

# Advances in **Cost-sensitive** Multiclass and Multilabel Classification

Hsuan-Tien Lin

htlin@csie.ntu.edu.tw

Professor  
Dept. of CSIE, National Taiwan Univ.



Chief Data Science Consultant  
Appier Inc.

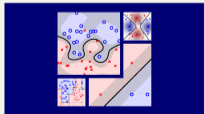
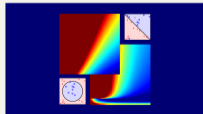


**Tutorial for KDD @ Anchorage, Alaska, USA  
August 4, 2019**

# More About Me



- co-author of textbook '*Learning from Data: A Short Course*'
- instructor of two Coursera Mandarin-teaching ML MOOCs on Coursera



## goal: make ML more realistic

- online/active learning: in ICML '12, ICML '14, AAAI '15, ...
- **cost-sensitive classification**: in ICML '10, KDD '12, IJCAI '16, ...
- **multi-label classification**: in NeurIPS '12, ICML '14, AAAI '18, ...
- large-scale data mining: co-led **KDDCup world-champion** NTU teams 2010–2013
- AI for digital marketing: Chief Data Scientist of **growing startup Appier** 2016–2019

# Outline

## 1 Cost-Sensitive Multiclass Classification

- CSMC Motivation and Setup
- CSMC by Bayesian Perspective
- CSMC by (Weighted) Binary Classification
- CSMC by Regression

## 2 Cost-Sensitive Multilabel Classification

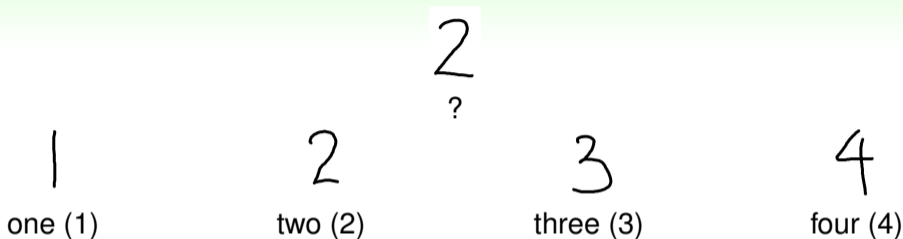
- CSML Motivation and Setup
- CSML by Bayesian Perspective
- CSML by (CS) Multiclass Classification
- CSML by (Weighted) Binary Classification
- CSML by Regression

## 3 CSMC & CSML: Selected Applications

- Bacteria Classification with Doctor-Annotated Costs
- Network Connection Classification with 'Danger' Costs on Imbalanced Data
- Social Tagging with Costs from Tag Counts

## 4 Summary

# Which Digit Did You Write?



- a **multiclass classification** problem: grouping 'pictures' into different 'categories'

**C'mon, we know about  
multiclass classification all too well! :-)**

# Performance Evaluation



- ZIP code recognition:  
1: **w**rong; 2: **r**ight; 3: **w**rong; 4: **w**rong
- check value recognition:  
1: **o**ne-dollar mistake; 2: **n**o mistake;  
3: **o**ne-dollar mistake; 4: **t**wo-dollar mistake

different applications: **evaluate mis-predictions differently**

# ZIP Code Recognition

2  
?

1: wrong; 2: right; 3: wrong; 4: wrong

- **regular** multiclass classification: only right or wrong
- **wrong cost: 1**; **right cost: 0**
- prediction error of  $h$  on some  $(\mathbf{x}, y)$ :

$$\text{classification cost} = \mathbb{1}[y \neq h(\mathbf{x})]$$

regular multiclass classification:

**well-studied**, many good algorithms

# Check Value Recognition

2  
?

- 1: one-dollar mistake; 2: no mistake;  
3: one-dollar mistake; 4: two-dollar mistake

- **cost-sensitive** multiclass classification: **different costs** for different mis-predictions
- e.g. prediction error of  $h$  on some  $(\mathbf{x}, y)$ :

$$\text{absolute cost} = |y - h(\mathbf{x})|$$

next: more about **cost-sensitive multiclass** classification (CSMC)

# What is the Status of the Patient?



?

(image by mcmurryjulie from Pixabay)



bird flu



cold



healthy

(images by Clker-Free-Vector-Images from Pixabay)

- another **classification** problem: grouping 'patients' into different 'status'

**are all mis-prediction costs equal?**



# Patient Status Prediction

error measure = society cost

actual \ predicted	bird flu	cold	healthy
bird flu	0	1000	<b>100000</b>
cold	100	0	3000
healthy	100	30	0

- bird flu mis-predicted as healthy: **very high cost**
- cold mis-predicted as healthy: **high cost**
- cold correctly predicted as cold: **no cost**

human doctors consider costs of decision;  
**can computer-aided diagnosis do the same?**

## Setup: Class-Dependent Cost-Sensitive Classification

## Given

$N$  examples, each (input  $\mathbf{x}_n$ , label  $y_n$ )  $\in \mathcal{X} \times \{1, 2, \dots, K\}$

and cost matrix  $C \in \mathbb{R}^{K \times K}$  with  $C(y, y) = 0 = \min_{1 \leq k \leq K} C(y, k)$

patient diagnosis  
with society cost

$$C = \begin{pmatrix} 0 & 1000 & 100000 \\ 100 & 0 & 3000 \\ 100 & 30 & 0 \end{pmatrix}$$

check digit recognition  
with absolute cost

(cost function)

$$C(y, k) = |y - k|$$

## Goal

a classifier  $g(\mathbf{x})$  that pays a small cost  $C(y, g(\mathbf{x}))$  on future **unseen** example  $(\mathbf{x}, y)$

includes **regular classification**  $C_c$  like  $\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$  as **special case**

## Which Age-Group?



?



infant (1)



child (2)



teen (3)



adult (4)

(images by Tawny van Breda, Pro File, Mieczysław Samol, lisa runnels, vontoba from Pixabay)

- small mistake—classify child as teen; **big mistake—classify infant as adult**

- cost matrix  $\mathcal{C}(y, g(x))$  for embedding ‘order’:  $\mathcal{C} = \begin{pmatrix} 0 & 1 & 4 & 5 \\ 1 & 0 & 1 & 3 \\ 3 & 1 & 0 & 2 \\ 5 & 4 & 1 & 0 \end{pmatrix}$

CSMC can help solve many other problems like **ordinal ranking**

# Cost Vector

cost vector  $\mathbf{c}$ : a row of cost components

- society cost for a bird flu patient:  $\mathbf{c} = (0, 1000, 100000)$
- absolute cost for digit 2:  $\mathbf{c} = (1, 0, 1, 2)$
- age-ranking cost for a teenager:  $\mathbf{c} = (3, 1, 0, 2)$
- 'regular' classification cost for label 2:  $\mathbf{c}_c^{(2)} = (1, 0, 1, 1)$
- movie recommendation
  - someone who loves romance movie but **hates terror**:

$$\mathbf{c} = (\text{romance} = 0, \text{fiction} = 5, \text{terror} = 100)$$

- someone who loves romance movie but **fine with terror**:

$$\mathbf{c} = (\text{romance} = 0, \text{fiction} = 5, \text{terror} = 3)$$

cost vector:

representation of **personal preference** in many applications

# Setup: Example-Dependent Cost-Sensitive Classification

## Given

$N$  examples, each (input  $\mathbf{x}_n$ , label  $y_n$ )  $\in \mathcal{X} \times \{1, 2, \dots, K\}$

and cost vector  $\mathbf{c}_n \in \mathbb{R}^K$

—will assume  $\mathbf{c}_n[y_n] = 0 = \min_{1 \leq k \leq K} \mathbf{c}_n[k]$

## Goal

a classifier  $g(\mathbf{x})$  that pays a small cost  $\mathbf{c}[g(\mathbf{x})]$  on future **unseen** example  $(\mathbf{x}, y, \mathbf{c})$

- will assume  $\mathbf{c}[y] = 0 = c_{\min} = \min_{1 \leq k \leq K} \mathbf{c}[k]$
- note:  $y$  not really needed in evaluation

**example-dependent**  $\supset$  class-dependent  $\supset$  regular

# Outline

## 1 Cost-Sensitive Multiclass Classification

- CSMC Motivation and Setup
- **CSMC by Bayesian Perspective**
- CSMC by (Weighted) Binary Classification
- CSMC by Regression

## 2 Cost-Sensitive Multilabel Classification

- CSML Motivation and Setup
- CSML by Bayesian Perspective
- CSML by (CS) Multiclass Classification
- CSML by (Weighted) Binary Classification
- CSML by Regression

## 3 CSMC & CSML: Selected Applications

- Bacteria Classification with Doctor-Annotated Costs
- Network Connection Classification with 'Danger' Costs on Imbalanced Data
- Social Tagging with Costs from Tag Counts

## 4 Summary

# Key Idea: Conditional Probability Estimator

## Goal (Class-Dependent Setup)

a classifier  $g(\mathbf{x})$  that pays a small cost  $\mathcal{C}(y, g(\mathbf{x}))$  on future **unseen** example  $(\mathbf{x}, y)$

if  $P(y|\mathbf{x})$  known

Bayes optimal  $g^*(\mathbf{x}) =$

$$\operatorname{argmin}_{1 \leq k \leq K} \sum_{y=1}^K P(y|\mathbf{x}) \mathcal{C}(y, k)$$

if  $q(\mathbf{x}, y) \approx P(y|\mathbf{x})$  well

approximately good  $g_q(\mathbf{x}) =$

$$\operatorname{argmin}_{1 \leq k \leq K} \sum_{y=1}^K q(\mathbf{x}, y) \mathcal{C}(y, k)$$

how to get **conditional probability estimator**  $q$ ?  
**logistic regression, Naïve Bayes, ...**

# Approximate Bayes-Optimal Decision

if  $q(\mathbf{x}, y) \approx P(y|\mathbf{x})$  well

(Domingos, 1999)

$$\text{approximately good } g_q(\mathbf{x}) = \operatorname{argmin}_{1 \leq k \leq K} \sum_{y=1}^K q(\mathbf{x}, y) \mathcal{C}(y, k)$$

## Approximate Bayes-Optimal Decision (ABOD) Approach

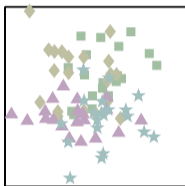
- 1 use your favorite algorithm on  $\{(\mathbf{x}_n, y_n)\}$  to get  $q(\mathbf{x}, y) \approx P(y|\mathbf{x})$
- 2 for each new input  $\mathbf{x}$ , **predict its class using  $g_q(\mathbf{x})$**  above

**ABOD: probability estimate + Bayes-optimal decision**

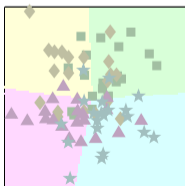
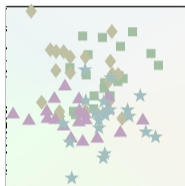


## ABOD on Artificial Data

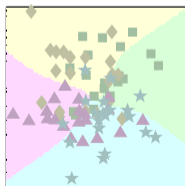
- 1 use your favorite algorithm on  $\{(\mathbf{x}_n, y_n)\}$  to get  $q(\mathbf{x}, y) \approx P(y|\mathbf{x})$
- 2 for each new input  $\mathbf{x}$ , **predict its class using  $g_q(\mathbf{x})$**  above



LogReg  
→



regular



rotate

		$g(\mathbf{x})$			
		1	2	3	4
$y$	1	0	1	2	4
	2	4	0	1	2
	3	2	4	0	1
	4	1	2	4	0

# ABOD for Binary Classification

Given  $N$  examples, each (input  $\mathbf{x}_n$ , label  $y_n$ )  $\in \mathcal{X} \times \{-1, +1\}$   
 and weights  $w_+$ ,  $w_-$  representing **two entries of cost matrix**

		$g(\mathbf{x})$	
		+1	-1
$y$	+1	0	$w_+$
	-1	$w_-$	0

if  $q(\mathbf{x}) \approx P(+1|\mathbf{x})$  well

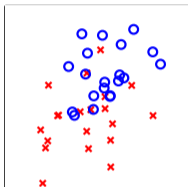
approximately good  $g_q(\mathbf{x}) = \text{sign}\left(w_+q(\mathbf{x}) - w_-(1 - q(\mathbf{x}))\right)$ , i.e. (Elkan, 2001),

$$g_q(\mathbf{x}) = +1 \iff w_+q(\mathbf{x}) - w_-(1 - q(\mathbf{x})) > 0 \iff q(\mathbf{x}) > \frac{w_-}{w_+ + w_-}$$

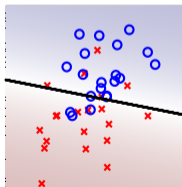
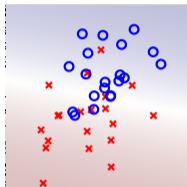
ABOD for binary classification:  
**probability estimate** + **threshold changing**

# ABOD for Binary Classification on Artificial Data

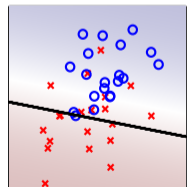
- 1 use your favorite algorithm on  $\{(\mathbf{x}_n, y_n)\}$  to get  $q(\mathbf{x}) \approx P(+1|\mathbf{x})$
- 2 for each new input  $\mathbf{x}$ , predict its class using  $g_q(\mathbf{x}) = \text{sign}(q(\mathbf{x}) - \frac{w_-}{w_+ + w_-})$



LogReg  
→



regular



positive emphasis

		$g(\mathbf{x})$	
		+1	-1
$y$	+1	0	10
	-1	1	0

# Pros and Cons of ABOD

## Pros

- optimal if good **probability estimate**  $q$
- **prediction** easily adapts to different  $\mathcal{C}$  without modifying **training** (probability estimate)

## Cons

- ‘difficult’: good **probability estimate** often more difficult  
than good **multiclass classification**
- ‘restricted’: only applicable to **class-dependent setup**  
—need ‘full picture’ of cost matrix
- ‘slower prediction’ (for multiclass): **more calculation** at prediction stage

can we use any **multiclass classification algorithm** for ABOD?

# MetaCost Approach

## Approximate Bayes-Optimal Decision (ABOD) Approach

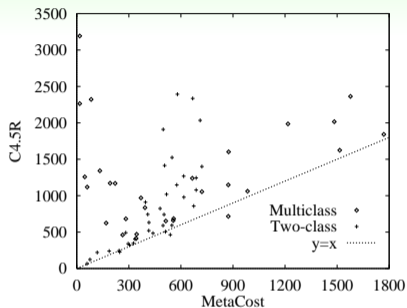
- 1 use your favorite algorithm on  $\{(\mathbf{x}_n, y_n)\}$  to get  $q(y, \mathbf{x}) \approx P(y|\mathbf{x})$
- 2 for each new input  $\mathbf{x}$ , predict its class using  $g_p(\mathbf{x})$

## MetaCost Approach (Domingos, 1999)

- 1 use your favorite **multiclass classification** algorithm on **bootstrapped**  $\{(\mathbf{x}_n, y_n)\}$  **and aggregate** the classifiers to get  $q(y, \mathbf{x}) \approx P(y|\mathbf{x})$
- 2 for each given input  $\mathbf{x}_n$ , **relabel it to**  $y'_n$  using  $g_q(\mathbf{x})$
- 3 run your favorite **multiclass classification** algorithm on **relabeled**  $\{(\mathbf{x}_n, y'_n)\}$  to get final classifier  $g$
- 4 for each new input  $\mathbf{x}$ , predict its class using  $g(\mathbf{x})$

pros: any **multiclass classification** algorithm can be used

# MetaCost on Semi-Real Data



(Domingos, 1999)

- some 'artificial' cost with UCI data
- MetaCost+C4.5: **cost-sensitive**
- C4.5: regular

not surprisingly,

**considering the cost properly does help**

# Outline

## 1 Cost-Sensitive Multiclass Classification

- CSMC Motivation and Setup
- CSMC by Bayesian Perspective
- **CSMC by (Weighted) Binary Classification**
- CSMC by Regression

## 2 Cost-Sensitive Multilabel Classification

- CSML Motivation and Setup
- CSML by Bayesian Perspective
- CSML by (CS) Multiclass Classification
- CSML by (Weighted) Binary Classification
- CSML by Regression

## 3 CSMC & CSML: Selected Applications

- Bacteria Classification with Doctor-Annotated Costs
- Network Connection Classification with 'Danger' Costs on Imbalanced Data
- Social Tagging with Costs from Tag Counts

## 4 Summary

## Key Idea: Cost Transformation

(heuristic) **relabeling** useful in MetaCost: a more **principled** way?

Yes, by Connecting Cost Vector to Regular Costs!

$$\underbrace{(1 \ 0 \ 1 \ 2)}_{\mathbf{c} \text{ of interest}} \xrightarrow{\text{shift equivalence}} \underbrace{(3 \ 2 \ 3 \ 4)}_{\text{shifted cost}} = \underbrace{(1 \ 2 \ 1 \ 0)}_{\text{mixture weights } u_\ell} \cdot \underbrace{\begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}}_{\text{regular costs}}$$

i.e.  $\mathbf{x}$  with  $\mathbf{c} = (1, 0, 1, 2)$  equivalent to

a weighted mixture  $\{(\mathbf{x}, y, u)\} = \{(\mathbf{x}, 1, 1), (\mathbf{x}, 2, 2), (\mathbf{x}, 3, 1)\}$

**cost equivalence** (Lin, 2014): for any classifier  $h$ ,

$$\mathbf{c}[h(\mathbf{x})] + \text{constant} = \sum_{\ell=1}^K u_\ell \mathbb{I}[\ell \neq h(\mathbf{x})]$$



# Meaning of Cost Equivalence

$$\mathbf{c}[h(\mathbf{x})] + \text{constant} = \sum_{\ell=1}^K u_{\ell} \mathbb{I}[\ell \neq h(\mathbf{x})]$$

on one  $(\mathbf{x}, y, \mathbf{c})$ :

wrong prediction charged by  
 $\mathbf{c}[h(\mathbf{x})]$

on all  $\{(\mathbf{x}, \ell, u_{\ell})\}$ :

wrong prediction charged by  
**total weighted classification error**  
of **relabeled data**

$$\begin{aligned} & \min_h \text{expected LHS} && \text{(original CSMC problem)} \\ = & \min_h \text{expected RHS} && \text{(weighted classification when } u_{\ell} \geq 0) \end{aligned}$$

Calculation of  $u_\ell$ Smallest Non-Negative  $u_\ell$ 's (Lin, 2014)

when constant =  $(K - 1) \max_{1 \leq k \leq K} \mathbf{c}[k] - \sum_{k=1}^K \mathbf{c}[k]$ ,

$$u_\ell = \max_{1 \leq k \leq K} \mathbf{c}[k] - \mathbf{c}[\ell]$$

e.g.  $\underbrace{(1 \ 0 \ 1 \ 2)}_{\mathbf{c} \text{ of interest}} \rightarrow \underbrace{(1 \ 2 \ 1 \ 0)}_{\text{mixture weights } u_\ell}$

- **largest  $\mathbf{c}[\ell]$** :  $u_\ell = 0$  (**least preferred relabel**)
- **smallest  $\mathbf{c}[\ell]$** :  $u_\ell = \text{largest}$  (**original label & most preferred relabel**)

$\ell$ 's and  $u_\ell$ 's **embed the cost**

# Data Space Expansion Approach

## Data Space Expansion (DSE) Approach (Abe, 2004)

- 1 for each  $(\mathbf{x}_n, y_n, \mathbf{c}_n)$  and  $\ell$ , let  $u_{n,\ell} = \max_{1 \leq k \leq K} \mathbf{c}_n[k] - \mathbf{c}_n[\ell]$
- 2 apply your favorite **multiclass classification algorithm** on the weighted mixtures  $\bigcup_{n=1}^N \{(\mathbf{x}_n, \ell, u_{n,\ell})\}_{\ell=1}^K$  to get  $g(\mathbf{x})$
- 3 for each new input  $\mathbf{x}$ , predict its class using  $g(\mathbf{x})$

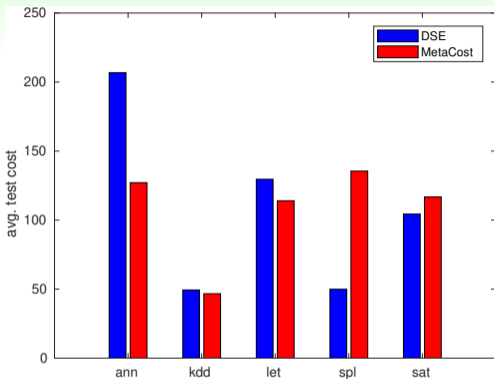
- by **cost equivalence**,

$$= \begin{array}{l} \text{good } g \text{ for } \text{new (weighted) regular classification problem} \\ \text{good } g \text{ for } \text{original cost-sensitive classification problem} \end{array}$$

- **weighted** regular classification: special case of CSMC  
but more easily solvable by, e.g., **sampling** + regular classification (Zadrozny, 2003)

pros: any **multiclass classification** algorithm can be used

# DSE versus MetaCost on Semi-Real Data



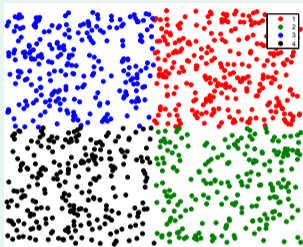
(Abe, 2004) some 'artificial' cost with UCI data

- use **sampling** + C4.5 for weighted regular classification

**DSE competitive to MetaCost**

# Cons of DSE: Unavoidable Noise

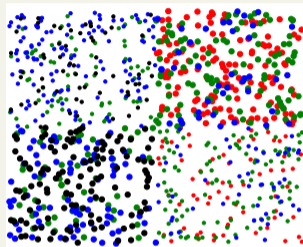
## Original Cost-Sensitive Classification Problem



individual examples  
without noise

+ absolute  
cost =

## New Regular Classification Problem



mixtures with  
relabeling noise

- cost embedded as weight + **noisy labels**
- new problem usually **harder** than original one

**need robust multiclass classification algorithm  
to deal with noise**

# Key Idea: Design Robust Multiclass Algorithm

## One-Versus-One: A Popular Classification Meta-Method

- for all different class pairs  $(i, j)$ ,
  - 1 take all examples  $(\mathbf{x}_n, y_n)$ 
    - that  $y_n = i$  or  $j$  (**original one-versus-one**)
    - that  $u_{n,i} \neq u_{n,j}$  with the larger- $u$  label and weight  $|u_{n,i} - u_{n,j}|$  (**robust one-versus-one**)
  - 2 train a binary classifier  $\hat{g}^{(i,j)}$  using those examples
- return  $g(\mathbf{x})$  that predicts using the votes from  $\hat{g}^{(i,j)}$
- **un-shifting** inside the meta-method to **remove noise**
- **robust step** makes it suitable for DSE

**cost-sensitive one-versus-one: DSE + robust one-versus-one**

# Cost-Sensitive One-Versus-One (CSOVO)

## Cost-Sensitive One-Versus-One (Lin, 2014)

- for all different class pairs  $(i, j)$ ,
  - 1 **robust one-versus-one** + calculate from  $\mathbf{c}_n$ : take all examples  $(\mathbf{x}_n, y_n)$  that  $\mathbf{c}_n[i] \neq \mathbf{c}_n[j]$  with **smaller-c** label and **weight**  $u_n^{(i,j)} = |\mathbf{c}_n[i] - \mathbf{c}_n[j]|$
  - 2 train a binary classifier  $\hat{g}^{(i,j)}$  using those examples
- return  $g(\mathbf{x})$  that predicts using the votes from  $\hat{g}^{(i,j)}$

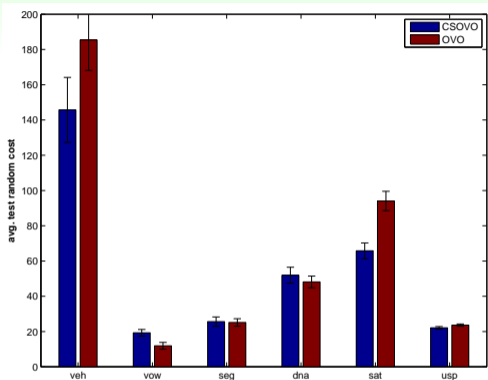
- comes with **good theoretical guarantee**:

$$\text{test cost of } g \leq 2 \sum_{i < j} \text{test cost of } \hat{g}^{(i,j)}$$

- sibling to Weighted All-Pairs (WAP) approach: even tighter guarantee (Beygelzimer, 2005) with **more sophisticated construction** of  $u_n^{(i,j)}$

physical meaning: each  $\hat{g}^{(i,j)}$  answers yes/no question 'prefer  $i$  or  $j$ '?

## CSOVO on Semi-Real Data



(Lin, 2014) some 'artificial' cost with UCI data

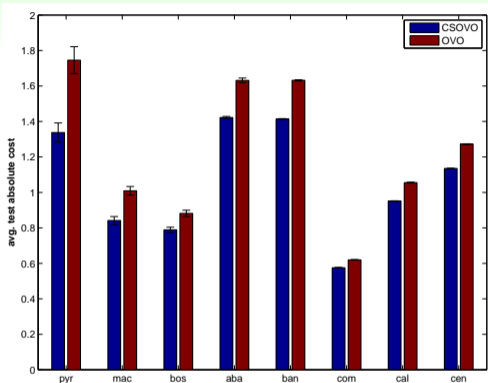
- CSOVO-SVM: cost-sensitive
- OVO-SVM: regular

not surprisingly **again**,

**considering the cost properly does help**



## CSOVO for Ordinal Ranking



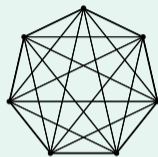
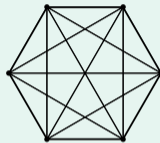
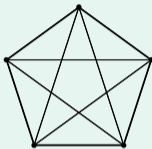
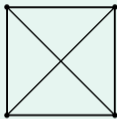
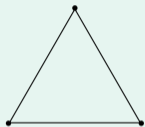
(Lin, 2014) absolute cost with benchmark ordinal ranking data

- CSOVO-SVM: cost-sensitive
- OVO-SVM: regular

**CSOVO significantly better for ordinal ranking**

# Cons of CSOVO: Many Binary Classifiers

$K$  classes  $\xrightarrow{\text{CSOVO}}$   $\frac{K(K-1)}{2}$  binary classifiers



**time-consuming** in both

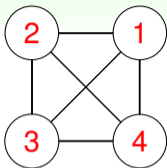
- training, especially with **many different  $c_n[i]$  and  $c_n[j]$**
- prediction

—parallization helps a bit, but **generally not feasible for large  $K$**

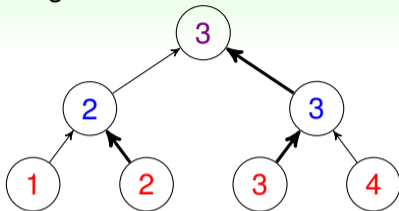
**CSOVO**: a simple meta-method **for median  $K$  only**

Key Idea: OVO  $\equiv$  Round-Robin Tournament

Round-Robin Tournament



Single-Elimination Tournament

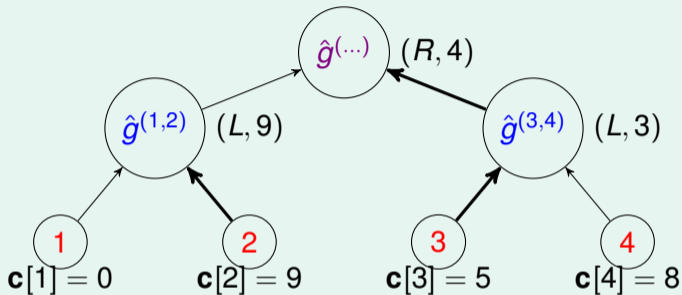


- prediction  $\equiv$  deciding **tournament winner** for each  $\mathbf{x}$
- (CS)OVO:  $\frac{K(K-1)}{2}$  games for prediction (and hence training)
- **single-elimination tournament** (for  $K = 2^\ell$ ):
  - $K - 1$  games for prediction via bottom-up: real-world
  - $\log_2 K$  **games** for prediction via top-down: **computer-world :-)**

next: **single-elimination tournament** for CSMC

## Filter Tree (FT) Approach

## Filter Tree (Beygelzimer, 2009) Training: from bottom to top



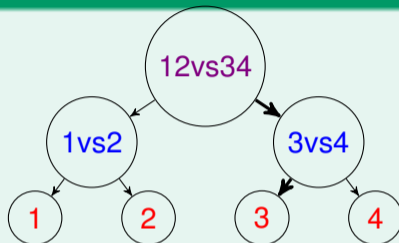
- $\hat{g}^{(1,2)}$  and  $\hat{g}^{(3,4)}$  trained like CSOVO: **smaller- $\mathbf{c}$**  label and **weight**  $u_n^{(i,j)} = |\mathbf{c}_n[i] - \mathbf{c}_n[j]|$
- $\hat{g}^{(\dots)}$  trained with  $(k_L, k_R)$  filtered by sub-trees  
 — **smaller- $\mathbf{c}$**  sub-tree direction and **weight**  $u_n^{(\dots)} = |\mathbf{c}_n[k_L] - \mathbf{c}_n[k_R]|$

FT: top classifiers **aware of bottom-classifier mistakes**

# Pros and Cons of FT

## Pros

- efficient:  
 $O(K)$  training,  $O(\log K)$  prediction
- **strong theoretical guarantee:**  
small-**regret** binary classifiers  
 $\implies$  small-**regret** CSMC classifier



## Cons

- 'asymmetric' to labels: **non-trivial structural decision**
- 'hard' **sub-tree dependent** top-classification tasks

next: **other reductions** to (weighted) binary classification

# Other Approaches via Weighted Binary Classification

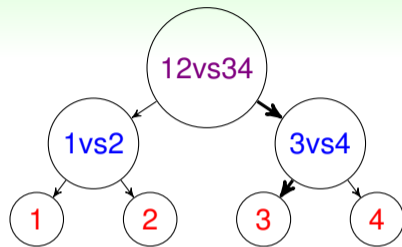
**FT:** with regret bound (Beygelzimer, 2009)

the lowest **achievable** cost within  $\{1, 2\}$  or  $\{3, 4\}$ ?

**Divide&Conquer Tree (TREE):**

**without** regret bound (Beygelzimer, 2009)

the lowest **ideal** cost within  $\{1, 2\}$  or  $\{3, 4\}$ ?



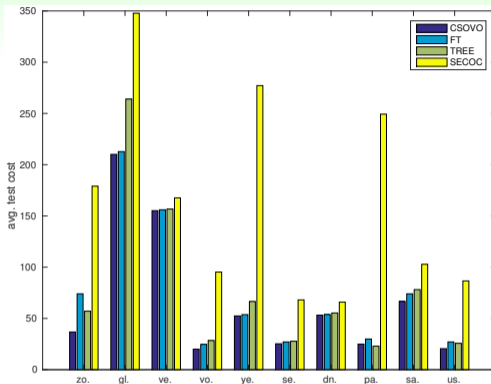
**Sensitive Err. Correcting Output Code (SECOC):** with regret bound (Langford, 2005)

$\mathbf{c}[1] + \mathbf{c}[3] + \mathbf{c}[4]$  greater than some  $\theta$ ?

training time:

$\text{SECOC } (O(T \cdot K)) > \text{FT } (O(K)) \approx \text{TREE } (O(K))$

# Comparison of Reductions to Weighted Binary Classification



(Lin, 2014) couple all meta-methods with SVM

- round-robin tournament (CSOVO)
- single-elimination tournament (FT, TREE)
- error-correcting-code (SECOC)

**CSOVO often among the best;  
FT somewhat competitive**

# Outline

## 1 Cost-Sensitive Multiclass Classification

- CSMC Motivation and Setup
- CSMC by Bayesian Perspective
- CSMC by (Weighted) Binary Classification
- **CSMC by Regression**

## 2 Cost-Sensitive Multilabel Classification

- CSML Motivation and Setup
- CSML by Bayesian Perspective
- CSML by (CS) Multiclass Classification
- CSML by (Weighted) Binary Classification
- CSML by Regression

## 3 CSMC & CSML: Selected Applications

- Bacteria Classification with Doctor-Annotated Costs
- Network Connection Classification with 'Danger' Costs on Imbalanced Data
- Social Tagging with Costs from Tag Counts

## 4 Summary



# Key Idea: Cost Estimator

## Goal

a classifier  $g(\mathbf{x})$  that pays a small cost  $\mathbf{c}[g(\mathbf{x})]$  on future **unseen** example  $(\mathbf{x}, y, \mathbf{c})$

if every  $\mathbf{c}[k]$  known

optimal

$$g^*(\mathbf{x}) = \operatorname{argmin}_{1 \leq k \leq K} \mathbf{c}[k]$$

if  $r_k(\mathbf{x}) \approx \mathbf{c}[k]$  well

approximately good

$$g_r(\mathbf{x}) = \operatorname{argmin}_{1 \leq k \leq K} r_k(\mathbf{x})$$

how to get cost estimator  $r_k$ ? **regression**

# Cost Estimator by Per-class Regression

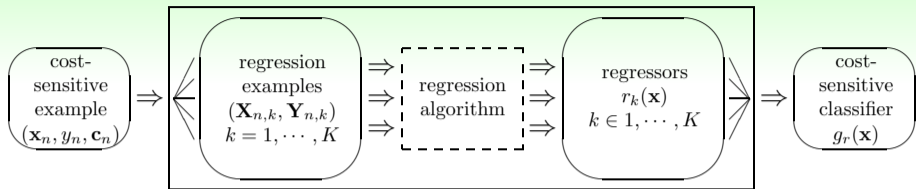
## Given

$N$  examples, each (input  $\mathbf{x}_n$ , label  $y_n$ , cost  $\mathbf{c}_n$ )  $\in \mathcal{X} \times \{1, 2, \dots, K\} \times R^K$

input	$\mathbf{c}_n[1]$	input	$\mathbf{c}_n[2]$	...	input	$\mathbf{c}_n[K]$
$\mathbf{x}_1$	0,	$\mathbf{x}_1$	2,		$\mathbf{x}_1$	1
$\mathbf{x}_2$	1,	$\mathbf{x}_2$	3,		$\mathbf{x}_2$	5
...						
$\mathbf{x}_N$	6,	$\mathbf{x}_N$	1,		$\mathbf{x}_N$	0
$\underbrace{\hspace{10em}}$		$\underbrace{\hspace{10em}}$			$\underbrace{\hspace{10em}}$	
$r_1$		$r_2$			$r_K$	

**want:**  $r_k(\mathbf{x}) \approx \mathbf{c}[k]$  for all future  $(\mathbf{x}, y, \mathbf{c})$  and  $k$

# The Reduction-to-Regression Framework



- 1 **encode**: transform cost-sensitive examples  $(\mathbf{x}_n, y_n, \mathbf{c}_n)$  to regression examples  $(\mathbf{x}_{n,k}, Y_{n,k}) = (\mathbf{x}_n, \mathbf{c}_n[k])$
- 2 **learn**: use your favorite algorithm on regression examples to get estimators  $r_k(\mathbf{x})$
- 3 **decode**: for each new input  $\mathbf{x}$ , predict its class using  $g_r(\mathbf{x}) = \operatorname{argmin}_{1 \leq k \leq K} r_k(\mathbf{x})$

the reduction-to-regression framework:

**systematic & easy to implement**

# Theoretical Guarantees (1/2)

$$g_r(\mathbf{x}) = \operatorname{argmin}_{1 \leq k \leq K} r_k(\mathbf{x})$$

## Theorem (Absolute Loss Bound)

For any set of cost estimators  $\{r_k\}_{k=1}^K$  and for any example  $(\mathbf{x}, y, \mathbf{c})$  with  $\mathbf{c}[y] = 0$ ,

$$\mathbf{c}[g_r(\mathbf{x})] \leq \sum_{k=1}^K |r_k(\mathbf{x}) - \mathbf{c}[k]|.$$

**low-cost classifier  $\Leftarrow$  accurate estimators**

## Theoretical Guarantees (2/2)

$$g_r(\mathbf{x}) = \operatorname{argmin}_{1 \leq k \leq K} r_k(\mathbf{x})$$

### Theorem (Squared Loss Bound)

For any set of cost estimators  $\{r_k\}_{k=1}^K$  and for any example  $(\mathbf{x}, y, \mathbf{c})$  with  $\mathbf{c}[y] = 0$ ,

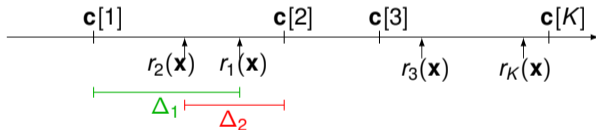
$$\mathbf{c}[g_r(\mathbf{x})] \leq \sqrt{2 \sum_{k=1}^K (r_k(\mathbf{x}) - \mathbf{c}[k])^2}.$$

applies to common **least-square regression**

## A Pictorial Proof

$$\mathbf{c}[g_r(\mathbf{x})] \leq \sum_{k=1}^K |r_k(\mathbf{x}) - \mathbf{c}[k]|$$

- assume  $\mathbf{c}$  ordered and not degenerate:  $y = 1$ ;  $0 = \mathbf{c}[1] < \mathbf{c}[2] \leq \dots \leq \mathbf{c}[K]$
- assume mis-prediction  $g_r(\mathbf{x}) = 2$ :  $r_2(\mathbf{x}) = \min_{1 \leq k \leq K} r_k(\mathbf{x}) \leq r_1(\mathbf{x})$

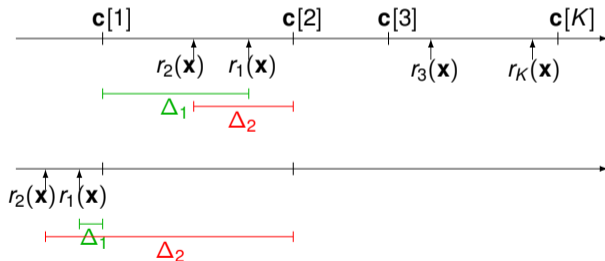


$$\mathbf{c}[2] - \underbrace{\mathbf{c}[1]}_0 \leq |\Delta_1| + |\Delta_2| \leq \sum_{k=1}^K |r_k(\mathbf{x}) - \mathbf{c}[k]|$$

## An Even Closer Look

let  $\Delta_1 \equiv r_1(\mathbf{x}) - \mathbf{c}[1]$  and  $\Delta_2 \equiv \mathbf{c}[2] - r_2(\mathbf{x})$

- 1  $\Delta_1 \geq 0$  and  $\Delta_2 \geq 0$ :  $\mathbf{c}[2] \leq \Delta_1 + \Delta_2$
- 2  $\Delta_1 \leq 0$  and  $\Delta_2 \geq 0$ :  $\mathbf{c}[2] \leq \Delta_2$
- 3  $\Delta_1 \geq 0$  and  $\Delta_2 \leq 0$ :  $\mathbf{c}[2] \leq \Delta_1$



$$\mathbf{c}[2] \leq \max(\Delta_1, 0) + \max(\Delta_2, 0) \leq |\Delta_1| + |\Delta_2|$$

## Tighter Bound with One-sided Loss

Define **one-sided loss**  $\xi_k \equiv \max(\Delta_k, 0)$

with  $\Delta_k \equiv (r_k(\mathbf{x}) - \mathbf{c}[k])$  if  $\mathbf{c}[k] = c_{\min} = 0$

$\Delta_k \equiv (\mathbf{c}[k] - r_k(\mathbf{x}))$  if  $\mathbf{c}[k] \neq c_{\min}$

## Intuition

- $\mathbf{c}[k] = c_{\min}$ : wish to have  $r_k(\mathbf{x}) \leq \mathbf{c}[k]$
- $\mathbf{c}[k] \neq c_{\min}$ : wish to have  $r_k(\mathbf{x}) \geq \mathbf{c}[k]$

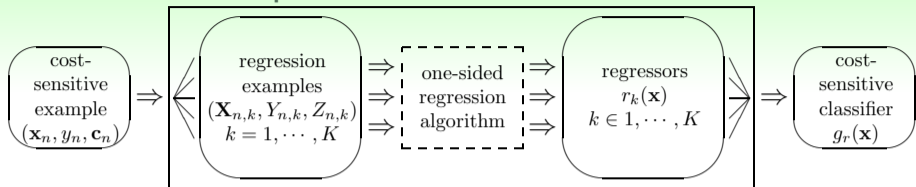
—both wishes same as  $\Delta_k \leq 0 \iff \xi_k = 0$

One-sided Loss Bound:

$$\mathbf{c}[g_r(\mathbf{x})] \leq \sum_{k=1}^K \xi_k \leq \sum_{k=1}^K |\Delta_k|$$



## The Improved Reduction Framework



(Tu, 2010)

- 1 **encode**: transform cost-sensitive examples  $(\mathbf{x}_n, y_n, \mathbf{c}_n)$  to **one-sided regression examples**  $(\mathbf{x}_n^{(k)}, Y_n^{(k)}, Z_n^{(k)}) = (\mathbf{x}_n, \mathbf{c}_n[k], 2 \llbracket \mathbf{c}_n[k] = 0 \rrbracket - 1)$
- 2 **learn**: use a **one-sided regression algorithm** to get estimators  $r_k(\mathbf{x})$
- 3 **decode**: for each new input  $\mathbf{x}$ , predict its class using  $g_r(\mathbf{x}) = \operatorname{argmin}_{1 \leq k \leq K} r_k(\mathbf{x})$

the reduction-to-OSR framework:

**need a good OSR algorithm**

## Regularized One-Sided Hyper-Linear Regression

Given

$$(\mathbf{x}_{n,k}, Y_{n,k}, Z_{n,k}) = (\mathbf{x}_n, \mathbf{c}_n[k], 2 \llbracket \mathbf{c}_n[k] = \mathbf{0} \rrbracket - 1)$$

Training Goal

all training  $\xi_{n,k} = \max \left( \underbrace{Z_{n,k} (r_k(\mathbf{x}_{n,k}) - Y_{n,k})}_{\Delta_{n,k}}, 0 \right)$  small

—will drop  $k$ 

$$\min_{\mathbf{w}, b} \quad \frac{\lambda}{2} \langle \mathbf{w}, \mathbf{w} \rangle + \sum_{n=1}^N \xi_n$$

to get  $r_k(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b$

## One-Sided Support Vector Regression

## Regularized One-Sided Hyper-Linear Regression

$$\min_{\mathbf{w}, b} \quad \frac{\lambda}{2} \langle \mathbf{w}, \mathbf{w} \rangle + \sum_{n=1}^N \xi_n$$

$$\xi_n = \max(Z_n \cdot (r_k(\mathbf{x}_n) - Y_n), 0)$$

## Standard Support Vector Regression

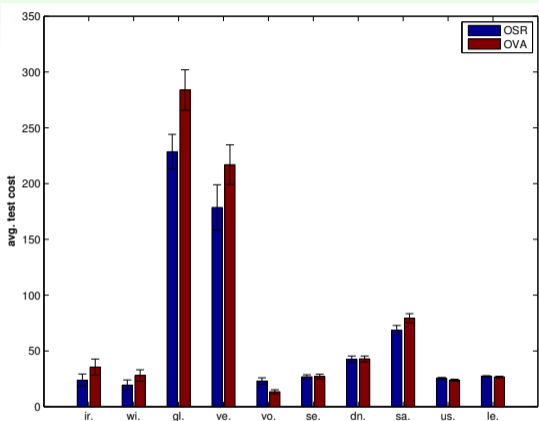
$$\min_{\mathbf{w}, b} \quad \frac{1}{2C} \langle \mathbf{w}, \mathbf{w} \rangle + \sum_{n=1}^N (\xi_n + \xi_n^*)$$

$$\xi_n = \max(+1 \cdot (r_k(\mathbf{x}_n) - Y_n - \epsilon), 0)$$

$$\xi_n^* = \max(-1 \cdot (r_k(\mathbf{x}_n) - Y_n + \epsilon), 0)$$

**OSR-SVM** = SVR + ( $\epsilon \leftarrow 0$ ) + (keep  $\xi_n$  or  $\xi_n^*$  by  $Z_n$ )

## OSR-SVM on Semi-Real Data

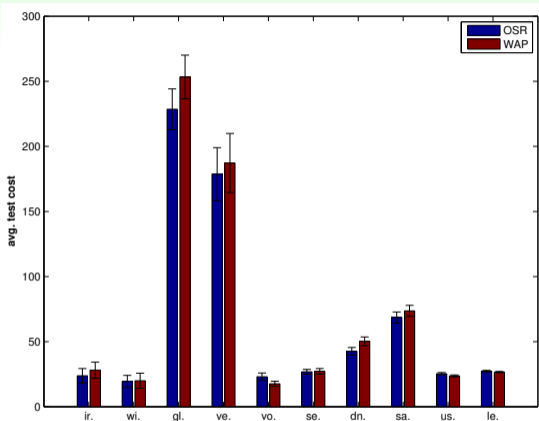


(Tu, 2010) some 'artificial' cost with UCI data

- OSR: cost-sensitive SVM
- OVA: regular one-versus-all SVM

**OSR often significantly better than OVA**

## OSR versus WAP on Semi-Real Data



(Tu, 2010) some 'artificial' cost with UCI data

- **OSR** (per-class):  
 $O(K)$  training,  $O(K)$  prediction
- **WAP**  $\approx$  **CSOVO** (pairwise):  $O(K^2)$  training,  $O(K^2)$  prediction

**OSR faster and competitive performance**

## From OSR-SVM to AOSR-DNN

$$\begin{array}{ll}
 \text{OSR-SVM} & \min_{\mathbf{w}, b} \quad \frac{\lambda}{2} \langle \mathbf{w}, \mathbf{w} \rangle + \sum_{n=1}^N \xi_n \\
 & \text{with} \quad r_k(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b \\
 & \quad \xi_n = \max(Z_n \cdot (r_k(\mathbf{x}_n) - Y_n), 0)
 \end{array}$$

$$\begin{array}{ll}
 \text{Appro. OSR-DNN} & \min_{\text{NNet}} \quad \text{regularizer} + \sum_{n=1}^N \delta_n \\
 & \text{with} \quad r_k(\mathbf{x}) = \text{NNet}(\mathbf{x}) \\
 & \quad \delta_n = \ln(1 + \exp(Z_n \cdot (r_k(\mathbf{x}_n) - Y_n)))
 \end{array}$$

AOSR-DNN (Chung, 2016a) = Deep Learning + OSR +

**smoother upper bound**  $\delta_n \geq \xi_n$  because  $\ln(1 + \exp(\bullet)) \geq \max(\bullet, 0)$

## From AOSR-DNN to CSDNN

## Cons of AOSR-DNN

**c** affects both classification **and feature-extraction** in DNN  
but hard to do effective **cost-sensitive feature extraction**

idea 1: pre-training with **c**

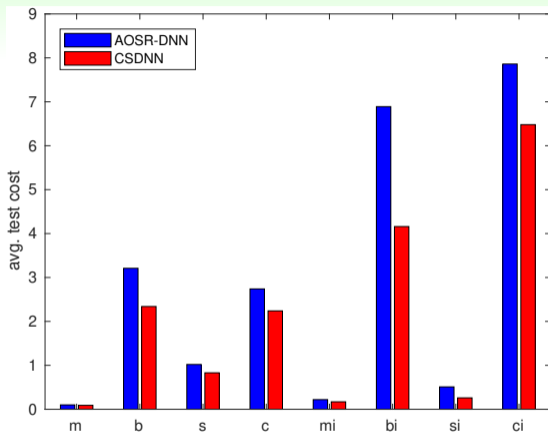
- layer-wise pre-training with **cost-sensitive autoencoders**  
loss = reconstruction + **AOSR**
- CSDNN (Chung, 2016a)  
= AOSR-DNN + cost-sens.  
pre-training

## idea 2: auxiliary cost-sensitive nodes

- **auxiliary nodes** to  
predict costs per layer  
loss = AOSR for classification  
+ **AOSR for auxiliary**
- applies to any deep learning model  
(Chung, 2016b)

CSDNN: world's **first successful CSMC deep learning model**

## AOSR-DNN versus CSDNN



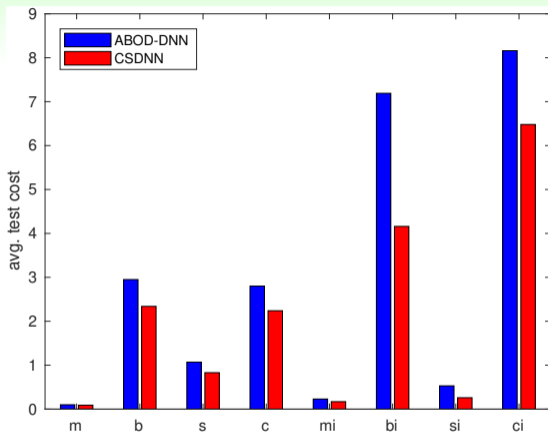
(Chung, 2016a)

- **AOSR-DNN**: cost-sensitive training
- **CSDNN**: **AOSR-DNN** + cost-sensitive feature extraction

**CSDNN wins**, justifying **cost-sensitive feature extraction**



## ABOD-DNN versus CSDNN



(Chung, 2016a)

- **ABOD-DNN:**  
probability estimate + cost-sensitive prediction
- **CSDNN:**  
cost-sensitive training + cost-sensitive feature extraction

**CSDNN still wins**, hinting **difficulty of probability estimate** without **cost-sensitive feature extraction**

# Outline

## 1 Cost-Sensitive Multiclass Classification

- CSMC Motivation and Setup
- CSMC by Bayesian Perspective
- CSMC by (Weighted) Binary Classification
- CSMC by Regression

## 2 Cost-Sensitive Multilabel Classification

- **CSML Motivation and Setup**
- CSML by Bayesian Perspective
- CSML by (CS) Multiclass Classification
- CSML by (Weighted) Binary Classification
- CSML by Regression

## 3 CSMC & CSML: Selected Applications

- Bacteria Classification with Doctor-Annotated Costs
- Network Connection Classification with 'Danger' Costs on Imbalanced Data
- Social Tagging with Costs from Tag Counts

## 4 Summary

# Which Fruit?



?

(image by Robert-Owen-Wahl from Pixabay)



apple



orange



strawberry



kiwi

(images by Pexels, PublicDomainPictures, 192635, Rob van der Meijden from Pixabay)

multiclass classification:  
classify input (picture) to **one category** (label), **remember? :-)**

# Which Fruits?



?: {apple, orange, kiwi}

(image by Michal Jarmoluk from Pixabay)



apple



orange



strawberry



kiwi

(images by Pexels, PublicDomainPictures, 192635, Rob van der Meijden from Pixabay)

**multilabel** classification:  
classify input to **multiple (or no)** categories

# Label Powerset: Multilabel Classification via Multiclass (Tsoumakas, 2007)

multiclass w/  $L = 4$   
classes

4 possible outcomes  
{a, o, s, k}

multilabel w/  $L = 4$  classes

$2^4 = 16$  possible outcomes  
 $2^{\{a, o, s, k\}}$

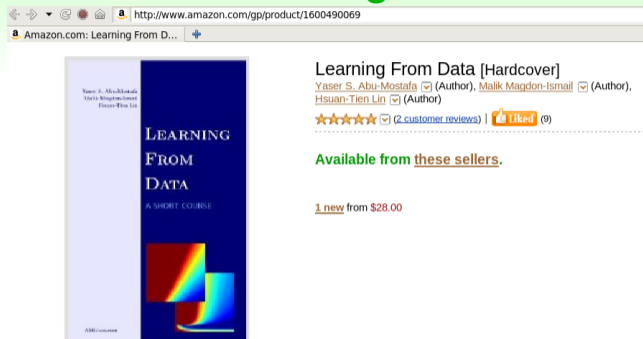


{  $\phi$ , a, o, s, k, ao, as, ak, os, ok, sk,  
aos, aok, ask, osk, aosk }

- **Label Powerset (LP):** reduction to multiclass classification
- difficulties for large  $L$ :
  - **computation:**  $2^L$  extended classes
  - **sparsity:** no or few example for some combinations

**LP:** feasible only for **small  $L$**

# What Tags?



http://www.amazon.com/gp/product/1600490069

Amazon.com: Learning From D...

Learning From Data [Hardcover]  
 Yaser S. Abu-Mostafa (Author), Malik Magdon-Ismael (Author),  
 Hsuan-Tien Lin (Author)

★★★★★ (2 customer reviews) | Liked (9)

Available from [these sellers](#).

1 new from \$28.00

?: { machine learning, data-structure, data mining, object oriented-programming, artificial intelligence, compiler, architecture, chemistry, textbook, children book, ... etc. }

another **multilabel** classification problem:  
**tagging** input to multiple categories

# Binary Relevance: Multilabel Classification via Yes/No

binary  
classification

{yes, no}

multilabel w/  $L$  classes:  $L$  yes/no  
questions

machine learning (Y), data structure (N), data mining (Y), OOP (N), AI (Y), compiler (N), architecture (N), chemistry (N), textbook (Y), children book (N), *etc.*

- **Binary Relevance (BR)**: reduction to **multiple isolated binary classification**
- disadvantages:
  - **isolation**—hidden relations not exploited (e.g. ML and DM **highly correlated**, ML **subset of** AI, textbook & children book **disjoint**)
  - **unbalanced**—few **yes**, many **no**

**BR**: simple (& strong) benchmark with known disadvantages

# Multilabel Classification Setup

## Given

$N$  examples (input  $\mathbf{x}_n$ , labelset  $\mathcal{Y}_n$ )  $\in \mathcal{X} \times 2^{\{1,2,\dots,L\}}$

- fruits:  $\mathcal{X} = \text{encoding}(\text{pictures})$ ,  $\mathcal{Y}_n \subseteq \{1, 2, \dots, 4\}$
- tags:  $\mathcal{X} = \text{encoding}(\text{merchandise})$ ,  $\mathcal{Y}_n \subseteq \{1, 2, \dots, L\}$

## Goal

a multilabel classifier  $g(\mathbf{x})$  that **closely predicts** the labelset  $\mathcal{Y}$  associated with some **unseen** inputs  $\mathbf{x}$  (by **exploiting hidden relations/combinations between labels**)

**multilabel classification:**  
**hot and important** with many real-world applications



# From Labelset to Coding View

	labelset	apple	orange	strawberry	binary code
	$\mathcal{Y}_1 = \{o\}$	0 (N)	1 (Y)	0 (N)	$\mathbf{y}_1 = [0, 1, 0]$
	$\mathcal{Y}_2 = \{a, o\}$	1 (Y)	1 (Y)	0 (N)	$\mathbf{y}_2 = [1, 1, 0]$
	$\mathcal{Y}_3 = \{o, s\}$	0 (N)	1 (Y)	1 (Y)	$\mathbf{y}_3 = [0, 1, 1]$
	$\mathcal{Y}_4 = \{\}$	0 (N)	0 (N)	0 (N)	$\mathbf{y}_4 = [0, 0, 0]$

(images by PublicDomainPictures, Narin Seandag, GiltonF, nihatyetkin from Pixabay)

subset  $\mathcal{Y}$  of  $2^{\{1,2,\dots,L\}}$   $\iff$  length- $L$  binary code  $\mathbf{y}$

# LP Approach: What Performance Measure?

Goal: a classifier  $g(\mathbf{x})$  that **closely predicts** the labelset  $\mathcal{Y}$  (**code**  $\mathbf{y}$ ) associated w/  $\mathbf{x}$

## LP Approach

- 1 encode**: transform multilabel examples  $(\mathbf{x}_n, \mathbf{y}_n)$  to multiclass examples  $(\mathbf{x}_n, Y_n)$ ,  
 where  $Y_n = \text{binary number of } \mathbf{y}_n$   
 $\mathbf{y} \rightarrow Y \mid [0, 0, 0] \rightarrow 0 \quad [0, 0, 1] \rightarrow 1 \quad [0, 1, 0] \rightarrow 2 \quad [0, 1, 1] \rightarrow 3$
- 2 learn**: use any (regular) algorithm on multiclass examples to get classifier  $\hat{g}(\mathbf{x})$
- 3 decode**: for each new input  $\mathbf{x}$ , predict its code using  
 $g(\mathbf{x}) = \text{binary representation of } \hat{g}(\mathbf{x})$

## Measuring 'Closely Predict'

- **regular** multiclass algorithm: optimizes  $\llbracket Y \neq \hat{g}(\mathbf{x}) \rrbracket$
- LP: correspondingly optimizes  $\llbracket \mathbf{y} \neq g(\mathbf{x}) \rrbracket$ , called **subset 0/1 error**

subset 0/1 error: a **strict** measure for multilabel classification

# BR Approach: What Performance Measure?

Goal: a classifier  $g(\mathbf{x})$  that **closely predicts** the labelset  $\mathcal{Y}$  (**code  $\mathbf{y}$** ) associated w/  $\mathbf{x}$

## BR Approach

- 1 **encode**: transform multilabel examples  $(\mathbf{x}_n, \mathbf{y}_n)$  to binary examples  $(\mathbf{x}_n, \mathbf{y}_n^{[\ell]})$
- 2 **learn**: use any algorithm on binary classification examples to get classifier  $\hat{g}_\ell(\mathbf{x})$
- 3 **decode**: for each new input  $\mathbf{x}$ , predict its code using

$$g(\mathbf{x}) = [\hat{g}_1(\mathbf{x}), \hat{g}_2(\mathbf{x}), \dots, \hat{g}_L(\mathbf{x})]$$

## Measuring 'Closely Predict'

- **regular** binary classification algorithm: optimizes  $\mathbb{I}[\mathbf{y}^{[\ell]} \neq \hat{g}_\ell(\mathbf{x})]$
- BR: correspondingly optimizes  $\frac{1}{L} |g(\mathbf{x}) \Delta \mathcal{Y}|$ , called **Hamming error**

Hamming error: a **simple** measure for multilabel classification

# Evaluating Multilabel Classifiers

## Different Approaches Optimizes Different Measure

- LP: subset 0/1 error
- BR: Hamming error

## Different (Assumed) Dependence Associated with Different Measure

(Dembczyński, 2012)

- strong conditional dependence: subset 0/1 error (need 'joint' optimization)
- no conditional dependence: Hamming error (can use 'marginal' optimization)

## Different Applications Needs Different Measure

- information retrieval: F1 score (harmonic mean of precision & recall)
- tag recommendation: ranking error

Cost-Sensitive Multilabel Classification (CSML):  
**design one approach for 'any' measure, just like CSMC**

# Setup: Cost-Sensitive Multilabel Classification (CSML)

## Given

$N$  examples, each (input  $\mathbf{x}_n$ , code  $\mathbf{y}_n$ )  $\in \mathcal{X} \times \{1, 2, \dots, L\}$

and cost function  $\mathcal{C} \in \mathbb{R}^{2^L \times 2^L}$  with  $\mathcal{C}(\mathbf{y}, \mathbf{y}) = \mathbf{0} = \min_{\mathbf{k} \in \{0,1\}^L} \mathcal{C}(\mathbf{y}, \mathbf{k})$

## Goal

a classifier  $g(\mathbf{x})$  that pays a small cost  $\mathcal{C}(\mathbf{y}, g(\mathbf{x}))$  on future **unseen** example  $(\mathbf{x}, \mathbf{y})$

- called **instance-based** CSML: each instance evaluated separately  
—more complicated to solve other kinds of CSML (Hsieh, 2018)
- possible extension to **example-dependent costs**  $\mathcal{C}_{\mathbf{x}}$  like CSMC

will focus on **'class'-dependent instance-based CSML**

# Outline

## 1 Cost-Sensitive Multiclass Classification

- CSMC Motivation and Setup
- CSMC by Bayesian Perspective
- CSMC by (Weighted) Binary Classification
- CSMC by Regression

## 2 Cost-Sensitive Multilabel Classification

- CSML Motivation and Setup
- **CSML by Bayesian Perspective**
- CSML by (CS) Multiclass Classification
- CSML by (Weighted) Binary Classification
- CSML by Regression

## 3 CSMC & CSML: Selected Applications

- Bacteria Classification with Doctor-Annotated Costs
- Network Connection Classification with 'Danger' Costs on Imbalanced Data
- Social Tagging with Costs from Tag Counts

## 4 Summary

# Approximate Bayes-Optimal Decision Revisited for CSML

if  $q(\mathbf{x}, \mathbf{y}) \approx P(\mathbf{y}|\mathbf{x})$  well

AOBD for CSML

$$\text{approximately good } g_q(\mathbf{x}) = \operatorname{argmin}_{\mathbf{k} \in \{0,1\}^L} \sum_{\mathbf{y} \in \{0,1\}^L} q(\mathbf{x}, \mathbf{y}) \mathcal{C}(\mathbf{y}, \mathbf{k})$$

difficulty of directly using AOBD

- difficulty in **probability estimation**:
- difficulty in **cost calculation**:
- difficulty in **inference**:

$2^L$  outputs per  $\mathbf{x}$  for  $q(\mathbf{x}, \mathbf{y})$   
 $2^L$  possible  $\mathbf{y}$  in  $\sum$  to compute per  $\mathbf{k}$   
 $\operatorname{argmin}$  over  $2^L$  possible candidates  $\mathbf{k}$

**ABOD: even harder** for CSML than CSMC

## Key Idea: Estimate Probability with Decomposition

$$\begin{array}{ccccccc}
 P(\mathbf{y} \mid \mathbf{x}) & = & P(\mathbf{y}[1] \mid \mathbf{x}) & \cdot & P(\mathbf{y}[2] \mid \mathbf{x}, \mathbf{y}[1]) & \cdot & P(\mathbf{y}[3] \mid \mathbf{x}, \mathbf{y}[1], \mathbf{y}[2]) & \cdot & \dots \\
 \Downarrow & & \Downarrow & & \Downarrow & & \Downarrow & & \\
 q(\mathbf{x}, \mathbf{y}) & & q_1(\mathbf{x}) & & q_2(\mathbf{x}, \mathbf{y}[1]) & & q_3(\mathbf{x}, \mathbf{y}[1], \mathbf{y}[2]) & & \dots
 \end{array}$$

- $q_\ell(\mathbf{x}, \mathbf{y}[1], \dots, \mathbf{y}[\ell - 1])$ : estimates  $P(\mathbf{y}[\ell] = 1 \mid \mathbf{x}, \mathbf{y}[1], \dots, (\ell - 1))$   
—learned with your favorite estimation algorithm, such as **logistic regression**
- if each  $q_\ell$  accurate, multiplied  $q$  also accurate

$$q(\mathbf{x}, \mathbf{y}) = \prod_{\ell=1}^L q_\ell^{\mathbf{y}[\ell]} (1 - q_\ell)^{(1 - \mathbf{y}[\ell])}$$

Probabilistic Classifier Chain (PCC) (Dembczyński, 2010):  
estimate  $q_\ell$ 's to conquer difficulty in **probability estimation**



## Key Idea: Sparsify $q(\mathbf{x}, \mathbf{y})$ with Representative $\mathbf{y}$ 's

conjecture:  $P(\mathbf{y}|\mathbf{x})$  usually small for most  $\mathbf{y}$ 's, hence **many  $q(\mathbf{x}, \mathbf{y})$  also small**

### draw typical $\mathbf{y}$ from $q(\mathbf{x}, \mathbf{y})$

(Dembczyński, 2011)

- Monte Carlo sampling from  $q_1, q_2, \dots, q_\ell$  sequentially

### keep most probable $\mathbf{y}[1, \dots, \ell]$

(Kumar, 2013)

- beam search: keep  $B$  most probable predictions

calculate **necessary costs/statistics with representative  $\mathbf{y}$**   
to conquer difficulty in **cost calculation**

# Key Idea: Derive Efficient Inference Rule for Specific $\mathcal{C}$

AOBD for CSML: approximately good  $g_q(\mathbf{x}) = \operatorname{argmin}_{\mathbf{k} \in \{0,1\}^L} \sum_{\mathbf{y} \text{ representative}} q(\mathbf{x}, \mathbf{y}) \mathcal{C}(\mathbf{y}, \mathbf{k})$

—even with representative  $\mathbf{y}$ , still exponentially many  $\mathbf{k}$

some  $\mathcal{C}$  allows efficient inference

- subset 0/1 error:  $\mathcal{C}(\mathbf{y}, \mathbf{k}) = 0$  iff  $\mathbf{y} = \mathbf{k}$  & 1 otherwise

optimal  $\mathbf{k} = \operatorname{argmax}_{\mathbf{y}} q(\mathbf{x}, \mathbf{y})$  over representative  $\mathbf{y}$

- Hamming error:  $\mathcal{C}(\mathbf{y}, \mathbf{k}) = \frac{1}{L} \sum_{\ell=1}^L \mathbb{I}[\mathbf{y}[\ell] \neq \mathbf{k}[\ell]]$

optimal  $\mathbf{k}[\ell] =$  majority bit of  $\mathbf{y}[\ell]$  over representative  $\mathbf{y}$

AOBD for CSML: generally **restricted to specific  $\mathcal{C}$**   
where difficulty in **inference** can be resolved

## Mini-Summary of Key Ideas

AOBD for CSML: if  $q(\mathbf{x}, \mathbf{y}) \approx P(\mathbf{y}|\mathbf{x})$  well

approximately good  $g_q(\mathbf{x}) = \operatorname{argmin}_{\mathbf{k} \in \{0,1\}^L} \sum_{\mathbf{y} \in \{0,1\}^L} q(\mathbf{x}, \mathbf{y}) \mathcal{C}(\mathbf{y}, \mathbf{k})$

difficulty of directly using AOBD revisited

- difficulty in **probability estimation**:
- difficulty in **cost calculation**:
- difficulty in **inference**:

$2^L$  outputs per  $\mathbf{x}$  for  $q(\mathbf{x}, \mathbf{y})$

$2^L$  possible  $\mathbf{y}$  in  $\sum$  to compute per  $\mathbf{k}$

$\operatorname{argmin}$  over  $2^L$  possible candidates  $\mathbf{k}$

corresponding key ideas

- estimate probability **with decomposition**
- sparsify probability estimation **with representative  $\mathbf{y}$ 's**
- derive **efficient inference rule** for specific  $\mathcal{C}$

next: **concrete approaches** that combine key ideas

# Putting It All Together: PCC for Subset 0/1 Error

- **training**: calculate  $q_1, q_2, \dots, q_L$ ,  
where  $q_\ell$  learned with extended inputs  $(\mathbf{x}_n, \mathbf{y}_n[1, \dots, \ell - 1])$  and outputs  $\mathbf{y}_n[\ell]$
- **inference**: for each  $\mathbf{x}$ ,
  - get  $B$  most **representative**  $\mathbf{y}$  by beam search (Kumar, 2013)
  - return  $g(\mathbf{x}) = \underset{\text{representative } \mathbf{y}}{\operatorname{argmax}} q(\mathbf{x}, \mathbf{y})$

## Cons of PCC

- ‘assymetric’ to labels: **non-trivial structural decision** of label order  
—often coupled with uniform aggregation (Ensemble PCC) to improve performance
- somewhat time consuming during **inference**

special case of  $B = 1$ :  
a classic approach called Classifier Chain (CC) (Read, 2009)

## Putting It All Together: PCC for F1 Score

- **training**: calculate  $q_1, q_2, \dots, q_L$ ,  
where  $q_\ell$  learned with extended inputs  $(\mathbf{x}_n, \mathbf{y}_n[1, \dots, \ell - 1])$  and outputs  $\mathbf{y}_n[\ell]$
- **inference**: for each  $\mathbf{x}$ ,
  - get  $B$  most **representative  $\mathbf{y}$**  by **sampling** (Dembczyński, 2011)
  - estimate **necessary statistics**

$$\delta_{\ell,k} = \mathbb{E}\left\{\mathbf{y}[\ell] \text{ given } |\mathbf{y}| = k\right\}$$

- return  $g(\mathbf{x})$  with exact inference from  $\delta_{\ell,k}$  using  $O(L^3)$  computation

strength depends on whether  $\delta_{\ell,k}$  estimated well enough  
with **probability estimation** + **representative sampling**

# Mini-Summary of the PCC Family

## with Efficient Inference Rules

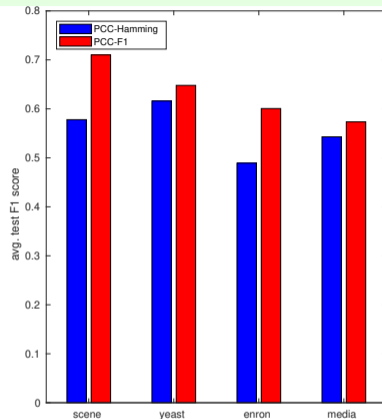
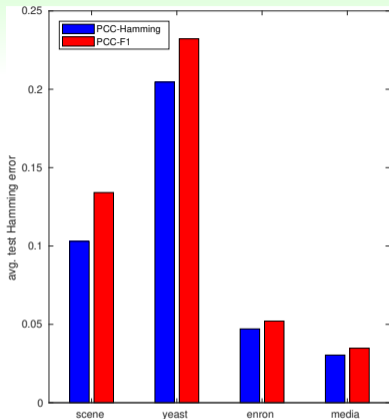
measure	inference rule
subset 0/1	mode of $q$
Hamming	threshold of 'marginal' $q$
ranking	sorted 'marginal' $q$
$\mathcal{C}(\mathbf{y}, \mathbf{k}) = -\frac{2 \mathbf{y} \cap \mathbf{k} }{ \mathbf{y}  +  \mathbf{k} }$ : F1	statistics $\delta_{\ell, k}$ from $q$

## without Efficient Inference Rules

measure	equation
accuracy	$\mathcal{C}(\mathbf{y}, \mathbf{k}) = -\frac{ \mathbf{y} \cap \mathbf{k} }{ \mathbf{y} \cup \mathbf{k} }$
multi. objective combination	$\mathcal{C}(\mathbf{y}, \mathbf{k}) = \mathcal{C}_1(\mathbf{y}, \mathbf{k}) + \mathcal{C}_2(\mathbf{y}, \mathbf{k})$

PCC: CSML approach 'in principle'

# Cost-Sensitivity of PCC



(Dembczyński, 2011) cost-sensitive behavior:

PCC-Hamming better (↓) for Hamming; PCC-F1 better (↑) for F1

# Outline

## 1 Cost-Sensitive Multiclass Classification

- CSMC Motivation and Setup
- CSMC by Bayesian Perspective
- CSMC by (Weighted) Binary Classification
- CSMC by Regression

## 2 Cost-Sensitive Multilabel Classification

- CSML Motivation and Setup
- CSML by Bayesian Perspective
- **CSML by (CS) Multiclass Classification**
- CSML by (Weighted) Binary Classification
- CSML by Regression

## 3 CSMC & CSML: Selected Applications

- Bacteria Classification with Doctor-Annotated Costs
- Network Connection Classification with 'Danger' Costs on Imbalanced Data
- Social Tagging with Costs from Tag Counts

## 4 Summary



# A Naïve CSML Approach: Cost-Sensitive Label Powerset

## Label Powerset (LP) Approach

multilabel  $\xrightarrow{\text{regular cost}}$  multiclass

## Cost-Sensitive LP (CSLP)

CSML  $\xrightarrow{\text{cost function } c}$  CSMC

## Cons of CSLP

- **complexity**, just like LP,

e.g. CSLP + CSOVO :  $O(2^L \cdot 2^L)$  classifiers  
 CSLP + FT :  $O(2^L)$  internal nodes

next: resolve complexity issue of CSLP

## Key Idea: Divide and Conquer

## Existing Multilabel Classification Approach for 'Speeding Up' LP:

Random  $k$ -Labelset (RAKEL) (Tsoumakas, 2007)

- in iteration  $t$  of training

- divide**: randomly choose  $k$  labels from  $\{1, 2, \dots, L\}$  as labelset  $S_t$

labelset	original label vector	restricted label vector
$S_1 = \{1, 2\}$	$\mathbf{y}_n = (0, 1, 0)$	$\mathbf{y}'_n = \mathbf{y}_n[S_1] = (0, 1)$
$S_2 = \{2, 3\}$	$\mathbf{y}_n = (0, 1, 0)$	$\mathbf{y}'_n = \mathbf{y}_n[S_2] = (1, 0)$

- train: **learn**  $\hat{g}_t(\mathbf{x})$  with LP on  $\{(\mathbf{x}_n, \mathbf{y}'_n = \mathbf{y}_n[S_t])\}$ , where  $\mathbf{y}[S]$  'restricts'  $\mathbf{y}$  to  $S$
- during prediction
  - conquer**: predict with  $g(\mathbf{x}) =$  voting per label from  $\hat{g}_t(\mathbf{x})$

labelset	predictions	votes
$S_1 = \{1, 2\}$	$\hat{g}_1(\mathbf{x}) = (0, 0)$	$(0, 0, -)$
$S_2 = \{2, 3\}$	$\hat{g}_2(\mathbf{x}) = (0, 1)$	$(-, 0, 1)$
$S_3 = \{1, 3\}$	$\hat{g}_3(\mathbf{x}) = (0, 1)$	$(0, -, 1)$
		voted $g(\mathbf{x}) = (0, 0, 1)$

RAKEL: training with LP **sufficiently efficient if  $k$  small**

## Cons of RAKEL

- **regular** internal classification, **cannot differentiate mistakes**

- small mistake:  $(0, 0)$  predicted as  $(0, 1)$
- big mistake:  $(0, 0)$  predicted as  $(1, 1)$

—can do **CSLP** for internal classification on **weighted Hamming error**:

Cost-Sensitive RAKEL (CS-RAKEL) (Lo, 2011)

- **voting**  $\approx$  Hamming error optimization, but **not a CSML approach**

- weaker performance for, e.g., F1 score

—can do **non-uniform voting** for **weighted Hamming error**:

Generalized  $k$ -Labelset Ensemble (GLE) (Lo, 2014)

next: a RAKEL-based approach that  
**handles 'any' cost function for CSML**

# Progressive $k$ -Labelset (PRAKEL) (Wu, 2017)

## cannot differentiate mistakes

- use **CSLP** for internal classification like (Lo, 2011)
- but with **general cost**

## not a CSML approach

- keep **voting** for simplicity
- but **calculate internal CSLP-costs from CSML-cost  $\mathcal{C}$**

## PRAKEL Approach

- in iteration  $t$  of training
  - **divide**: randomly choose  $k$  labels from  $\{1, 2, \dots, L\}$  as labelset  $S_t$  like RAKEL
  - **train**: learn  $\hat{g}_t(\mathbf{x})$  with **CSLP** on  $\{(\mathbf{x}_n, \mathbf{y}'_n = \mathbf{y}_n[S_t], \mathbf{c}'_n)\}$ ,  
where  $\mathbf{y}[S]$  'restricts'  $\mathbf{y}$  to  $S$  &  $\mathbf{c}'_n$  relates to  $\mathcal{C}$
- during prediction
  - **conquer**: predict with  $g(\mathbf{x}) = \text{voting per label from } \hat{g}_t(\mathbf{x})$

remaining issue: **defining  $\mathbf{c}'_n$**

## Key Idea: Reference Vector

## Goal

define  $\mathbf{c}'_n$  from  $(\mathbf{x}_n, \mathbf{y}_n)$  for  $S_t$  such that  $\mathbf{c}'_n$  relates to  $\mathcal{C}$

## FT

cost of **top-level prediction**  
calculated from **bottom-level prediction**

## PRAKEL

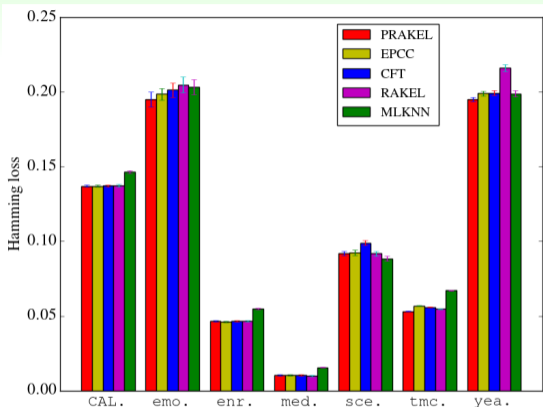
costs of **in- $S_t$  prediction**  
calculated from **out-of- $S_t$  prediction**

labelset	predictions	votes
$S_1 = \{2, 3\}$	$\hat{g}_1(\mathbf{x}) = (?, ?)$	$(-, ?, ?, -, -)$

- define **reference vector  $\mathbf{r}$**  containing the ‘-’ parts  
—proposed  $\mathbf{r}$  (Wu, 2017): predicted labels from  $\{\hat{g}_1, \hat{g}_2, \dots, \hat{g}_{t-1}\}$
- goal of  $\hat{g}_1$ : trade-off of ‘?’ parts within  $\mathcal{C}$
- $\mathbf{c}'_n$ : entries  $\mathcal{C}(\mathbf{y}_n, \mathbf{k})$  where  **$\mathbf{k}$  matches  $\mathbf{r}$  in ‘-’**

PRAKEL: RAKEL + CSPA + **progressive  $\mathbf{c}'_n$**  from  $(\mathbf{y}_n, \mathcal{C}, \text{reference } \mathbf{r}_n)$

# PRAKEL versus Others on Hamming Error

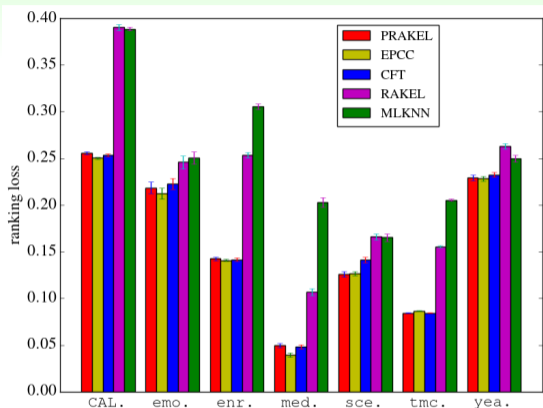


(Wu, 2017)

- all approaches fairly good
- **PRAKEL** and **EPCC** relatively stable

**PRAKEL: competitive for 'easier' measure**

# PRAKEL versus Others on Ranking Error

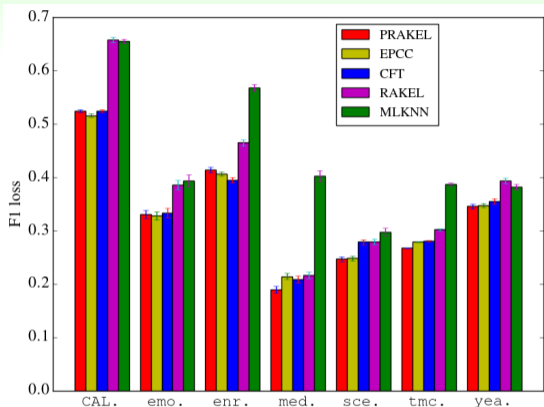


(Wu, 2017)

- **PRAKEL** better than cost-insensitive (**RAKEL** & **MLKNN**)
- **PRAKEL** competitive to cost-sensitive (**EPCC** & **CFT**)

**PRAKEL: strong for ranking error as well**

# PRAKEL versus Others on F1 Score



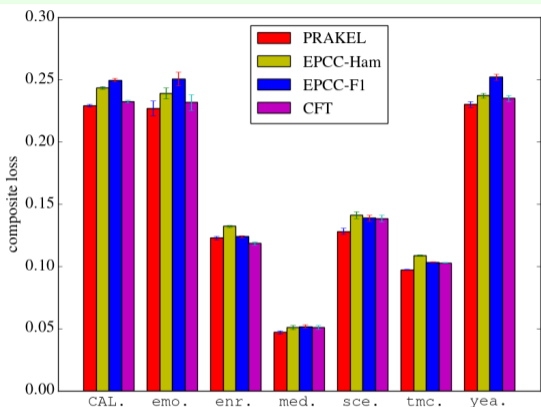
(Wu, 2017)

- **PRAKEL** sometimes the best

**PRAKEL** can be better than **EPCC**



# PRAKEL versus Others on Composite Error



(Wu, 2017)

- composite:  
 $\alpha$  Hamming -  $\beta$  F1
- **PRAKEL** significantly better than both
  - EPCC-Hamming
  - EPCC-F1

**PRAKEL** more useful for 'general' CSML than EPCC

# Outline

## 1 Cost-Sensitive Multiclass Classification

- CSMC Motivation and Setup
- CSMC by Bayesian Perspective
- CSMC by (Weighted) Binary Classification
- CSMC by Regression

## 2 Cost-Sensitive Multilabel Classification

- CSML Motivation and Setup
- CSML by Bayesian Perspective
- CSML by (CS) Multiclass Classification
- **CSML by (Weighted) Binary Classification**
- CSML by Regression

## 3 CSMC & CSML: Selected Applications

- Bacteria Classification with Doctor-Annotated Costs
- Network Connection Classification with 'Danger' Costs on Imbalanced Data
- Social Tagging with Costs from Tag Counts

## 4 Summary

# CSLP Revisited

## Label Powerset (LP) Approach

multilabel  $\xrightarrow{\text{regular cost}}$  multiclass

## Cost-Sensitive LP (CSLP)

CSML  $\xrightarrow{\text{cost function } c}$  CSMC

## Cons of CSLP

- **complexity**, just like LP,

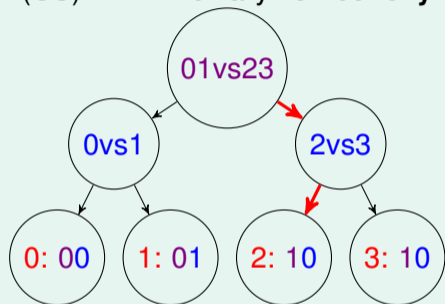
CSLP + ABOD :  $O(2^L)$  estimates  
 CSLP + CSOVO :  $O(2^L \cdot 2^L)$  classifiers  
 CSLP + FT :  $O(2^L)$  internal nodes

conquered by  
 decomposition + special inference  
 sampling + efficient decoding (Yang, 2018)  
 sampling + node sharing (Li, 2014)

next: CSLP + FT

## CSLP + FT for Prediction

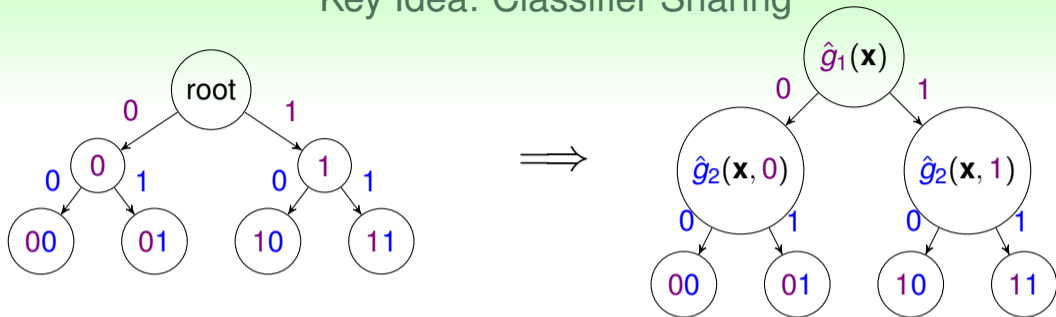
(CS)LP:  $Y = \text{binary number of } \mathbf{y}$



- with 'binary number encoding' (proper ordering):
  - $\ell$ -th layer nodes (classifier)  $\Leftrightarrow \ell$ -th label
- FT:  $O(\log K)$  prediction,  $O(K)$  training
  - $\log_2(2^L) = L$  **predictions only :-)**
  - still  $O(2^L)$  training complexity
  - actually,  $2^L - 1$  internal nodes

next: representing  $2^L - 1$  internal nodes efficiently

## Key Idea: Classifier Sharing



$2^L - 1$  nodes  $\Rightarrow L$  classifiers

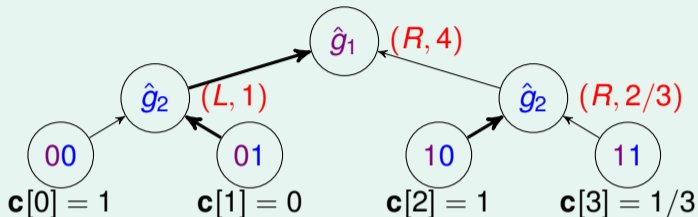
- root node  $\hat{g}_1(\mathbf{x})$ : just like BR on 1-st label
- 2-nd layer nodes 'shared' in  $\hat{g}_2(\mathbf{x}, \tilde{\mathbf{y}}[1])$ , where  $\tilde{\mathbf{y}}[1] = \hat{g}_1(\mathbf{x})$

CSLP + FT similar to PCC in prediction

PCC	$q_1(\mathbf{x})$	$q_2(\mathbf{x}, \tilde{\mathbf{y}}[1])$	$q_3(\mathbf{x}, \tilde{\mathbf{y}}[1], \tilde{\mathbf{y}}[2])$	...
CSLP + FT	$\hat{g}_1(\mathbf{x})$	$\hat{g}_2(\mathbf{x}, \tilde{\mathbf{y}}[1])$	$\hat{g}_3(\mathbf{x}, \tilde{\mathbf{y}}[1], \tilde{\mathbf{y}}[2])$	...

## CSLP + FP Training

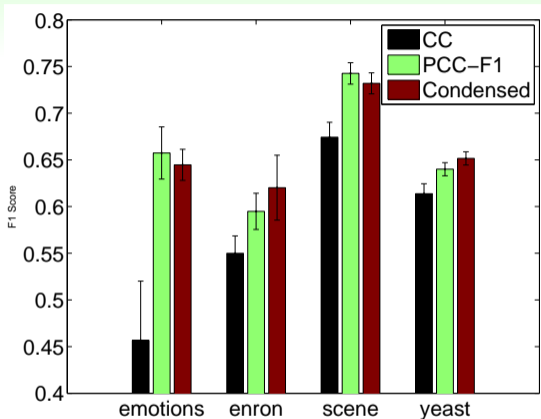
even with classifier sharing,  $2^L - 1$  weighted binary examples per  $(\mathbf{x}_n, \mathbf{y}_n)$  in FT



- not **all binary examples** relevant to training  $\hat{g}_\ell$   
—prediction goes through **one path** anyway
- Condensed FT (CFT) (Li, 2014) :  
keeping only those examples **near prediction path** for training  $\hat{g}_\ell$

CFT = CSLP + FT  
+ Proper Ordering + **Classifier Sharing** + **Example Sampling**

## CFT versus PCC on F1 Score



(Li, 2014)

- CFT certainly better than CC
- CFT can be better than PCC

**CFT competitive within 'chaining approaches' for CSML**

# Outline

## 1 Cost-Sensitive Multiclass Classification

- CSMC Motivation and Setup
- CSMC by Bayesian Perspective
- CSMC by (Weighted) Binary Classification
- CSMC by Regression

## 2 Cost-Sensitive Multilabel Classification

- CSML Motivation and Setup
- CSML by Bayesian Perspective
- CSML by (CS) Multiclass Classification
- CSML by (Weighted) Binary Classification
- **CSML by Regression**

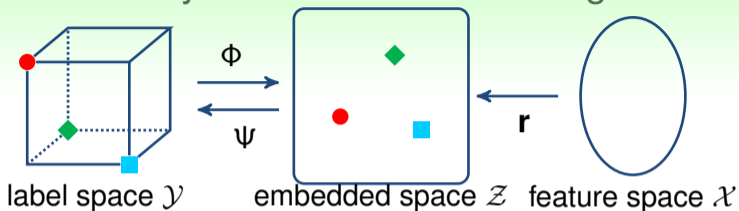
## 3 CSMC & CSML: Selected Applications

- Bacteria Classification with Doctor-Annotated Costs
- Network Connection Classification with 'Danger' Costs on Imbalanced Data
- Social Tagging with Costs from Tag Counts

## 4 Summary



## Key Idea: Label Embedding



## Training Stage

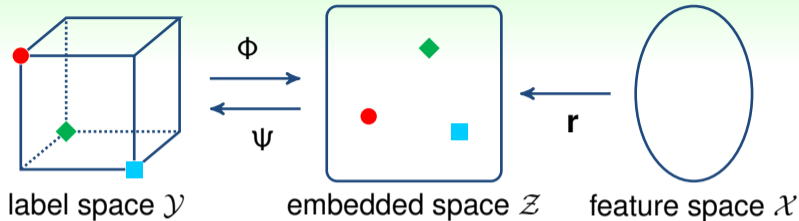
- **embedding function  $\phi$** : label vector  $\mathbf{y} \rightarrow$  embedded vector  $\mathbf{z}$
- learn a regressor  $\mathbf{r}$  from  $\{(\mathbf{x}_n, \mathbf{z}_n)\}_{n=1}^N$

## Predicting Stage

- for testing instance  $\mathbf{x}$ , predicted embedded vector  $\tilde{\mathbf{z}} = \mathbf{r}(\mathbf{x})$
- **decoding function  $\Psi$** :  $\tilde{\mathbf{z}} \rightarrow$  predicted label vector  $\tilde{\mathbf{y}}$

label embedding: popular for extracting **joint** information of labels

# Cost-Sensitive Label Embedding

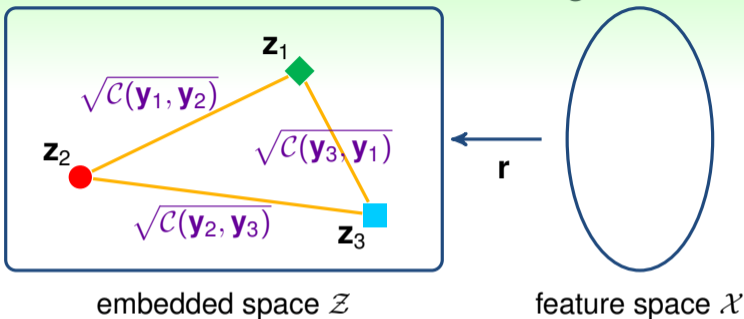


## Existing Works

- **label embedding**: PLST (Tai, 2012), CPLST (Chen, 2012), FaIE (Lin, 2014), RAKEL (Tsoumakas 2007), etc.
- **cost-sensitivity**: CFT (Li, 2014), PCC (Dembczyński, 2010), etc.
- **cost-sensitivity + label embedding**: ongoing

Cost-Sensitive Label Embedding: considering  $\mathcal{C}$  in  $\Phi$  and  $\Psi$

## Cost-Sensitive Encoding

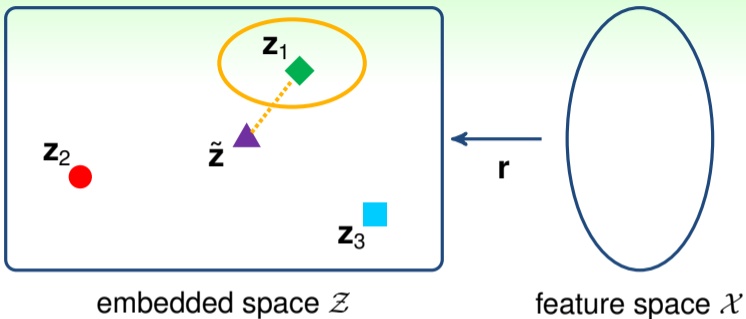


## Training Stage

- distances between embedded vectors  $\Leftrightarrow$  cost information
- larger distance  $d(\mathbf{z}_i, \mathbf{z}_j) \Leftrightarrow$  higher cost  $\mathcal{C}(\mathbf{y}_i, \mathbf{y}_j)$

$d(\mathbf{z}_i, \mathbf{z}_j) \approx \sqrt{\mathcal{C}(\mathbf{y}_i, \mathbf{y}_j)}$ : by  
 multidimensional scaling (Huang, 2017) or deep learning (Chiu, 2018)

## Cost-Sensitive Decoding



## Predicting Stage

- for testing instance  $\mathbf{x}$ , predicted embedded vector  $\tilde{\mathbf{z}} = \mathbf{r}(\mathbf{x})$
- find **nearest embedded vector  $\mathbf{z}_q$**  of  $\tilde{\mathbf{z}}$

cost-sensitive decoding:  $g(\mathbf{x}) = \text{corresponding } \mathbf{y}_q$

# Theoretical Explanation

## Cost Bound Theorem (Huang, 2017)

$$C(\mathbf{y}, \tilde{\mathbf{y}}) \leq 5 \left( \underbrace{(d(\mathbf{z}, \mathbf{z}_q) - \sqrt{C(\mathbf{y}, \mathbf{y}_q)})^2}_{\text{embedding error}} + \underbrace{\|\mathbf{z} - \mathbf{r}(\mathbf{x})\|^2}_{\text{regression error}} \right)$$

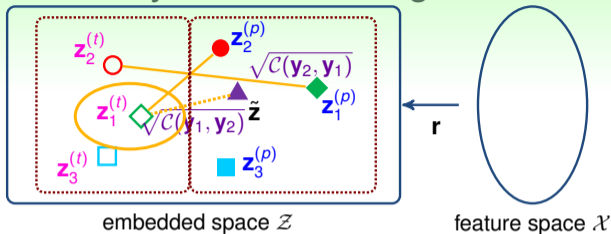
## Optimization

- **embedding error** → multidimensional scaling
- **regression error** → any regression  $\mathbf{r}$

challenge: **asymmetric cost function** vs. **symmetric distance**?

i.e.  $C(\mathbf{y}_i, \mathbf{y}_j) \neq C(\mathbf{y}_j, \mathbf{y}_i)$  vs.  $d(\mathbf{z}_i, \mathbf{z}_j)$

## Key Idea: Mirroring Trick



- two roles of  $\mathbf{y}_i$ : **ground truth role**  $\mathbf{y}_i^{(t)}$  and **prediction role**  $\mathbf{y}_i^{(p)}$ 
  - $\sqrt{C(\mathbf{y}_i, \mathbf{y}_j)} \Rightarrow$  predict  $\mathbf{y}_i$  as  $\mathbf{y}_j \Rightarrow$  for  $\mathbf{z}_i^{(t)}$  and  $\mathbf{z}_j^{(p)}$
  - $\sqrt{C(\mathbf{y}_j, \mathbf{y}_i)} \Rightarrow$  predict  $\mathbf{y}_j$  as  $\mathbf{y}_i \Rightarrow$  for  $\mathbf{z}_i^{(p)}$  and  $\mathbf{z}_j^{(t)}$
- learn **regression function**  $\mathbf{r}$  from  $\mathbf{z}_1^{(p)}, \mathbf{z}_2^{(p)}, \dots, \mathbf{z}_L^{(p)}$
- find **nearest embedded vector** of  $\tilde{\mathbf{z}}$  from  $\mathbf{z}_1^{(t)}, \mathbf{z}_2^{(t)}, \dots, \mathbf{z}_L^{(t)}$

mirroring trick: handle asymmetric cost with **embedding flexibility**

# Cost-Sensitive Label Embedding with Multidimensional Scaling (CLEMS)

## training stage of CLEMS (Huang, 2017)

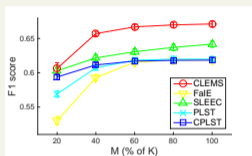
- given training instances  $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$  and cost function  $\mathcal{C}$
- apply mirroring trick to set up  $\mathbf{z}_n^{(t)}$  and  $\mathbf{z}_n^{(p)}$  for label vector  $\mathbf{y}_n$
- compute embedding function  $\Phi: \mathbf{y}_n \rightarrow \mathbf{z}_n^{(p)}$  by multidimensional scaling such that  $d(\mathbf{z}_m^{(t)}, \mathbf{z}_n^{(p)}) \approx \sqrt{\mathcal{C}(y_n, y_m)}$
- learn a regression function  $\mathbf{r}$  from  $\{(\mathbf{x}_n, \mathbf{z}_n^{(p)} = \Phi(\mathbf{y}_n))\}_{n=1}^N$

## predicting stage of CLEMS

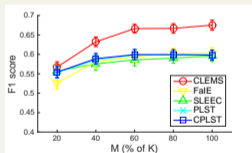
- given the testing instance  $\mathbf{x}$
- obtain the predicted embedded vector by  $\tilde{\mathbf{z}} = \mathbf{r}(\mathbf{x})$
- prediction  $\tilde{\mathbf{y}} = \Psi(\tilde{\mathbf{z}}) = \Phi^{-1}(\text{nearest neighbor}) = \Phi^{-1}(\text{argmin } d(\mathbf{z}_n^{(t)}, \tilde{\mathbf{z}}))$

minor details: embed **subset of**, rather than 'all',  $\{0, 1\}^L$  for efficiency

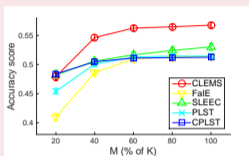
## Comparison with Label Embedding Approaches

F1 score ( $\uparrow$ )

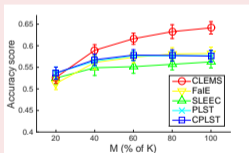
yeast



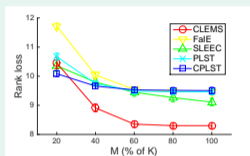
birds

Accuracy score ( $\uparrow$ )

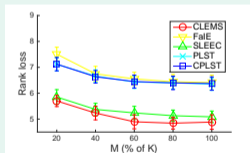
yeast



birds

Rank loss ( $\downarrow$ )

yeast



birds

CLEMS: best label embedding approach across different criteria



# Comparison with Cost-Sensitive Algorithms

data	F1 score ( $\uparrow$ )			Accuracy score ( $\uparrow$ )			Rank loss ( $\downarrow$ )		
	CLEMS	CFT	PCC	CLEMS	CFT	PCC	CLEMS	CFT	PCC
emot.	<b>0.676</b>	0.640	0.643	<b>0.589</b>	0.557	–	1.484	1.563	<b>1.467</b>
scene	<b>0.770</b>	0.703	0.745	<b>0.760</b>	0.656	–	0.672	0.723	<b>0.645</b>
yeast	<b>0.671</b>	0.649	0.614	<b>0.568</b>	0.543	–	<b>8.302</b>	8.566	8.469
birds	<b>0.677</b>	0.601	0.636	<b>0.642</b>	0.586	–	4.886	4.908	<b>3.660</b>
med.	<b>0.814</b>	0.635	0.573	<b>0.786</b>	0.613	–	5.170	5.811	<b>4.234</b>
enron	<b>0.606</b>	0.557	0.542	<b>0.491</b>	0.448	–	29.40	26.64	<b>25.11</b>
lang.	<b>0.375</b>	0.168	0.247	<b>0.327</b>	0.164	–	31.03	34.16	<b>19.11</b>
flag	<b>0.731</b>	0.692	0.706	<b>0.615</b>	0.588	–	2.930	3.075	<b>2.857</b>
slash	<b>0.568</b>	0.429	0.503	<b>0.538</b>	0.402	–	4.986	5.677	<b>4.472</b>
CAL.	<b>0.419</b>	0.371	0.391	<b>0.273</b>	0.237	–	1247	1120	<b>993</b>
arts	<b>0.492</b>	0.334	0.349	<b>0.451</b>	0.281	–	9.865	10.07	<b>8.467</b>
EUR.	<b>0.670</b>	0.456	0.483	<b>0.650</b>	0.450	–	89.52	129.5	<b>43.28</b>

- **generality for CSML:**  $\text{CLEMS} = \text{CFT} > \text{PCC}$
- **performance:**  $\text{CLEMS} \approx \text{PCC} > \text{CFT}$
- **speed:**  $\text{CLEMS} \approx \text{PCC} > \text{CFT}$

CLEMS: very **competitive** for CSML

# Outline

## 1 Cost-Sensitive Multiclass Classification

- CSMC Motivation and Setup
- CSMC by Bayesian Perspective
- CSMC by (Weighted) Binary Classification
- CSMC by Regression

## 2 Cost-Sensitive Multilabel Classification

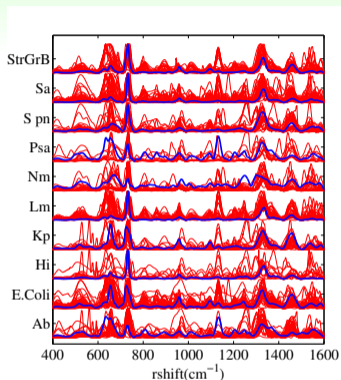
- CSML Motivation and Setup
- CSML by Bayesian Perspective
- CSML by (CS) Multiclass Classification
- CSML by (Weighted) Binary Classification
- CSML by Regression

## 3 CSMC & CSML: Selected Applications

- **Bacteria Classification with Doctor-Annotated Costs**
- Network Connection Classification with 'Danger' Costs on Imbalanced Data
- Social Tagging with Costs from Tag Counts

## 4 Summary

# A Real Medical Application: Bacteria Classification



automatic classification from spectrum to bacterium (Jan, 2011)

**are all mis-prediction costs equal?**

# A Real Medical Application of CSMC: Classifying Bacteria

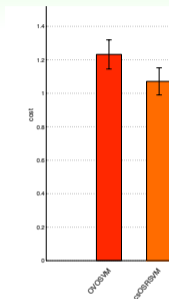
## The Problem

- Gram-positive as Gram-positive: small cost  
Gram-positive as Gram-negative: **big cost**
- cost matrix averaged from two doctors:

	Ab	Ecoli	HI	KP	LM	Nm	Psa	Spn	Sa	GBS
Ab	0	1	10	7	9	9	5	8	9	1
Ecoli	3	0	10	8	10	10	5	10	10	2
HI	10	10	0	3	2	2	10	1	2	10
KP	7	7	3	0	4	4	6	3	3	8
LM	8	8	2	4	0	5	8	2	1	8
Nm	3	10	9	8	6	0	8	3	6	7
Psa	7	8	10	9	9	7	0	8	9	5
Spn	6	10	7	7	4	4	9	0	4	7
Sa	7	10	6	5	1	3	9	2	0	7
GBS	2	5	10	9	8	6	5	6	8	0

issue 1: is cost-sensitive classification **really useful**?

# Cost-Sensitive vs. Traditional on Bacteria Data

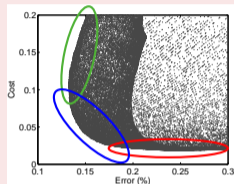


(Jan, 2011)

**cost-sensitive** better than **traditional**;  
but why are people **still not**  
using those cool ML works for their AI? :-)

## Issue 2: Error Rate of Cost-Sensitive Classifiers

### The Problem



- cost-sensitive classifier: **low cost but high error rate**
- traditional classifier: **low error rate but high cost**
- how can we get the **blue** classifiers?: **low error rate and low cost**

cost-and-error-sensitive:  
more suitable for **real-world medical needs**

# Improved Classifier for Both Cost and Error

(Jan, 2012)

Cost	
iris	≈
wine	≈
glass	≈
vehicle	≈
vowel	○
segment	○○
dna	○○
satimage	≈
usps	○○
zoo	○○
splice	≈
ecoli	≈
soybean	≈

Error	
iris	○
wine	○
glass	○
vehicle	○
vowel	○
segment	○○
dna	○○
satimage	○
usps	○
zoo	○
splice	○
ecoli	○
soybean	○

now, are people using those cool ML works for their AI? :-)

# Lessons Learned from CSMC Research in Applications



?



H7N9-infected



cold-infected



healthy

- 1 more realistic (generic) in academia  $\neq$  more realistic (feasible) in application  
e.g. the 'cost' of **inputting a cost matrix?** :-)
- 2 **cross-domain collaboration** important  
e.g. getting the 'cost matrix' from **domain experts**
- 3 not easy to win **human trust**  
—humans are somewhat **multi-objective**

important yet **challenging**  
to use CSMC/CSML in practical applications



# Outline

## 1 Cost-Sensitive Multiclass Classification

- CSMC Motivation and Setup
- CSMC by Bayesian Perspective
- CSMC by (Weighted) Binary Classification
- CSMC by Regression

## 2 Cost-Sensitive Multilabel Classification

- CSML Motivation and Setup
- CSML by Bayesian Perspective
- CSML by (CS) Multiclass Classification
- CSML by (Weighted) Binary Classification
- CSML by Regression

## 3 CSMC & CSML: Selected Applications

- Bacteria Classification with Doctor-Annotated Costs
- **Network Connection Classification with 'Danger' Costs on Imbalanced Data**
- Social Tagging with Costs from Tag Counts

## 4 Summary

# A Real Security Application: Intrusion Detection

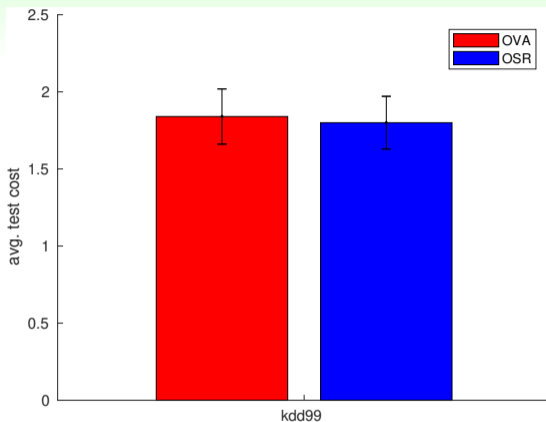
## KDDCup 2009 Problem

	normal	probe	DOS	U2R	R2L
normal	0	1	2	2	2
probe	1	0	2	2	2
DOS	2	1	0	2	2
U2R	3	2	2	0	2
R2L	4	2	2	2	0

- 'usual' mis-prediction: 2
- dangerous intrusion mis-predicted as normal: 4
- suspicious intrusion mis-predicted as normal: 1

KDDCup 1999: 'earliest' public data for CSMC

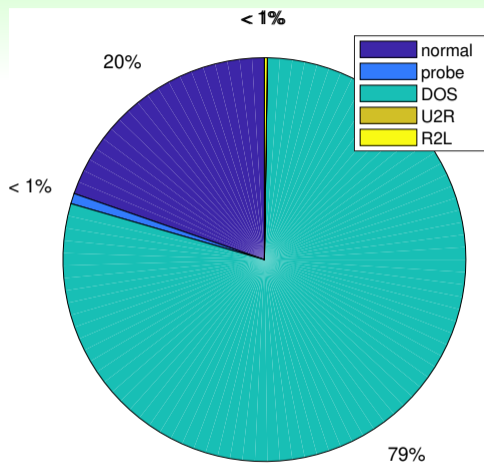
# Experiment Result for KDDCup 1999 (Jan, 2012)



- OVA: regular
- OSR: CSMC

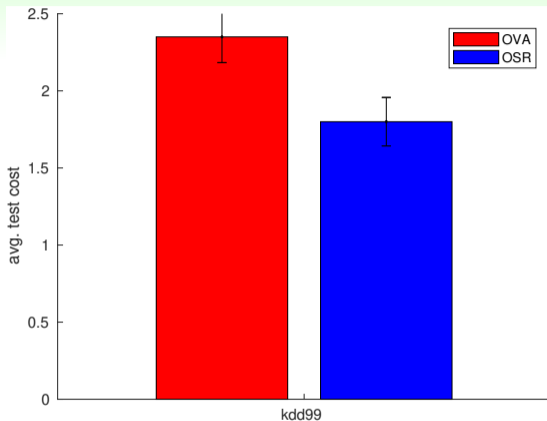
difference between **regular** & CSMC **small**, why? :-)

# Distribution of KDDCup 1999 Data



**KDDCup 1999: highly imbalanced!!**

# Experiment Result for Balanced KDDCup 1999 (Jan, 2012)



balanced cost  
= inverse class ratio  $\cdot \mathcal{C}$

- OVA: regular
- OSR: CSMC

**balancing by inverse class ratio** (special case of CSMC) helpful!

# Lessons Learned on CSMC for Imbalanced Classification

- 1 balancing (re-weighting) by **inverse class ratio**: helpful and strong baseline
- 2 'cost = relative class ratio': **the first benchmark of CSMC** (Domingos, 1999), also helpful
- 3 CSMC helps **tune other measure** of imbalanced classification, such as G-mean (Khan, 2018)

## Usage of Cost (in General)

- **representing real application needs**
  - human-annotated: for bacteria classification
  - heuristic: KDDCup 1999 cost matrix
- **adjusting learning condition**
  - imbalanced classification
  - 'easy'/'hard' classes
- **tuning for other metric of interest**

CSMC: a '**swiss knife**' for imbalanced classification

# Outline

## 1 Cost-Sensitive Multiclass Classification

- CSMC Motivation and Setup
- CSMC by Bayesian Perspective
- CSMC by (Weighted) Binary Classification
- CSMC by Regression

## 2 Cost-Sensitive Multilabel Classification

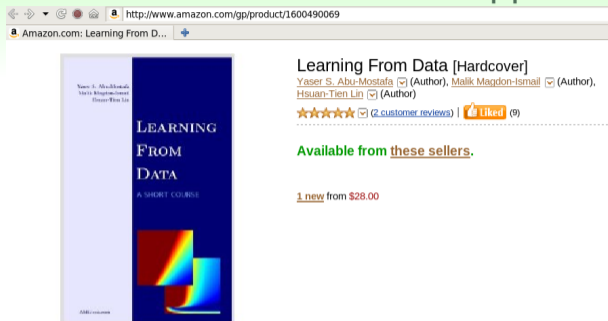
- CSML Motivation and Setup
- CSML by Bayesian Perspective
- CSML by (CS) Multiclass Classification
- CSML by (Weighted) Binary Classification
- CSML by Regression

## 3 CSMC & CSML: Selected Applications

- Bacteria Classification with Doctor-Annotated Costs
- Network Connection Classification with 'Danger' Costs on Imbalanced Data
- Social Tagging with Costs from Tag Counts

## 4 Summary

# A Real Internet Application: Social Tagging



http://www.amazon.com/gp/product/1600490069

Amazon.com: Learning From D...

**Learning From Data [Hardcover]**  
 Yaser S. Abu-Mostafa (Author), Malik Magdon-Ismael (Author),  
 Hsuan-Tien Lin (Author)

★★★★★ (2 customer reviews) | Liked (9)

Available from [these sellers](#).

1 new from \$28.00

? : { machine learning, data structure, data mining, object-oriented programming, artificial intelligence, compiler, architecture, chemistry, textbook, children-book, ... etc. }

social tagging: **tags** given by **internet users**

- **high-count** tags: popularly recognized, more **salient**
- **low-count** tags: weakly related or even **noisy**

social tagging: instance-tag count reflect **importance** of **tag**



# Cost of Social Tagging Data

## Social Tagging for URL Categories (Lo, 2014)

url instance $\mathbf{x}$	high count	low count
www.tvguide.com	entertain	art
www.ebookee.com	ebook	education
www.python.com/doc	python	c

$$C_{\mathbf{x}}(\mathbf{y}, \mathbf{k}) = \sum_{\ell=1}^L c_{\ell}[\mathbf{x}] \mathbb{I}[\mathbf{y}[\ell] \neq \mathbf{k}[\ell]],$$

where per-(label, instance) cost  $c_{\ell}[\mathbf{x}] = \begin{cases} \text{for tag: count for } \mathbf{x} \\ \text{for no-tag: constant to balance tag \& no-tag costs} \end{cases}$

$C_{\mathbf{x}}$ : weighted Hamming error per  $\mathbf{x}$

# Per-Instance Weighted Hamming Error

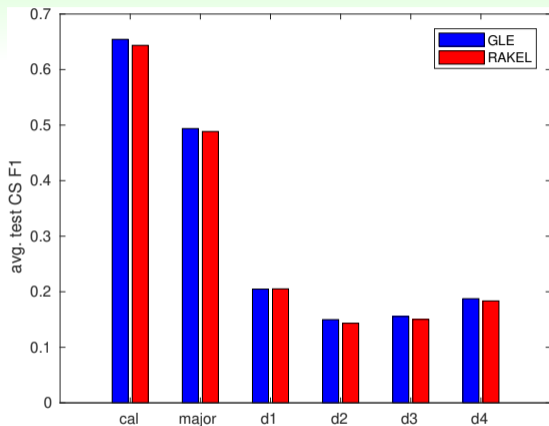
$$C_{\mathbf{x}}(\mathbf{y}, \mathbf{k}) = \sum_{\ell=1}^L c_{\ell}[\mathbf{x}] \mathbb{I}[\mathbf{y}[\ell] \neq \mathbf{k}[\ell]],$$

where per-(label, instance) cost  $c_{\ell}[\mathbf{x}] = \begin{cases} \text{for tag: count for } \mathbf{x} \\ \text{for no-tag: constant to balance tag \& no-tag costs} \end{cases}$

- $C$  per instance: more general than our CSML setting
  - could use **non-Bayesian CSML approaches**
- each  $C_{\mathbf{x}}$  just weighted Hamming error: less general than our CSML setting
  - GLE** (Lo, 2014) suffices, might not need **PRAKEL** (Wu, 2017)

is CSML really useful?

# GLE versus RAKEL for Social Tagging



(Lo, 2014)

- GLE training measure: weighted Hamming
- evaluation measure: **weighted F1** instead of weighted Hamming
- GLE better than RAKEL

CSML helps, even when training/evaluation not fully matching

# Summary

- **cost-sensitive multiclass classification: class/example-dependent**
  - Bayesian: MetaCost (Domingos, 1999)
  - non-Bayesian: Data Space Expansion (Abe, 2004) (to multiclass), Cost-Sensitive One-Versus-One (Lin, 2014), Filter Tree (Beygelzimer, 2009), . . . (to binary), One-Sided Regression (Tu, 2010) (to regression) —some SVM implementations here:  
<http://www.csie.ntu.edu.tw/~htlin/program/cssvm/>
- **cost-sensitive multilabel classification:**
  - Bayesian: PCC (Dembczyński, 2010)
  - non-Bayesian: Progressive  $k$ -Labelsets (Wu, 2017) (to multiclass), Condensed Filter Tree (Li, 2014) (to binary), CLEMS (Huang, 2017) (to regression)
- **applicatons & beyond:**
  - cost-and-error-sensitive learning for bacteria classification (Jan, 2012)
  - cost-sensitive learning for imbalanced learning (Jan, 2012)
  - cost-sensitive learning for social tagging prediction (Lo, 2014)

discussion welcomed on algorithm and **application** opportunities

## Giants' Shoulder (1/2)

- Domingos, MetaCost: A General Method for Making Classifiers Cost-Sensitive, 1999
- Elkan, The Foundations of Cost-Sensitive Learning, 2001
- Zadrozny, Cost-Sensitive Learning by Cost-Proportionate Example Weighting, 2003
- Abe, An Iterative Method for Multi-Class Cost-Sensitive Learning, 2004
- Beygelzimer, Error Limiting Reductions Between Classification Tasks, 2005
- Langford, Sensitive Error Correcting Output Codes, 2005
- Tsoumakas, Random  $k$ -Labelsets: An Ensemble Method for Multilabel Classification, 2007
- Beygelzimer, Error-Correcting Tournaments, 2009
- Read, Classifier Chains for Multi-label Classification, 2009
- Lin, A Simple Cost-Sensitive Multiclass Classification Algorithm Using One-Versus-One Comparisons, 2014
- Tu, One-Sided Support Vector Regression for Multiclass Cost-Sensitive Classification, 2010
- Dembczyński, Bayes Optimal Multilabel Classification via Probabilistic Classifier Chains, 2010
- Dembczyński, An Exact Algorithm for F-Measure Maximization, 2011
- Jan, Cost-Sensitive Classification on Pathogen Species of Bacterial Meningitis by Surface Enhanced Raman Scattering, 2011
- Lo, Cost-sensitive Multi-label Learning for Audio Tag Annotation and Retrieval, 2011

## Giants' Shoulder (2/2)

- Dembczyński, On Label Dependence and Loss Minimization in Multi-Label Classification, 2012
- Chen, Feature-Aware Label Space Dimension Reduction for Multi-Label Classification, 2012
- Jan, A Simple Methodology for Soft Cost-Sensitive Classification, 2012
- Tai, Multilabel Classification with Principal Label Space Transformation, 2012
- Kumar, Beam Search Algorithms for Multilabel Learning, 2013
- Li, Condensed Filter Tree for Cost-Sensitive Multi-Label Classification, 2014
- Lin, Multi-label Classification via Feature-Aware Implicit Label Space Encoding, 2014
- Lo, Generalized  $k$ -Labelsets Ensemble for Multi-Label and Cost-Sensitive Classification, 2014
- Chung, Cost-Aware Pre-Training for Multiclass Cost-Sensitive Deep Learning, 2016a
- Chung, Cost-Sensitive Deep Learning with Layer-Wise Cost Estimation, 2016b
- Huang, Cost-Sensitive Label Embedding for Multi-Label Classification, 2017
- Khan, Cost-Sensitive Learning of Deep Feature Representations From Imbalanced Data, 2018
- Wu, Progressive  $k$ -Labelsets for Cost-Sensitive Multi-Label Classification, 2017
- Chiu, Multi-Label Classification with Feature-Aware Cost-Sensitive Label Embedding, 2018
- Hsieh, A Deep Model with Local Surrogate Loss for General Cost-Sensitive Multi-Label Learning, 2018
- Yang, Cost-Sensitive Reference Pair Encoding for Multi-Label Learning, 2018

# Acknowledgments

- KDD Organizers!
- Computational Learning Lab @ NTU and Learning Systems Group @ Caltech for discussions

final advertisement

The Appier logo is displayed in a blue, stylized serif font on a white rectangular background.The CrossX logo consists of the word "CrossX" in white, bold, sans-serif font, centered within a blue rectangular background.The AIXON logo features the word "AIXON" in white, bold, sans-serif font, centered within a blue rectangular background.The AIQUA logo shows the word "AIQUA" in white, bold, sans-serif font, centered within a blue rectangular background.

(yes, we are hiring!!)

Thank you. Questions?