{Y}_l, \mathcal{Y}_i = \{+1, -1\}.$$

- ▶ For example, in document classification, a document \mathbf{x} can be classified as “news”, “movie”, and “science”

$$\mathbf{y} = (\underbrace{+1}_{\text{news}}, \underbrace{+1}_{\text{movie}}, \underbrace{-1}_{\text{sports}}, \underbrace{-1}_{\text{politics}}, \underbrace{-1}_{\text{finance}}, \underbrace{+1}_{\text{science}}, \underbrace{-1}_{\text{art}}).$$

- ▶ The goal is to find a mapping function $f \in \mathcal{H}$ that predicts the best values of an output given an input $f : \mathcal{X} \rightarrow \mathcal{Y}$.

Central problems in multilabel classification

- ▶ The size of the output space (searching space) is exponential in the number of microlabels.

$$\mathcal{Y} = \mathcal{Y}_1 \times \cdots \times \mathcal{Y}_l, \mathcal{Y}_i = \{+1, -1\} \quad |\mathcal{Y}| = 2^l.$$

- ▶ The dependency of microlabels needs to be exploited to improve the prediction performance.
 - ▶ If a document is about “movie”, then it is more likely to be about “art” than “science”.

Real world applications

- Social network, information can spread through multiple users.



$$\mathbf{y} = (\underbrace{+1}_{\text{Ted}}, \underbrace{-1}_{\text{Alice}}, \underbrace{+1}_{\text{David}}, \underbrace{-1}_{\text{Mark}}, \underbrace{+1}_{\text{Alex}}, \underbrace{-1}_{\text{Zoe}}, \underbrace{-1}_{\text{Frank}})$$

- Image annotation, an image can associate with multiple tags.



$$\mathbf{y} = (\underbrace{+1}_{\text{boat}}, \underbrace{+1}_{\text{sea}}, \underbrace{-1}_{\text{sun}}, \underbrace{-1}_{\text{beach}}, \underbrace{-1}_{\text{people}}, \underbrace{+1}_{\text{ice}}, \underbrace{+1}_{\text{land}})$$

- Document classification, an article can be assigned to multiple categories.



$$\mathbf{y} = (\underbrace{+1}_{\text{news}}, \underbrace{+1}_{\text{economics}}, \underbrace{-1}_{\text{sports}}, \underbrace{-1}_{\text{politics}}, \underbrace{-1}_{\text{movie}}, \underbrace{-1}_{\text{science}}, \underbrace{-1}_{\text{art}})$$

- Drug discovery, a drug can be effective for multiple symptoms.



$$\mathbf{y} = (\underbrace{+1}_{\text{heart}}, \underbrace{+1}_{\text{stroke}}, \underbrace{+1}_{\text{blood}}, \underbrace{+1}_{\text{fever}}, \underbrace{-1}_{\text{digest}}, \underbrace{-1}_{\text{liver}}, \underbrace{+1}_{\text{swelling}})$$

Flat multilabel classification approaches

- ▶ The categorization is proposed in [?]
- ▶ Problem transformation
 - ▶ Model the multilabel classification as a collection of single-label classification problems and solve each problem independently.
 - ▶ For example, ML-KNN [?], CC [?, ?], IBLR [?].
- ▶ Algorithm adaptation
 - ▶ Modify the single-label classification algorithm for multilabel classification problems.
 - ▶ For example, ADABOOST.MH [?, ?], CORRLOG [?], MTL [?].
- ▶ These approaches does not model the dependency structure explicitly.

Structured output prediction

- ▶ Model the dependency structure with an output graph defined on microlabels.
- ▶ The categorization is proposed in [?].
- ▶ Hierarchical classification
 - ▶ The output graph is a rooted tree or a DAG defining different levels of granularities.
 - ▶ For example, SSVM [?, ?].
- ▶ Graph labeling
 - ▶ The output graph takes a more general form (e.g., a tree, a chain).
 - ▶ For example, CRF [?, ?], M^3N [?], MMCRF [?, ?], SPIN [?].
- ▶ These approaches assume the output graph is known *apriori*.

Contributions

- ▶ Structured output prediction models when the output graph is known.
 - ▶ SPIN for network influence prediction [?].
 - ▶ MMCRF to work with general output graph structures [?].
- ▶ Structured output prediction models working with unknown output graph.
 - ▶ MVE to combine multiple structured output predictors with ensemble [?].
 - ▶ AMM and MAM to aggregate the inference results from multiple structured output predictors [?, ?].
 - ▶ RTA to perform joint learning and inference over a collection of random spanning trees [?].
- ▶ Codes for developed models are available from <http://hongyusu.github.io>.

Outline

- ▶ Preliminaries
- ▶ Structured output learning with known output graph
- ▶ Structured output learning with unknown output graph
- ▶ Future work
- ▶ Experimental results

Preliminaries

- ▶ Training examples come in pairs $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$.
- ▶ $\mathbf{x} \in \mathcal{X}$ is an arbitrary input space.
- ▶ \mathcal{Y} is an output space of a collection of ℓ -dimensional *multilabels*.

$$\mathbf{y} = (y_1, \dots, y_\ell) \in \mathcal{Y}.$$

- ▶ y_i is a *microlabel* and $y_i \in \{1, \dots, r_i\}$, $r_i \in \mathbb{Z}$.
- ▶ For example, multilabel binary classification $y_i \in \{-1, +1\}$.
- ▶ We are given a set of m training examples $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^m$.
- ▶ Each example (\mathbf{x}, \mathbf{y}) is mapped into a joint feature space $\phi(\mathbf{x}, \mathbf{y})$.
- ▶ \mathbf{w} is the weight vector in the joint feature space.
- ▶ Define a linear score function $F(\mathbf{w}, \mathbf{x}, \mathbf{y}) = \langle \mathbf{w}, \phi(\mathbf{x}, \mathbf{y}) \rangle$.
- ▶ \mathbf{w} makes sure example \mathbf{x} with correct multilabel \mathbf{y} achieves higher score than with any other incorrect multilabel $\mathbf{y}' \in \mathcal{Y}$.

Inference problem

- ▶ The prediction $\mathbf{y}_w(\mathbf{x})$ of an input \mathbf{x} is the multilabel \mathbf{y} that maximizes the score function

$$\mathbf{y}_w(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{w}, \phi(\mathbf{x}, \mathbf{y}) \rangle. \quad (1)$$

- ▶ Search space $|\mathcal{Y}| = 2^\ell$ is exponential in size.
- ▶ (??) is called *inference* problem which is \mathcal{NP} -hard for most output feature maps.
- ▶ We aim at using an output feature map in which the inference can be solved with a polynomial algorithm, e.g., dynamic programming.

Input-output feature maps

- ▶ We assume that the joint feature map ϕ is a potential function on a Markov network $G = (E, V)$.
- ▶ A vertex $v_i \in V$ corresponds to a microlabel y_i , an edge $(v_i, v_j) \in E$ corresponds to the pairwise correlation of the microlabel y_i and y_j .
- ▶ G models potential pairwise correlations.

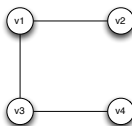


- ▶ $\varphi(\mathbf{x}) \in \mathbb{R}^d$ is the input feature map, e.g., bag-of-words of a document.
- ▶ $\psi(\mathbf{y}) \in \mathbb{R}^{|E|}$ is the output feature map which maps the multilabel \mathbf{y} into a collection of edges and labels

$$\varphi(\mathbf{y}) = (u_e)_{e \in E}, u_e \in \{-1, +1\}^2.$$

An example of output feature map

- ▶ Markov network $G = (E, V)$



- ▶ Multilabel \mathbf{y}

$$\mathbf{y} = (y_1, y_2, y_3, y_4) = (+1, -1, +1, +1)$$

- ▶ Output feature map $\psi(\mathbf{y})$

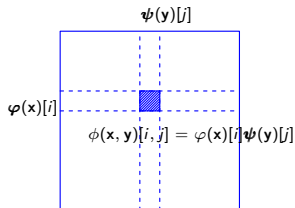
$$\psi(\mathbf{y}) = (\underbrace{0}_{--}, \underbrace{0}_{-+}, \underbrace{1}_{+-}, \underbrace{0}_{++}, \underbrace{0}_{--}, \underbrace{0}_{-+}, \underbrace{0}_{+-}, \underbrace{1}_{++}, \underbrace{0}_{--}, \underbrace{0}_{-+}, \underbrace{0}_{+-}, \underbrace{1}_{++})$$

$\underbrace{\hspace{10em}}_{(v_1, v_3)} \quad \underbrace{\hspace{10em}}_{(v_1, v_2)} \quad \underbrace{\hspace{10em}}_{(v_3, v_4)}$

Joint feature map

- The joint feature is the Kronecker product of $\varphi(\mathbf{x})$ and $\psi(\mathbf{y})$

$$\phi(\mathbf{x}, \mathbf{y}) = (\phi_e(\mathbf{x}, \mathbf{y}))_{e \in E} = (\varphi(\mathbf{x}) \otimes \psi(\mathbf{y}))_{e \in E}.$$

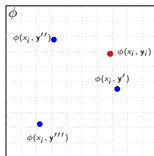


- The score function can be factorized by the output graph G

$$F(\mathbf{w}, \mathbf{x}, \mathbf{y}) = \langle \mathbf{w}, \phi(\mathbf{x}, \mathbf{y}) \rangle = \sum_{e \in E} \langle \mathbf{w}_e, \phi_e(\mathbf{x}, \mathbf{y}_e) \rangle.$$

Optimization problem

- ▶ To learn parameter \mathbf{w} , we aim to maximize the margin between correct pair $(\mathbf{x}_i, \mathbf{y}_i)$ and all the other incorrect pairs $(\mathbf{x}_i, \mathbf{y})$, $\mathbf{y} \in \mathcal{Y}/\mathbf{y}_i$ in the joint feature space ϕ .



- ▶ The model is max-margin conditional random field MMCRF [?, ?].
- ▶ The primal optimization problem is defined as

$$\min_{\mathbf{w}, \xi_k} \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{k=1}^m \xi_k \quad (2)$$

$$\begin{aligned} \text{s.t.} \quad & \langle \mathbf{w}, \phi(\mathbf{x}_k, \mathbf{y}_k) \rangle - \langle \mathbf{w}, \phi(\mathbf{x}_k, \mathbf{y}) \rangle \geq \ell(\mathbf{y}_k, \mathbf{y}) - \xi_k, \\ & \xi_k \geq 0, \forall \mathbf{y} \in \mathcal{Y}, k \in \{1, \dots, m\}. \end{aligned}$$

- ▶ $\ell(\mathbf{y}, \mathbf{y}_i)$ scales the margin according to the multilabel \mathbf{y} .

Marginal-dual optimization

- ▶ (??) is difficult as the number of the constraints is $m \times |\mathcal{Y}|$.
- ▶ The dual optimization problem is defined as

$$\begin{aligned} \max_{\alpha \geq 0} \quad & \alpha^\top \ell - \frac{1}{2} \alpha^\top K \alpha \\ \text{s.t.} \quad & \sum_{\mathbf{y} \in \mathcal{Y}} \alpha(i, \mathbf{y}) \leq C, \forall i \in \{1, \dots, m\}. \end{aligned} \tag{3}$$

- ▶ (??) is also challenging due to the exponential number of dual variables.
- ▶ We use edge marginals to replace the dual variables [?]

$$\mu(i, e, u_e) = \sum_{\mathbf{y}} \mathbf{1}_{\{\psi_e(\mathbf{y}) = u_e\}} \alpha(i, \mathbf{y}).$$

- ▶ The margin-dual optimization problem is

$$\max_{\mu \in \mathcal{M}} \quad \mu^\top \ell - \frac{1}{2} \mu^\top K \mu. \tag{4}$$

- ▶ The number of marginal-dual variable is $m \times 4|E|$.

Conditional gradient optimization

- (??) is optimized by conditional gradient descent which optimizes μ_k that corresponds to a single example while keeps others ($\mu_j, j \neq k$) fixed

$$\max_{\mu_k \in \mathcal{M}} \mu_k^\top \ell_k - \frac{1}{2} \sum_j \mu_k^\top K \mu_j, \forall k.$$

- Current gradient of μ_k is given by $g_i = \ell_i - \sum_j K \mu_j$.
- Compute a feasible solution μ_k^* as an update direction

$$\mu_k^* = \operatorname{argmax}_{\mu_k \in \mathcal{M}} \mu_k^\top g_k = \operatorname{argmax}_{\mu_k \in \mathcal{M}} \sum_e \mu(k, e)^\top g(k, e). \quad (5)$$

- (??) is an instantiation of MAP problem
 - G is tree, exact inference with polynomial algorithm, e.g. dynamic programming in [?]
 - G is general graph, approximate inference, e.g. loopy belief propagation in [?]
- Perform the update via exact line search $\mu_k \leftarrow \mu_k + \tau(\mu_k^* - \mu_k)$.

Compute duality gap

- ▶ We use duality gap to measure the progress of the optimization.
- ▶ Primal and marginal-dual objective functions

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{k=1}^m (\ell_k - \langle \mathbf{w}, \Delta \phi(\mathbf{x}_k, \mathbf{y}_k) \rangle)$$

$$g(\mu) = \sum_{k=1}^m \mu_k \ell_k - \frac{1}{2} \sum_{k=1}^m \sum_{j=1}^m \mu_k K^{\Delta \phi}(\mathbf{x}_k, \mathbf{y}_k; \mathbf{x}_j, \mathbf{y}_j) \mu_j$$

- ▶ $\max_{\mu} g(\mu) \leq \min_{\mathbf{w}} f(\mathbf{w})$, gap is minimized at optimal.
- ▶ Duality gap at μ^t

$$\begin{aligned} f(\mathbf{w}^t) - g(\mu^t) &= C \left(\ell - K^{\Delta \phi} \mu^t \right) - \mu^t \left(\ell - K^{\Delta \phi} \mu^t \right) \\ &= C^T \nabla g(\mu^t) - \mu^{t^T} \nabla g(\mu^t) \end{aligned}$$

1. Estimate the marginal-dual objective by linear approximation $\nabla g(\mu^t)$.
2. Marginal-dual objective value at μ^t is computed by $\mu^{t^T} \nabla g(\mu^t)$.
3. Primal objective value is estimate by $C^T \nabla g(\mu^t)$.









Bibliography



Argyriou, A., Evgeniou, T., and Pontil, M. (2008).

Convex multi-task feature learning.

Machine Learning, 73(3):243–272.



Bian, W., Xie, B., and Tao, D. (2012).

Corrlog: Correlated logistic models for joint prediction of multiple labels.

Journal of Machine Learning Research - Proceedings Track, pages 109–117.



Cheng, W. and Hüllermeier, E. (2009).

Combining instance-based learning and logistic regression for multilabel classification.

Machine Learning, 76(2-3):211–225.



Esuli, A., Fagni, T., and Sebastiani, F. (2008).

Boosting multi-label hierarchical text categorization.

Information Retrieval, 11(4):287–313.

Bibliography (cont.)



Lafferty, J. D., McCallum, A., and Pereira, F. C. N. (2001).

Conditional random fields: Probabilistic models for segmenting and labeling sequence data.

In *Proceedings of the 8th International Conference on Machine Learning (ICML 2001)*, pages 282–289. Morgan Kaufmann Publishers Inc.



Marchand, M., Su, H., Morvant, E., Rousu, J., and Shawe-Taylor, J. (2014).

Multilabel structured output learning with random spanning trees of max-margin markov networks.

In *Advances in Neural Information Processing System NIPS2014*, page to appear.



Read, J., Pfahringer, B., Holmes, G., and Frank, E. (2009).

Classifier chains for multi-label classification.

In Buntine, W., Grobelnik, M., Mladenić, D., and Shawe-Taylor, J., editors, *Machine Learning and Knowledge Discovery in Databases*, volume 5782, pages 254–269. Springer Berlin Heidelberg.

Bibliography (cont.)



Read, J., Pfahringer, B., Holmes, G., and Frank, E. (2011).

Classifier chains for multi-label classification.

Machine Learning, 85(3):333–359.



Rousu, J., Saunders, C., Szedmak, S., and Shawe-Taylor, J. (2007).

Efficient algorithms for max-margin structured classification.

Predicting Structured Data, pages 105–129.



Schapire, R. and Singer, Y. (1999).

Improved boosting algorithms using confidence-rated predictions.

Machine Learning, 37(3):297–336.



Su, H. (2015).

Multilabel Classification through Structured Output Learning - Methods and Applications.

PhD thesis, Department of Information and Computer Science, Aalto University.

Bibliography (cont.)



Su, H., Gionis, A., and Rousu, J. (2014).

Structured prediction of network response.

In Proceedings, 31th International Conference on Machine Learning ICML2014, volume 32 of *Journal of Machine Learning Research WCP*, pages 442–450.



Su, H., Heinonen, M., and Rousu, J. (2010).

Structured output prediction of anti-cancer drug activity.

In Proceedings, 5th IAPR International Conference on Pattern Recognition in Bioinformatics (PRIB2010), volume 6282 of *Lecture Note in Computer Science*, pages 38–49.



Su, H. and Rousu, J. (2011).

Multi-task drug bioactivity classification with graph labeling ensembles.

In Proceedings, 6th IAPR International Conference on Pattern Recognition in Bioinformatics (PRIB2011), volume 7035 of *Lecture Note in Computer Science*, pages 157–167.

Bibliography (cont.)



Su, H. and Rousu, J. (2013).

Multilabel classification through random graph ensembles.

In Proceedings, 5th Asian Conference on Machine Learning (ACML2013), volume 29 of Journal of Machine Learning Research WCP, pages 404–418.



Su, H. and Rousu, J. (2015).

Multilabel classification through random graph ensembles.

Machine Learning, 99(2):231–256.



Taskar, B., Abbeel, P., and Koller, D. (2002).

Discriminative probabilistic models for relational data.

In Proceedings of the Eighteenth Conference on Uncertainty in Artificial Intelligence (UAI 2002), pages 485–492, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.

Bibliography (cont.)



Taskar, B., Guestrin, C., and Koller, D. (2004).

Max-margin markov networks.

In Thrun, S., Saul, L., and Schölkopf, B., editors, *Advances in Neural Information Processing Systems 16*, pages 25–32. MIT Press.



Tsochantaridis, I., Hofmann, T., Joachims, T., and Altun, Y. (2004).

Support vector machine learning for interdependent and structured output spaces.

In *Proceedings of the 21th International Conference on Machine Learning (ICML 2004)*, pages 823–830. ACM.



Tsochantaridis, I., Joachims, T., Hofmann, T., and Altun, Y. (2005).

Large margin methods for structured and interdependent output variables.

Journal of Machine Learning Research, 6:1453–1484.

Bibliography (cont.)



Tsoumakas, G., Katakis, I., and Vlahavas, I. (2010).

Mining multi-label data.

In Maimon, O. and Rokach, L., editors, *Data Mining and Knowledge Discovery Handbook*, pages 667–685. Springer US.



Zhang, M. and Zhou, Z. (2007).

MI-knn: A lazy learning approach to multi-label learning.

Pattern Recognition, 40:2007.