

## Unleash the Power of Label Space: Label Enhancement

### Xin Geng

PAttern Learning and Mining (PALM) Lab http://palm.seu.edu.cn

School of Computer Science and Engineering Southeast University





## Outline

- Introduction
- A Theoretical View
- Label Enhancement Methods
  - Fuzzy Label
  - Probabilistic Label
  - Label Distribution
- Applications
- Conclusion







## 0/1 Labels

- Most existing data sets: a bipartition of the label set into relevant and irrelevant labels
  - 1: relevant label
  - 0: irrelevant label







## Fuzziness





maturity		
rare		
medium rare		
medium		
medium well		
well done		

The definition of the class labels is fuzzy



## Probability



## PALM

## Ambiguity

Relevant labels

Instances





Class labels of the instance are ambiguous



## $0/1 \text{ Label} \rightarrow \text{Fine Label}$

A real number  $d_x^y$  is assigned to the label y for the instance x

WLOG 
$$ightharpoondows d^y_x \in [0,1]$$
  
Complete label set  $ightharpoondows \sum_y d^y_x = 1$  Fine label





## Fuzziness – Membership Degree



maturity	membership degree
rare	0.05
medium rare	0.2
medium	0.25
medium well	0.35
well done	0.15

The definition of the class labels is fuzzy



## Probability – Probabilistic Labels





0123456789

0 — airplane	2 — bird	4 — deer	6 — frog	8 — ship
1 — automobile	3 — cat	5 — dog	7 — horse	9 — truck





## Ambiguity – Label Distribution



Class labels of the instance are ambiguous



 $\{0, 1, 0, 1, 0\}$ 

## **Practical Restrictions**

- Directly obtaining fine labels is difficult:
  - High cost
  - Difficult to quantify
- 0/1 Labels simplify the real world: a bipartition of the label set into relevant and irrelevant labels
  - 1: relevant label
  - 0: irrelevant label



# PALM

## **Problem Formulation**

The 0/1 label vector of  $x_i$  is denoted by  $l_i = (l_{x_i}^{y_1}, l_{x_i}^{y_2}, ..., l_{x_i}^{y_c})^{\mathrm{T}}$ , where  $l_{x_i}^{y_j} \in \{0,1\}$  represents whether  $y_j$  describes  $x_i$ , c is the number of labels. Then,  $l_i \in \{0,1\}^c$ .

The fine label vector of  $x_i$  is denoted by  $d_i = (d_{x_i}^{y_1}, d_{x_i}^{y_2}, \dots, d_{x_i}^{y_c})^T$ , where  $d_{x_i}^{y_j} \in [0,1]$  represents the fine label of  $y_j$  to  $x_i$ . Then,  $d_i \in [0,1]^c$ .

Label Enhancement can be defined as follows.

Given a training set  $S = \{(x_i, l_i) | 1 \le i \le n\}$ , label enhancement is to recover the fine label vector  $d_i$  of  $x_i$ from the 0/1 label vector  $l_i$ , and thus transform *S* into a fine label training set  $E = \{(x_i, d_i) | 1 \le i \le n\}$ .

## Outline

- Introduction
- A Theoretical View
- Label Enhancement Methods
  - Fuzzy Label
  - Probabilistic Label
  - Label Distribution
- Applications
- Conclusion









- l, x are generated from some conditional distribution p(l, x | d).
- d is generated from the posterior distribution p(d|l, x).
- A fixed-form distribution q(d|l, x) is utilized to approximate p(d|l, x).



*d* should maximize the lower bound of the joint probability density p(l, x):  $\log p(l, x) \ge \mathbb{E}_{q(d|l,x)}[\log p(l|d) + \log p(x|d)] - \mathrm{KL}[q(d|l, x)||p(d)]$ 

# PALM

### A Theoretical View

We formulate the label enhancement problem into an optimization framework and yields the target function for minimization:

$$T(\boldsymbol{\vartheta}, \boldsymbol{\eta}, \boldsymbol{w}) = \frac{1}{L} \sum_{m=1}^{L} \frac{1}{2} \|\boldsymbol{x} - \boldsymbol{\rho}^{(m)}\|_{2}^{2} - \sum_{i=1}^{c} l_{i} \log \tau_{i}^{(m)} + (1 - l_{i}) \cdot \log \left(1 - \tau_{i}^{(m)}\right) \\ + \frac{1}{2} \{ \operatorname{tr}(\boldsymbol{\Sigma}) + \boldsymbol{\mu}^{\top} \boldsymbol{\mu} - k - \log |\boldsymbol{\Sigma}| \},$$

where  $\Sigma = MLP_{\Sigma}(\boldsymbol{l}, \boldsymbol{x}; \boldsymbol{w}), \boldsymbol{\mu} = MLP_{\mu}(\boldsymbol{l}, \boldsymbol{x}; \boldsymbol{w}), \boldsymbol{\tau}^{(m)} = MLP_{\tau}(\boldsymbol{d}^{(m)}; \boldsymbol{\vartheta}),$  $\boldsymbol{\rho}^{(m)} = MLP_{\boldsymbol{\rho}}(\boldsymbol{d}^{(m)}; \boldsymbol{\eta}), \boldsymbol{d}^{(m)} = \boldsymbol{\mu} + \boldsymbol{\Sigma}\boldsymbol{\epsilon}^{(m)}, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}).$ 

## Outline

- Introduction
- A Theoretical View
- Label Enhancement Methods
  - Fuzzy Label
  - Probabilistic Label
  - Label Distribution
- Applications
- Conclusion







## Label Enhancement Methods

#### Label Enhancement for Fuzzy Label

- LE based on fuzzy clustering [Gayar et al., ANNPR'06]
- LE based on kernel method
   [Jiang et al., NCA'06]

#### Label Enhancement for Probabilistic Label

- LE based on probabilistic end-to-end noisy correction
   [Kun Yi et al., CVPR'19]
- LE based on label smoothing [Szegedy et al., CVPR'16]
- LE based on distillation [Hinton et al., arxiv'15]

#### Label Enhancement for Label Distribution

- LE based on manifold learning
   [Hou et al., AAAI'16]
- Graph laplacian label enhancement
  [Xu and Geng, IJCAI'18]
- LE based on reinforcement learning
   [Gao and Geng, IJCAI'20]



#### LE based on fuzzy clustering

[Gayar et al., ANNPR'06]



Fuzzy C-means clustering (The membership of the instance to the cluster) The memberships of the instances belonging to the same class are added up to form the cluster-class connection matrix.



By fuzzy composition operation, the memberships of instances to clusters are transformed into the memberships of instances to class labels using the connection matrix.



#### LE based on fuzzy clustering

[Gayar et al., ANNPR'06]

#### • Step 1: Fuzzy C-Means clustering (FCM)

- 1. Given the cluster number p, initialize the  $n \times p$  cluster membership matrix M ( $m_{ik}$  denotes the membership of  $x_i$  to the k-th cluster)
- 2. Calculate the cluster prototype

$$\boldsymbol{\mu}_k = \frac{\sum_{i=1}^n (m_{ik})^{\beta} \boldsymbol{x}_i}{\sum_{i=1}^n (m_{ik})^{\beta}}$$

3. Update the cluster membership matrix M

$$m_{ik} = \frac{1}{\sum_{j=1}^{p} \left(\frac{Dist(\boldsymbol{x}_{i}, \boldsymbol{\mu}_{k})}{Dist(\boldsymbol{x}_{i}, \boldsymbol{\mu}_{j})}\right)^{\frac{2}{\beta-1}}}$$

4. Repeat 2 and 3 until convergence

Each row of M,  $m_i$ , represents the membership of the instance  $x_i$  to each cluster



#### LE based on fuzzy clustering

[Gayar et al., ANNPR'06]

- Step 2: Calculate the cluster-class connection matrix
  - 1. Initialize  $c \times p$  zero matrix A
  - 2. Update each row  $A_j$  with

$$\pmb{A}_j = \pmb{A}_j + \pmb{m}_i$$
 ,

- 3. Normalized each column of A
- 4. Normalized each row of *A*

$$if \ l_{x_i}^{y_j} = 1$$

 $a_{jk}$  denotes the connection between class *j* and cluster *k* 

• Step 3: Calculate the fine labels of  $x_i$ 

1. 
$$D_i = A \circ m_i$$
 (fuzzy composition)  
 $D_i^j = \max_k (a_{jk} \times m_{ik})$ 

2. Normalize  $D_i$ 



#### LE based on kernel method

[Jiang et al., NCA'06]



Introduce nonlinearity via kernel method



#### LE based on kernel method

[Jiang et al., NCA'06]

• Step 1: For each label  $y_j$ , suppose  $C^{y_j}$  contains all the instances labeled by  $y_j$ , the size of  $C^{y_j}$  is  $n_j$ , then, the center of  $C^{y_j}$  is

where  $\phi(x_i)$  is a nonlinear function determined by the kernel function  $K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$ 

• Step 2: Calculate the class radius  $r_j = \max_{x_i \in C^{y_j}} \| \Psi^{y_j} - \phi(x_i) \|, \qquad r_j^2$ 

 $r_j^2$  can be calculated via inner product of  $\phi(x_i)$ 

• Step 3: Calculate the distance between instance  $x_i$  and class center

$$d_{ij}^2 = \| \boldsymbol{\phi}(\boldsymbol{x}_i) - \boldsymbol{\Psi}^{y_j} \|^2$$

 $d_{ij}^2$  can be calculated via inner product of  $\phi(x_i)$ 

#### LE based on kernel method

[Jiang et al., NCA'06]

• Step 4: calculate the membership of instance  $x_i$  to label  $y_j$ 

$$m_{x_{i}}^{y_{j}} = \begin{cases} 1 - \sqrt{\frac{\left\|d_{ij}^{2}\right\|}{\left(r_{j}^{2} + \delta\right)}} & \text{if } l_{x_{i}}^{y_{j}} = 1 \\ 0 & \text{if } l_{x_{i}}^{y_{j}} = 0 \end{cases} \begin{cases} \text{Cannot change the membership} \\ \text{of irrelevant labels} \end{cases}$$

• Step 5: Normalize  $m_{x_i} = [m_{x_i}^{y_1}, m_{x_i}^{y_2}, ..., m_{x_i}^{y_c}]$ 

## Outline

- Introduction
- A Theoretical View
- Label Enhancement Methods
  - Fuzzy Label
  - Probabilistic Label
  - Label Distribution
- Applications
- Conclusion







[Kun Yi et al., CVPR'19]

Step 1: Probabilistic label is initialized by the noisy label  $\hat{y}$  $\widetilde{y} = \widetilde{K} \widehat{y}$  K is a large constant

Step 2: Normalize the probabilistic label to a probability distribution

$$y^d = softmax(\widetilde{y})$$

Step 3: Update both the network and the probabilistic label  $y^d$ 

# Fion

#### LE based on probabilistic end-to-end noisy correction

[Yi and Wu, CVPR'19]

Loss 1: Compatibility loss  $\mathcal{L}_o$ : cross entropy between  $\widehat{\mathbf{Y}}$  and  $\mathbf{Y}^d$ 

$$\mathcal{L}_o(\widehat{\mathbf{Y}}, \mathbf{Y}^d) = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^c \widehat{\mathbf{y}}_{ij} \log \mathbf{y}_{ij}^d$$

Loss 2: Classification loss  $\mathcal{L}_c$ : K-L divergence been the output and  $Y^d$ 

$$\mathcal{L}_{c}(f(\boldsymbol{x};\boldsymbol{\theta}),\boldsymbol{Y}^{d}) = \frac{1}{n} \sum_{i=1}^{n} \mathrm{KL}(f(\boldsymbol{x}_{i};\boldsymbol{\theta})||\boldsymbol{y}_{i}^{d})$$

Loss 3: Entropy loss  $\mathcal{L}_e$ : a regularization term to force the network to peak at only one category rather than being flat

$$\mathcal{L}_e(f(\boldsymbol{x};\boldsymbol{\theta})) = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^c f_j(\boldsymbol{x};\boldsymbol{\theta}) \log f_j(\boldsymbol{x};\boldsymbol{\theta})$$

The over all framework

$$\mathcal{L} = \frac{1}{n} \mathcal{L}_c(f(\boldsymbol{x}; \boldsymbol{\theta}), \boldsymbol{Y}^d) + \alpha \mathcal{L}_o(\widehat{\boldsymbol{Y}}, \boldsymbol{Y}^d) + \frac{\beta}{c} \mathcal{L}_e(f(\boldsymbol{x}; \boldsymbol{\theta}))$$



Step 1: After label smoothing, the distribution changes as follows

$$P_i = \begin{cases} 1, if(i = y) \\ 0, if(i \neq y) \end{cases} \implies P_i = \begin{cases} 1 - \varepsilon, if(i = y) \\ \frac{\varepsilon}{K - 1}, if(i \neq y) \end{cases} \quad \mathsf{K} = \# \text{ labels} \end{cases}$$

• Step 2: After label smoothing, the loss function changes as follows

$$Loss = -\sum_{i=1}^{K} p_i logq_i \implies Loss_i = \begin{cases} (1-\varepsilon) * Loss, if(i=y) \\ \varepsilon * Loss, if(i\neq y) \end{cases}$$



#### LE based on knowledge distillation

[Hinton et al., arxiv'15]



Neural networks typically produce class probabilities by using a "softmax" output layer that converts the logit,  $z_i$ , computed for each class into a probability,  $q_i$ , by comparing  $z_i$  with the other logits.

$$q_i = \frac{\exp\left(\frac{Z_i}{T}\right)}{\sum_j \exp\left(\frac{Z_j}{T}\right)}$$

T is a temperature that is normally set to 1. Using a higher value for T produces a smoother probability distribution over classes.

## Outline

- Introduction
- A Theoretical View
- Label Enhancement Methods
  - Fuzzy Label
  - Probabilistic Label
  - Label Distribution
- Applications
- Conclusion







#### LE based on manifold learning

[Hou, Geng and Zhang, AAAI'16]

- Feature space: continuous Euclidean space
- Label space: discrete logical space



The manifold structure is transferred from the feature space to the label space.



### LE based on manifold learning

[Hou, Geng and Zhang, AAAI'16]

0

• Manifold learning in feature space [Roweis & Saul, Science, 2000]

$$\begin{aligned} \arg\min_{W} & \sum_{i=1}^{n} \| \boldsymbol{x}_{i} - \sum_{j \neq i} W_{i}^{j} \boldsymbol{x}_{j} \|^{2} \\ \text{s.t.} & \mathbf{1}^{\mathrm{T}} \boldsymbol{W}_{i} = 1 \\ \text{Manifold learning} & \text{Local topological structure} \\ \arg\min_{\mu} & \sum_{i=1}^{n} \| \boldsymbol{\mu}_{i} - \sum_{j \neq i} W_{i}^{j} \boldsymbol{\mu}_{j} \|^{2} \\ \text{s.t.} & \forall 1 \leq i \leq n, 1 \leq l \leq q \qquad y_{i}^{l} \boldsymbol{\mu}_{i}^{l} \geq \lambda, \lambda > \end{aligned}$$



#### **Graph Laplacian Label Enhancement**

#### [Xu and Geng, IJCAI'18]

Model

Nonlinear transformation  $D_i = W^{\mathsf{T}} \varphi(x_i) + b = \widehat{W} \phi_i$ 

**Goal** Determining the best parameter  $\widehat{W}^*$ 

Target function

 $\min_{\widehat{W}} L(\widehat{W}) + \lambda \Omega(\widehat{W})$ Feature space constraint



#### **Graph Laplacian Label Enhancement**

[Xu and Geng, IJCAI'18]

The first part of the target function

$$L(\widehat{W}) = \sum_{i=1}^{n} \|\widehat{W}\phi_{i} - L_{i}\|^{2} \text{ Least squares (LS)}$$

The second part of the target function



## LE based on reinforcement learning

[Gao and Geng, IJCAI'20]

Leveraging the prior knowledge





Prior knowledge in Emotion Relation [Mikels, et al. BRM, 2005]

The properties implied in the prior knowledge

🖌 ground-truth hard label

1. 
$$d_x^{\alpha} > d_x^{\beta}$$
,  $\alpha \neq \beta$   
1.  $d_x^{y_g} > d_x^{y_i}$ ,  $i \neq g$ 

2.  $d_x^{\alpha \pm i} > d_x^{\alpha \pm j}$ , j > i, i, j > 02.  $d_i^{y_i} > d_i^{y_j}$ , distance(j, g) > distance(i, g)



## LE based on reinforcement learning

[Gao and Geng, IJCAI'20]



Reinforcement learning for LE



$$\begin{split} Q(s,a;\theta,\phi,\beta) &= V(s;\theta,\beta) + \\ (A(s,a;\theta,\phi) - \frac{1}{|A|} \sum_{a^{,}} A(s,a^{,};\theta,\phi)), \end{split}$$

## Outline

- Introduction
- A Theoretical View
- Label Enhancement Methods
  - Fuzzy Label
  - Probabilistic Label
  - Label Distribution
- Applications
- Conclusion







### The Power of Label Space





0/1 Label Space

Fine Label Space









### Linear Discriminant Analysis

[Zhao et al., CVIU'14]

Sea

[Zhao et al. Neurocomputing, 2015]

Beach

#### Soft Label Linear Discriminant Analysis (SL-LDA)











## Outline

- Introduction
- A Theoretical View
- Label Enhancement Methods
  - Fuzzy Label
  - Probabilistic Label
  - Label Distribution
- Applications
- Conclusion





# PALM

## Conclusion

#### Label enhancement

- recovers fine labels (e.g., label distribution, probabilistic labels, fuzzy labels) from 0/1 labels.
- could be theoretically explained via variational inference.
- offers more possibilities for operations in the label space (e.g., linear discriminant analysis, label distribution learning, model compression, neural network regularization, label embedding).



## Interested in LDL & LE?

#### All the papers, codes and datasets are available at: http://palm.seu.edu.cn/xgeng/LDL/index.htm

LDL Home Download Contact Label Distribution For real applications where the overall distrib A more general learning framework which in	Learning bution of the importance of the labels matters. cludes both single-label and multi-label learning as its special cases.	
Introduction		Introduction
Label Distribution Learning is a novel machine learning paradigm. A label distribution covers a certain number of labels, representing the degree to which each label describes the instance. LDL is a general learning framework which includes both single-label and multi-label learning as its special cases.		
X. Geng, Label Distribution Learning, IEEE Transactions on Kno	wiedge and Data Engineering (IEEE TKDE), 2016, in press.	
Our alogrithms can be used freely for academic, non-profit purposes	If you intend to use it for commercial development, please contact us.	
In academic papers using our codes and data, the following reference	es will be appreciated:	
[1] X. Geng. Label Distribution Learning. IEEE Transactions on	Knowledge and Data Engineering (IEEE TKDE), 2016, in press.	
[2] X. Geng, C. Yin, and ZH. Zhou. Facial Age Estimation by Le	aming from Label Distributions. IEEE Transactions on Pattern Analysis and Machine Intelligence (IEEE TPAMI), 2013, 35(10): 2401-2412.	
Applications of LDL		
		1
	Facial Age Estimation	
	X Geno, O. Wang, and Y. Xia. Facial Age Estimation by Adaptive Label Distribution Learning. In: Proceedings of the 22nd International Conference on Pattern Recognition	
	(ICPR'14), Stockholm, Sweden, 2014, pp. 4465 - 4470.	
o fast 20 slow 30 fast 70	<ul> <li>A. Genig, C. Im, and ZH. Zhou, Facar Age Estimation by Learning non-Laber cisal Doubles. IEEE Transactions on Pattern Analysis and Machine Intelligence (IEEE TPAMI), 2013, 35(10): 2401-2412.</li> </ul>	
	<ul> <li>X. Geng, K. Smith-Miles, ZH. Zhou. Facial Age Estimation by Learning from Label Distributions. In: Proceedings of the 24th AAAI Conference on Artificial Intelligence</li> </ul>	1





#### http:// palm.seu.edu.cn