## From Ordinal Ranking to Binary Classification

#### Hsuan-Tien Lin

Department of Computer Science and Information Engineering National Taiwan University

#### Talk in the Applied Math Department at NDHU February, 2010

Joint work with Dr. Ling Li at Caltech (ALT'06, NIPS'06)





- 2 Ordinal Ranking Setup
- 3) The Reduction Framework
  - Key Ideas
  - Important Properties
  - Theoretical Usefulness
  - Algorithmic Usefulness
- Experimental Results



# Which Digit Did You Write?





Hsuan-Tien Lin (CSIE, NTU)



Hsuan-Tien Lin (CSIE, NTU)



• {yes, no}: binary classification

# new types of machine learning problems keep coming from new applications

Hsuan-Tien Lin (CSIE, NTU)



- 2 Ordinal Ranking Setup
- The Reduction Framework
  - Key Ideas
  - Important Properties
  - Theoretical Usefulness
  - Algorithmic Usefulness
- Experimental Results



Ordinal Ranking Setup

# Which Age-Group?





# Properties of Ordinal Ranking (1/2)



# general classification cannot properly use order information



Hsuan-Tien Lin (CSIE, NTU)

### Hot or Not?





#### rank: natural representation of human preferences

Hsuan-Tien Lin (CSIE, NTU)

# Properties of Ordinal Ranking (2/2)

ranks do **not** carry numerical information

rating 9 not 2.25 times "hotter" than rating 4

Select a rating to see the next picture. NOT 01 02 03 04 05 06 07 08 09 010 HOT

actual metric hidden



# general regression deteriorates without correct numerical information



# How Much Did You Like These Movies?

http://www.netflix.com



goal: use "movies you've rated" to automatically predict your **preferences (ranks)** on future movies



## **Ordinal Ranking Setup**

#### Given

*N* examples (input  $x_n$ , rank  $y_n$ )  $\in \mathcal{X} \times \mathcal{Y}$ 

- age-group:  $\mathcal{X} = encoding(human pictures), \mathcal{Y} = \{1, \cdots, 4\}$
- hotornot:  $\mathcal{X} = encoding(human pictures), \mathcal{Y} = \{1, \cdots, 10\}$
- netflix:  $\mathcal{X} = encoding(movies), \mathcal{Y} = \{1, \cdots, 5\}$

#### Goal

an ordinal ranker (decision function) r(x) that "closely predicts" the ranks *y* associated with some **unseen** inputs *x* 

ordinal ranking: a hot and important research problem



# Importance of Ordinal Ranking

- relatively new for machine learning (not so new for statisticians)
- connecting classification and regression
- matching human preferences—many applications in social science, information retrieval, psychology and recommendation systems





# Formalizing (Non-)Closeness: Cost

- ranks carry no numerical information: how to say "close"?
- artificially quantify the cost of being wrong

e.g. loss of customer loyalty when the system says ★★★★but you feel ★★☆☆☆☆

#### Yes! cost == loss for statisticians

cost vector c of example (x, y, c):
c[k] = cost when predicting (x, y) as rank k
e.g. for (Sweet Home Alabama,★★☆☆☆), a proper cost is c = (1,0,2,10,15)

#### closely predict: small cost during testing



#### Ordinal Ranking Setup

# **Ordinal Cost Vectors**

For an ordinal example  $(x, y, \mathbf{c})$ , the cost vector  $\mathbf{c}$  should

- be consistent with rank y:  $\mathbf{c}[y] = \min_k \mathbf{c}[k] (= 0)$
- respect order information: V-shaped (ordinal) or even convex (strongly ordinal)



# **Our Contributions**

T a theoretical and algorithmic foundation of ordinal ranking, which reduces ordinal ranking to binary classification, and ...

- provides a methodology for designing new ordinal ranking algorithms with any ordinal cost effortlessly
- takes many existing ordinal ranking algorithms as special cases
- introduces **new theoretical guarantee** on the generalization performance of ordinal rankers
- leads to superior experimental results







#### Ordinal Ranking Setup

# Central Idea: Reduction



complex ordinal ranking problems



simpler binary classification problems with well-known results on models, algorithms and theories

(cassette player)

# If I have seen further it is by standing on the shoulders of Giants—I. Newton



Hsuan-Tien Lin (CSIE, NTU)

3



2 Ordinal Ranking Setup

#### The Reduction Framework

- Key Ideas
- Important Properties
- Theoretical Usefulness
- Algorithmic Usefulness

#### Experimental Results





- 2 Ordinal Ranking Setup
- 3 The Reduction Framework
  - Key Ideas
  - Important Properties
  - Theoretical Usefulness
  - Algorithmic Usefulness
  - Experimental Results



# Threshold Ranker

if getting an ideal score s(x) of a movie x, how to construct the discrete r(x) from an analog s(x)?



quantize s(x) by **ordered** (non-uniform) thresholds  $\theta_k$ 

- ocommonly used in previous work:
  - threshold perceptrons
  - threshold hyperplanes
  - threshold ensembles

(PRank, Crammer and Singer, 2002) (SVOR, Chu and Keerthi, 2005)

(ORBoost, Lin and Li, 2006)

hreshold ranker: 
$$r(x) = \min \{k : s(x) < \theta_k\}$$



# The Reduction Framework Key Ideas Key Idea: Associated Binary Queries

getting the rank using a threshold ranker

- is  $s(x) > \theta_1$ ? Yes
- is  $s(x) > \theta_2$ ? No
- is  $s(x) > \theta_3$ ? No
- is  $s(x) > \theta_4$ ? No

generally, how do we query the rank of a movie *x*?

- is movie x better than rank 1? Yes
- Is movie x better than rank 2? No
- is movie x better than rank 3? No
  - is movie x better than rank 4? No

#### associated binary queries: is movie *x* better than rank *k*?



# More on Associated Binary Queries

say, the machine uses g(x, k) to answer the query "*is movie x better than rank k*?" e.g. for threshold ranker:  $g(x, k) = sign(s(x) - \theta_k)$ 





# Computing Ranks from Associated Binary Queries

Key Ideas

when g(x, k) answers "is movie x better than rank k?"

Consider  $(g(x, 1), g(x, 2), \cdots, g(x, K-1)),$ 

- concordant predictions: (Y, Y, N, N, N, N, N)
- extracting the rank from concordant predictions:

The Reduction Framework

- minimum index searching:  $r_g(x) = \min \{k : g(x, k) = \mathbb{N}\}$
- counting:  $r_g(x) = 1 + \sum_k \llbracket g(x,k) = Y \rrbracket$
- two approaches equivalent for concordant predictions
- mistaken/non-concordant predictions? e.g. (Y, N, Y, Y, N, N, Y)

#### counting: simpler to analyze and robust to mistake



# The Counting Approach

Say y = 5, i.e.,  $((z)_1, (z)_2, \cdots, (z)_7) = (Y, Y, Y, Y, N, N, N)$ 

if g<sub>1</sub>(x, k) reports concordant predictions (Y, Y, N, N, N, N, N)

Key Ideas

•  $g_1(x,k)$  made 2 binary classification errors

The Reduction Framework

•  $r_{g_1}(x) = 3$  by counting: the absolute cost is 2

absolute cost = # of binary classification errors

- if g<sub>2</sub>(x, k) reports non-concordant predictions (Y, N, Y, Y, N, N, Y)
  - $g_2(x, k)$  made 2 binary classification errors
  - $r_{g_2}(x) = 5$  by counting: the absolute cost is 0

absolute cost  $\leq$  # of binary classification errors

If 
$$(z)_k$$
 = desired answer &  $r_g$  computed by counting,  
 $|y - r_g(x)| \le \sum_{k=1}^{K-1} \left[ (z)_k \ne g(x,k) \right] .$ 

#### The Reduction Framework Key Ideas

### Binary Classification Error v.s. Ordinal Ranking Cost

## Say y = 5, i.e., $((z)_1, (z)_2, \cdots, (z)_7) = (Y, Y, Y, Y, N, N, N)$

- if  $g_1(x, k)$  reports concordant predictions (Y, Y, N, N, N, N, N)
  - $g_1(x,k)$  made 2 binary classification errors
  - $r_{g_1}(x) = 3$  by counting: the **squared** cost is 4
- if g<sub>3</sub>(x, k) reports concordant predictions (Y, N, N, N, N, N, N)
  - $g_3(x,k)$  made 3 binary classification errors
  - $r_{g_3}(x) = 2$  by counting: the **squared** cost is 9
    - 1 error in binary classification
  - $\implies$  5 cost in ordinal ranking



# Importance of Associated Binary Examples

• 
$$(w)_k \equiv \left| \mathbf{c}[k+1] - \mathbf{c}[k] \right|$$
: the importance of  $((x,k), (z)_k)$ 

per-example cost bound (Li and Lin, 2007): for **concordant predictions** or **strongly ordinal costs**  $\mathbf{c}[r_g(x)] \leq \sum_{k=1}^{K-1} (w)_k [[(z)_k \neq g(x,k)]]$ 



3



2 Ordinal Ranking Setup

#### The Reduction Framework

- Key Ideas
- Important Properties
- Theoretical Usefulness
- Algorithmic Usefulness

#### Experimental Results



### The Reduction Framework (1/2)



#### the reduction framework: systematic & easy to implement



Hsuan-Tien Lin (CSIE, NTU)

# The Reduction Framework (2/2)



• performance guarantee:

accurate binary predictions  $\Longrightarrow$  correct ranks

#### wide applicability: works with any ordinal c & any binary classification algorithm

#### • simplicity:

mild computation overheads with O(NK) binary examples

#### state-of-the-art:

allows new improvements in binary classification to be immediately inherited by ordinal ranking



The Reduction Framework

Important Properties

### Theoretical Guarantees of Reduction (1/3)

absolutely good binary classifier absolutely good ranker? YES!

error transformation theorem (Li and Lin, 2007)

For **concordant predictions** or **strongly ordinal costs**, if *g* makes test error  $\Delta$  in the induced binary problem, then  $r_g$  pays test cost at most  $\Delta$  in ordinal ranking.

- a one-step extension of the per-example cost bound
- conditions: general and minor
- performance guarantee in the absolute sense

what if no "absolutely good" binary classifier?



The Reduction Framework

Important Properties

### Theoretical Guarantees of Reduction (2/3)

- absolutely good binary classifier
  - $\implies$  absolutely good ranker? YES!
- relatively good binary classifier relatively good ranker? YES!

regret transformation theorem (Lin, 2008)

For **concordant predictions** or **strongly ordinal costs**, if *g* is  $\epsilon$ -close to the optimal binary classifier  $g_*$ , then  $r_g$  is  $\epsilon$ -close to the optimal ranker  $r_*$ .

# "reduction to binary" sufficient for algorithm design, **but necessary?**



The Reduction Framework

Important Properties

### Theoretical Guarantees of Reduction (3/3)

- absolutely good binary classifier
  - $\implies$  absolutely good ranker? **YES**!
- relatively good binary classifier relatively good ranker? YES!
- algorithm producing relatively good binary classifier algorithm producing relatively good ranker? YES!

#### equivalence theorem (Lin, 2008)

For a general family of **ordinal costs**, a good ordinal ranking algorithm exists **if & only if** a good binary classification algorithm exists for the corresponding learning model.

#### ordinal ranking is equivalent to binary classification



3



2 Ordinal Ranking Setup

#### The Reduction Framework

- Key Ideas
- Important Properties
- Theoretical Usefulness
- Algorithmic Usefulness

#### Experimental Results



Proving New Generalization Theorems



new ordinal ranking theorem = reduction + any cost + bin. thm. + math derivation



3



2 Ordinal Ranking Setup

#### The Reduction Framework

- Key Ideas
- Important Properties
- Theoretical Usefulness
- Algorithmic Usefulness

#### Experimental Results



# **Unifying Existing Algorithms**

ordinal ranking = reduction + cost + binary classification

ordinal ranking	cost	binary classification algorithm
PRank (Crammer and Singer, 2002)	absolute	modified perceptron rule
kernel ranking (Rajaram et al., 2003)	classification	modified hard-margin SVM
SVOR-EXP SVOR-IMC (Chu and Keerthi, 2005)	classification absolute	modified soft-margin SVM modified soft-margin SVM
ORBoost-LR ORBoost-All (Lin and Li, 2006)	classification absolute	modified AdaBoost modified AdaBoost

- development and implementation time could have been saved
- algorithmic structure revealed (SVOR, ORBoost)

# variants of existing algorithms can be designed quickly by tweaking reduction



Hsuan-Tien Lin (CSIE, NTU)

Designing New Algorithms Effortlessly

ordinal ranking = reduction + cost + binary classification

ordinal ranking	cost	binary classification algorithm
RED-SVM	absolute	standard soft-margin SVM
RED-C4.5	absolute	standard C4.5 decision tree
(Li and Lin, 2007)		

SVOR (modified SVM) v.s. RED-SVM (standard SVM):



# advantages of core binary classification algorithm inherited in the new ordinal ranking one



Hsuan-Tien Lin (CSIE, NTU)



- 2 Ordinal Ranking Setup
- 3) The Reduction Framework
  - Key Ideas
  - Important Properties
  - Theoretical Usefulness
  - Algorithmic Usefulness

#### Experimental Results



**Experimental Results** 

## Reduction-C4.5 v.s. SVOR



Experimental Results

### Reduction-SVM v.s. SVOR





- 2 Ordinal Ranking Setup
- 3) The Reduction Framework
  - Key Ideas
  - Important Properties
  - Theoretical Usefulness
  - Algorithmic Usefulness
- Experimental Results



#### • reduction framework: simple but useful

- establish equivalence to binary classification
- unify existing algorithms
- simplify design of new algorithms
- facilitate derivation of new theoretical guarantees
- superior experimental results:

better performance and faster training time

# reduction keeps ordinal ranking up-to-date with binary classification

