
DW-ML-kNN: A Dual Weighted Multi-label kNN Algorithm



Duoqian Miao, Zhifei Zhang*

Tongji University

Department of Computer Science and Technology

22 September, 2012



Outline

- n Multi-label Objects
- n Multi-label Learning
- n DW-ML-kNN Algorithm
- n Experiments
- n Conclusion



Outline

- n Multi-label Objects
- n Multi-label Learning
- n DW-ML-kNN Algorithm
- n Experiments
- n Conclusion

Multi-label Objects

n Text

Diaoyu Islands issue

CNC report from Beijing

Added On September 14, 2012

The spokesman called the purchase 'illegal and invalid,' warning that China has made its strong opposition clear - and is now taking military measures, to safeguard its sovereignty.

An official with the Ministry of Agriculture said Thursday that China will conduct routine patrols near the Diaoyu Islands to preserve the country's territorial integrity, and protect its fishermen.

A vice minister of commerce said Thursday the move will inevitably have a negative impact on Sino-Japan economic and trade ties.

The official said some Japanese enterprises have already begun to feel the effects, after the government carried out its "nationalization" of the Diaoyu Islands.

Chinese tourists and travel agencies are also canceling trips to Japan, which they would otherwise have visited over the highly lucrative, week-long upcoming October National holiday...

Shandong Traffic and Communication Tourism Group, a travel consultancy, on Thursday announced it would suspend all trips to Japan, in protest against the country's "purchase" of the Diaoyu Islands.

Politics
Military
Tour
Economy
.....

Multiple labels

Multi-label Objects (Cont')

n Image

Clouds

Mountains

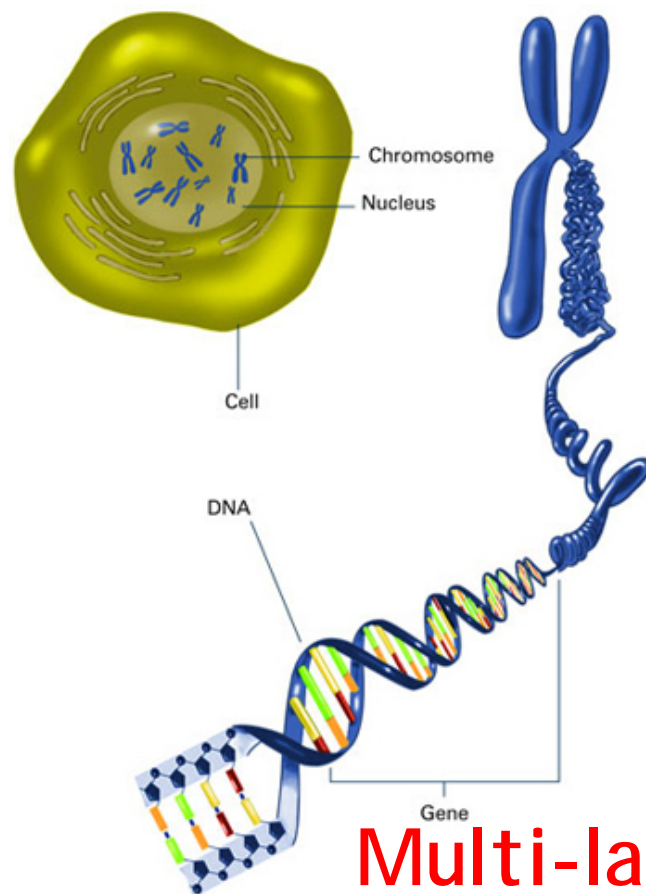
Trees



Lake

Multi-label Objects (Cont')

n Genomics



Metabolism
Transcription
Protein Synthesis

.....

Multi-label objects are ubiquitous!

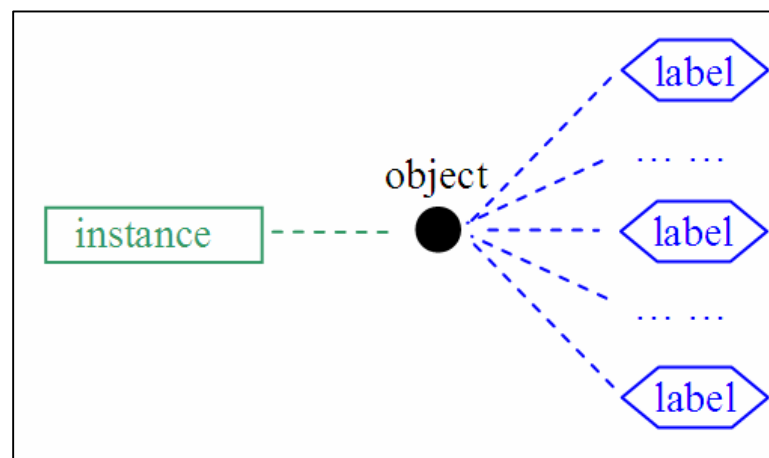
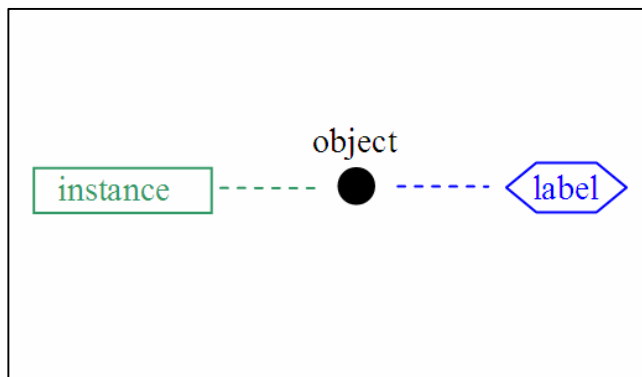


Outline

- n Multi-label Objects
- n **Multi-label Learning**
- n DW-ML-kNN Algorithm
- n Experiments
- n Conclusion

Multi-label Learning

n MLL VS SLL (Single-label Learning)



Multi-label Learning

n Formal Definition

- q d -dimensional feature space $\mathcal{C} = \mathbb{R}^d$
- q label space with q labels $\mathcal{Y} = \{l_1, l_2, \mathbf{L}, l_q\}$
- q Inputs: training set with m examples
$$S = \{(x_i, Y_i) \mid x_i \in \mathcal{C}, Y_i \subseteq \mathcal{Y}\} (i = 1, 2, \mathbf{L}, m)$$
- q Outputs:
 - n multi-label predictor $h: \mathcal{C} \rightarrow 2^{\mathcal{Y}}$
 - n a ranking function $f: \mathcal{C} \times \mathcal{Y} \rightarrow \mathbb{R}$



Multi-label Learning (Cont')

n Applications

- q Text classification
- q Image annotation
- q Functional genomics
- q

n Methods

- q Problem transformation methods Fit data to algorithm
- q Algorithm adaptation methods Fit algorithm to data



Outline

- n Multi-label Objects
- n Multi-label Learning
- n DW-ML-kNN Algorithm
- n Experiments
- n Conclusion

ML-kNN [Zhang&Zhou, PRJ07]

n Basic idea

- q kNN + MAP with neighbors' labeling information

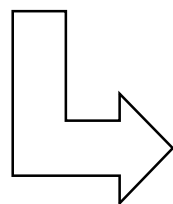
n Settings

- q the k nearest neighbors of x identified in the training set $N(x)$
- q the number of examples in $N(x)$ having the l -th label $C_x(l)$
- q the hypothesis that an example have (not) the l -th label $H_1^l(H_0^l)$
- q the event that there are exactly j examples in $N(x)$ having the l -th label E_j^l

ML-kNN (Cont')

n MAP

$$l \in \mathcal{Y} \text{ iff } P(H_1^l | E_{C_x(l)}^l) > P(H_0^l | E_{C_x(l)}^l)$$



$$P(H_b^l | E_{C_x(l)}^l) \propto P(H_b^l) P(E_{C_x(l)}^l | H_b^l)$$

Probabilities needed:

$$P(H_b^l) \ (1 \leq l \leq q, b \in \{0, 1\})$$

$$P(E_j^l | H_b^l) \ (0 \leq j \leq k)$$

directly estimated from
the training set based on
frequency counting

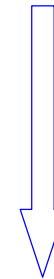
DW-ML-kNN

n Why?

- q ML-kNN performs poor when dealing with imbalanced data.



compute probabilities
by frequency counting



tend to assign labels
with high frequency

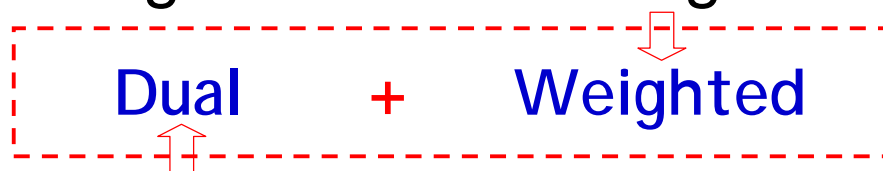
DW-ML-kNN (Cont')

n How?

- q convert distances to weights
- q find neighbors with or without a certain label

n Idea

- q assign higher weights to closer neighbors, and lower weights to farther neighbors


$$\text{Dual} + \text{Weighted}$$

- q higher probability with a label but lower without it, assign the example this label to a large extent

DW-ML-kNN (Cont')

n “Weighted”

- q the probability that example x with the same label of its one neighbor a

$$P(a | x) = \frac{1}{\sqrt{2p}} e^{-\frac{d^2(x,a)}{2}}$$

- q the weighted value of example x with the label l

$$w_x(l) = \frac{\sum_{a \in N(x) \wedge y_a(l)=1} P(a | x)}{\sum_{a \in N(x)} P(a | x)}$$

~~frequency counting~~

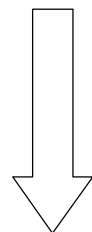
replace

DW-ML-kNN (Cont')

n “Dual”

q two posterior probabilities

$$y_t(l) = \arg \max_{b \in \{0,1\}} \{ P(H_b^l | E_{C_t(l)}^l) + P(H_b^l | E_{K-C_t(l)}^{\sim l}) \}$$



the event that there are exactly $K-C_t(l)$ examples in $N(t)$ without the l -th label

$$P(H_b^l) [P(E_{C_t(l)}^l | H_b^l) + P(E_{K-C_t(l)}^{\sim l} | H_b^l)]$$

Probabilities needed:

$$P(H_b^l)$$

$$P(E_{C_t(l)}^l | H_b^l)$$

$$P(E_{K-C_t(l)}^{\sim l} | H_b^l)$$

directly estimated from the training set based on distance weighting

DW-ML-kNN (Cont')

n Process

- q Step1: compute the prior probabilities from the training set



- q Step2: compute the posterior probabilities from the training set based on **distance weighting**



- q Step3: compute the label set and label ranking for the unseen example based on **dual posteriors**



Outline

- n Multi-label Objects
- n Multi-label Learning
- n DW-ML-kNN Algorithm
- n Experiments
- n Conclusion

Experiments

n Datasets

- q Yeast: biological dataset about protein function classification
- q Scene: image dataset about semantic indexing of still scenes
- q Emotions: music dataset about song classification by emotions

Dataset	Domain	$ S $	$D(S)$	$L(S)$	$LC(S)$	$LD(S)$
Yeast	Biology	2417	103	14	4.237	0.303
Scene	Multimedia	2712	294	6	1.074	0.179
Emotions	Music	593	72	6	1.869	0.311

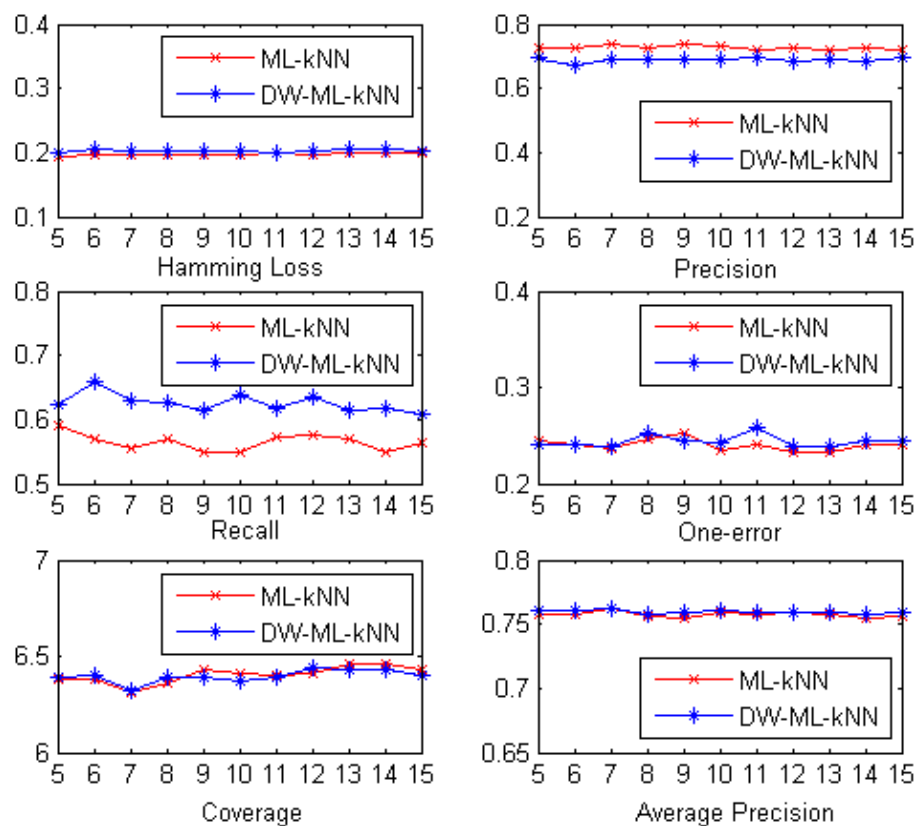
Experiments (Cont')

n Evaluation

- q One-error \downarrow : Average times the top-ranked label is not in the set of proper labels of the example
 - q Coverage \downarrow : Average steps are needed to go down the label list to cover the true label set
 - q Average precision \uparrow : Average fraction of labels ranked above a proper label in the true label set
 - q Hamming loss \downarrow : Average times an example-label pair is misclassified
 - q Precision \uparrow : Average fraction of truly predicted labels of the predicted labels
 - q Recall \uparrow : Average fraction of truly predicted labels of the true labels
- } Ranking-based
- } Classification-based

Experiments (Cont')

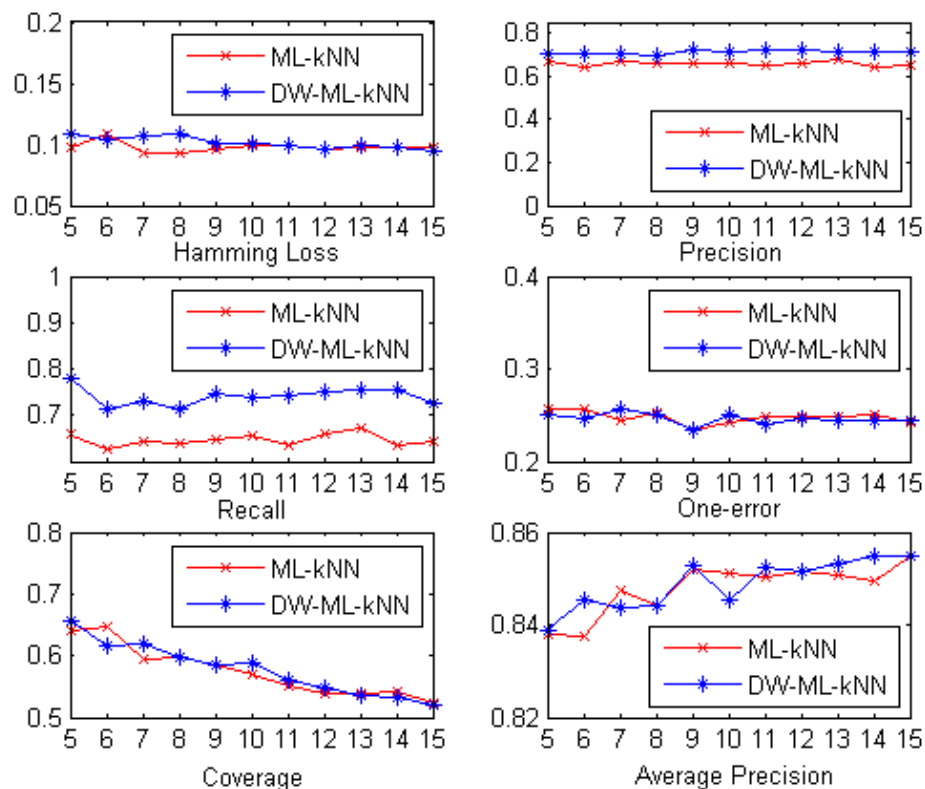
n Yeast



$k=7$

Experiments (Cont')

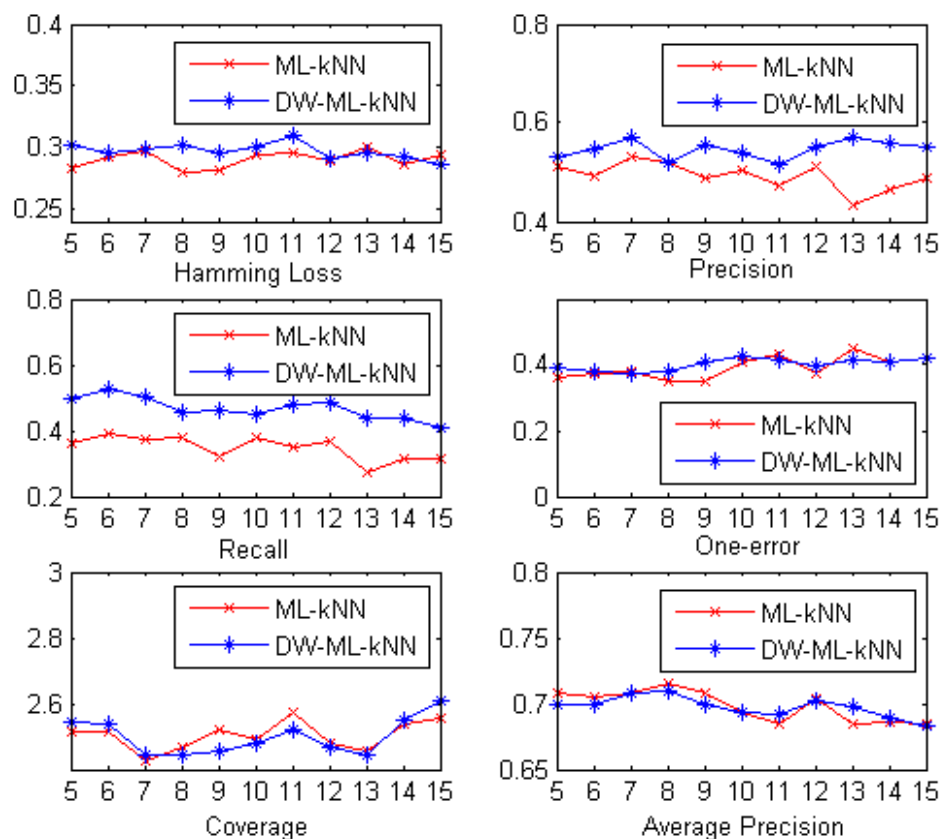
n Scene



$k=9$

Experiments (Cont')

n Emotions



$k=12$

Experiments (Cont')

n DW-ML-kNN VS ML-kNN

Metric	Yeast		Scene		Emotions	
	ML-kNN	DW-ML-kNN	ML-kNN	DW-ML-kNN	ML-kNN	DW-ML-kNN
Hamming loss	0.1973	0.2028	0.0978	0.1013	0.2900	0.2981
Precision	0.7273	0.6880	0.6583	0.7106	0.4928	0.5468
Recall	0.5653	0.6275	0.6471	0.7625	0.3510	0.4781
One-error	0.2405	0.2398	0.2481	0.2462	0.3896	0.4012
Coverage	6.4015	6.3956	0.5804	0.5828	2.4980	2.4888
Average precision	0.7576	0.7595	0.8472	0.8490	0.6898	0.6994
Win(s)	2(6)	4(6)	2(6)	4(6)	2(6)	4(6)

ü DW-ML-kNN achieves better performance

ü Recall is obviously improved

ü Hamming loss isn't improved, but close



Outline

- n Multi-label Objects
- n Multi-label Learning
- n DW-ML-kNN Algorithm
- n Experiments
- n Conclusion

Conclusion

- n The proposed algorithm
 - q DW-ML-kNN is better than ML-kNN on the whole
 - q Deal with imbalanced data better to some extent
- n Further exploiting label correlations
 - q a document labeled as *politics* would be **unlikely** labeled as *entertainment*
 - q an image labeled as *trees* and *lake* would be **likely** labeled as *clouds*



Thank you!

Q & A?

Email: zzf_tj01@126.com

Weibo: @同济志飞