



# Label Distribution Learning Forests

Wei Shen<sup>1,2</sup>, Kai Zhao<sup>1</sup>, Yilu Guo<sup>1</sup>, Alan Yuille<sup>2</sup> <sup>1</sup>School of Communication and Information Engineering, Shanghai University



<sup>2</sup> Department of Computer Science, Johns Hopkins University

## Method

 $\Box$  Our goal is to learn a mapping function g:  $\mathbf{x} \rightarrow \mathbf{d}$  between an input sample  $\mathbf{x}$  and its

 $\square$  We want to learn the mapping function g(x) by a decision tree based model  $\mathcal{T}$ 

 $\varphi(\cdot)$  is an index function to bring the  $\varphi(n)$ -th output of function  $f(x, \Theta)$  in correspondence with split node *n* 

 $s_2^{}(\mathbf{x}; \mathbf{\Theta})$ 

 $\mathbf{f}: \mathbf{x} \to \mathbb{R}^M$  is a real-valued feature learning function 

split nodes  $n \in \mathcal{N}$  each *n* defines a split function  $s_n(\cdot, \Theta): \mathbf{x} \to [0,1], s_n(\mathbf{x}; \Theta) = \sigma(f_{\varphi(n)}(\mathbf{x}; \Theta))$ leaf nodes  $\ell \in \mathcal{L}$  each  $\ell$  holds a distribution  $\mathbf{q}_{\ell} = (q_{\ell_1}, q_{\ell_2}, \dots, q_{\ell_C})^{\mathsf{T}}$ 

**\Box** The probability of the sample **x** falling into leaf node  $\ell$  is given by  $p(\ell | \mathbf{x}; \boldsymbol{\Theta}) = \prod s_n(\mathbf{x}; \boldsymbol{\Theta})^{\mathbf{1}(\ell \in \mathcal{L}_n^l)} (1 - s_n(\mathbf{x}; \boldsymbol{\Theta}))^{\mathbf{1}(\ell \in \mathcal{L}_n^r)}$ 

The output of the tree  $\mathcal{T}$  w.r.t. **x**, is defined by  $\mathbf{g}(\mathbf{x}; \boldsymbol{\Theta}, \mathcal{T}) = \sum p(\ell | \mathbf{x}; \boldsymbol{\Theta}) \mathbf{q}_{\ell}$ 

$$R(\mathbf{q}, \boldsymbol{\Theta}; \mathcal{S}) = -\frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{C} d_{\mathbf{x}_{i}}^{y_{c}} \log(g_{c}(\mathbf{x}_{i}; \boldsymbol{\Theta}, \mathcal{T})) = -\frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{C} d_{\mathbf{x}_{i}}^{y_{c}} \log\left(\sum_{\ell \in \mathcal{L}} p(\ell | \mathbf{x}_{i}; \boldsymbol{\Theta}) q_{\ell_{c}}\right)$$

**D** An alternating optimization strategy to address  $(\Theta^*, \mathbf{q}^*) = \arg \min_{\Theta \in \mathcal{S}} R(\mathbf{q}, \Theta; \mathcal{S})$ 

$$\frac{\partial R(\mathbf{q}, \boldsymbol{\Theta}; \mathcal{S})}{\partial \boldsymbol{\Theta}} = \sum_{i=1}^{N} \sum_{n \in \mathcal{N}} \frac{\partial R(\mathbf{q}, \boldsymbol{\Theta}; \mathcal{S})}{\partial f_{\varphi(n)}(\mathbf{x}_{i}; \boldsymbol{\Theta})} \frac{\partial f_{\varphi(n)}(\mathbf{x}_{i}; \boldsymbol{\Theta})}{\partial \boldsymbol{\Theta}},$$
  
$$\frac{\partial R(\mathbf{q}, \boldsymbol{\Theta}; \mathcal{S})}{\partial f_{\varphi(n)}(\mathbf{x}_{i}; \boldsymbol{\Theta})} = \frac{1}{N} \sum_{c=1}^{C} d_{\mathbf{x}_{i}}^{y_{c}} \Big( s_{n}(\mathbf{x}_{i}; \boldsymbol{\Theta}) \frac{g_{c}(\mathbf{x}_{i}; \boldsymbol{\Theta}, \mathcal{T}_{n}^{r})}{g_{c}(\mathbf{x}_{i}; \boldsymbol{\Theta}, \mathcal{T})} - \big(1 - s_{n}(\mathbf{x}_{i}; \boldsymbol{\Theta})\big) \frac{g_{c}(\mathbf{x}_{i}; \boldsymbol{\Theta}, \mathcal{T}_{n}^{l})}{g_{c}(\mathbf{x}_{i}; \boldsymbol{\Theta}, \mathcal{T})} \Big)$$

 $\min_{\mathbf{q}} R(\mathbf{q}, \boldsymbol{\Theta}; \mathcal{S}), \text{ s.t.}, \forall \ell, \sum_{\mathbf{q}} q_{\ell_c} = 1 \text{ Constrained Convex Optimization Problem!}$ We propose to address this optimization problem by Variational Bounding (VB) [2]

$$R(\mathbf{q}, \boldsymbol{\Theta}; \mathcal{S}) \leq -\frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{C} d_{\mathbf{x}_{i}}^{y_{c}} \sum_{\ell \in \mathcal{L}} \xi_{\ell}(\bar{q}_{\ell_{c}}, \mathbf{x}_{i}) \log\left(\frac{p(\ell|\mathbf{x}_{i}; \boldsymbol{\Theta})q_{\ell_{c}}}{\xi_{\ell}(\bar{q}_{\ell_{c}}, \mathbf{x}_{i})}\right), \ \xi_{\ell}(q_{\ell_{c}}, \mathbf{x}_{i}) = \frac{p(\ell|\mathbf{x}_{i}; \boldsymbol{\Theta})q_{\ell_{c}}}{g_{c}(\mathbf{x}_{i}; \boldsymbol{\Theta}, \mathcal{T})}$$
fine
$$\phi(\mathbf{q}, \bar{\mathbf{q}}) = -\frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{C} d_{\mathbf{x}_{i}}^{y_{c}} \sum_{\ell \in \mathcal{L}} \xi_{\ell}(\bar{q}_{\ell_{c}}, \mathbf{x}_{i}) \log\left(\frac{p(\ell|\mathbf{x}_{i}; \boldsymbol{\Theta})q_{\ell_{c}}}{\xi_{\ell}(\bar{q}_{\ell_{c}}, \mathbf{x}_{i})}\right)$$

 $\phi(\mathbf{q}, \bar{\mathbf{q}}) \geq \phi(\mathbf{q}, \mathbf{q}) = R(\mathbf{q}, \Theta; \mathcal{S})$  and  $\phi(\bar{\mathbf{q}}, \bar{\mathbf{q}}) = R(\bar{\mathbf{q}}, \Theta; \mathcal{S})$  hold the conditions for VB  $\mathbf{q}^{(t+1)} = \arg\min_{\mathbf{q}} \phi(\mathbf{q}, \mathbf{q}^{(t)}), \text{s.t.}, \forall \ell, \sum_{c=1}^{C} q_{\ell_c} = 1 \quad \blacksquare \quad \varphi(\mathbf{q}, \mathbf{q}^{(t)}) = \phi(\mathbf{q}, \mathbf{q}^{(t)}) + \sum_{\ell \in \mathcal{L}} \lambda_{\ell} (\sum_{c=1}^{C} q_{\ell_c} - 1) \text{ Lagrangian!}$  $\lambda_{\ell} = \frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{C} d_{\mathbf{x}_i}^{y_c} \xi_{\ell}(q_{\ell_c}^{(t)}, \mathbf{x}_i) \text{ and } q_{\ell_c}^{(t+1)} = \frac{\sum_{i=1}^{N} d_{\mathbf{x}_i}^{y_c} \xi_{\ell}(q_{\ell_c}^{(t)}, \mathbf{x}_i)}{\sum_{c=1}^{C} \sum_{i=1}^{N} d_{\mathbf{x}_i}^{y_c} \xi_{\ell}(q_{\ell_c}^{(t)}, \mathbf{x}_i)}$ 

 $\Box$  Learning a forest  $\mathcal{F} = \{\mathcal{T}_1, \dots, \mathcal{T}_K\} \mathbf{g}(\mathbf{x}; \mathbf{\Theta}, \mathcal{F}) = \frac{1}{K} \sum_{k=1}^K \mathbf{g}(\mathbf{x}; \mathbf{\Theta}, \mathcal{T}_k)$ 

$$R_{\mathcal{F}} = \frac{1}{K} \sum_{k=1}^{K} R_{\mathcal{T}_{k}} \qquad \qquad \frac{\partial R_{\mathcal{F}}}{\partial \Theta} = \frac{1}{K} \sum_{i=1}^{N} \sum_{k=1}^{K} \sum_{n \in \mathcal{N}_{k}} \frac{\partial R_{\mathcal{T}_{k}}}{\partial f_{\varphi_{k}(n)}(\mathbf{x}_{i}; \Theta)} \frac{\partial f_{\varphi_{k}(n)}(\mathbf{x}_{i}; \Theta)}{\partial \Theta}$$





### Reference

[1] P. Kontschieder, M. Fiterau, A. Criminisi, and S. R. Bulò. Deep neural decision forests. In Proc. ICCV, 2015.

[2] M. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. An introduction to variational methods for graphical models. Machine Learning, 37(2):183–233, 1999.