Slides adapted from Prof Carpuat and Duraiswami



#### CMSC 422 Introduction to Machine Learning Lecture 10 Multiclass Classification and Reductions

Furong Huang / furongh@cs.umd.edu





# Given an arbitrary method for binary classification, how can we learn to make multiclass predictions?

## Fundamental ML concept: reductions



# One Example of Reduction: Learning with Imbalanced Data

#### TASK: $\alpha$ -Weighted Binary Classification

Given:

- 1. An input space  $\mathcal{X}$
- 2. An unknown distribution  $\mathcal{D}$  over  $\mathcal{X} \times \{-1, +1\}$

*Compute:* A function *f* minimizing:  $\mathbb{E}_{(x,y)\sim \mathcal{D}} \left[ \alpha^{y=1} \left[ f(x) \neq y \right] \right]$ 

#### **Subsampling Optimality Theorem:**

If the binary classifier achieves a binary error rate of  $\epsilon$ , then the error rate of the  $\alpha$ -weighted classifier is  $\alpha \epsilon$ 



# Another Example of Reduction: Multiclass Classification

#### TASK: MULTICLASS CLASSIFICATION

Given:

- 1. An input space X and number of classes K
- 2. An unknown distribution  $\mathcal{D}$  over  $\mathcal{X} \times [K]$

*Compute:* A function *f* minimizing:  $\mathbb{E}_{(x,y)\sim\mathcal{D}}[f(x)\neq y]$ 



# **Reduction 1: OVA**

"One versus all" (aka "one versus rest")

- Train K-many binary classifiers
- classifier k predicts whether an example belong to class k or not
- At test time,
  - If only one classifier predicts positive, predict that class
  - Break ties randomly



#### Algorithm 12 ONEVERSUSALLTRAIN(D<sup>multiclass</sup>, BINARYTRAIN)

- 1: for *i* = 1 to *K* do
- 2:  $\mathbf{D}^{bin} \leftarrow \text{relabel } \mathbf{D}^{multiclass} \text{ so class } i \text{ is positive and } \neg i \text{ is negative}$
- $f_i \leftarrow \text{BINARYTRAIN}(\mathbf{D}^{bin})$
- 4: end for
- 5: **return**  $f_1, \ldots, f_K$

#### Algorithm 13 ONEVERSUSALLTEST $(f_1, \ldots, f_K, \hat{x})$

- 1: *score*  $\leftarrow \langle 0, 0, \dots, 0 \rangle$
- 2: for i = 1 to K do
- $y \leftarrow f_i(\hat{x})$
- $_{4:} \quad score_i \leftarrow score_i + y$
- 5: end for
- 6: return argmax<sub>k</sub> score<sub>k</sub>

// initialize *K*-many scores to zero



# **Time complexity**

- Suppose you have N training examples, in K classes. How long does it take to train an OVA classifier
  - if the base binary classifier takes O(N) time to learn?
  - if the base binary classifier takes O(N<sup>2</sup>) time to learn?



## **Error bound**

- Theorem: Suppose that the average error of the K binary classifiers is ε, then the error rate of the OVA multiclass classifier is at most (K-1) ε
- To prove this: how do different errors affect the maximum ratio of the probability of a multiclass error to the number of binary errors ("efficiency")?



# **Error bound proof**

 If we have a false negative on one of the binary classifiers (assuming all other classifiers correctly output negative)
 What is the probability that we will make an incorrect multiclass prediction?

(K - 1) / K

# Efficiency: [(K - 1) / K] / 1 = (K - 1) / K



# **Error bound proof**

# 2. If we have m false positives with the binary classifiersWhat is the probability that we will make an incorrect multiclass prediction?

If there is also a false negative: 1

Efficiency = 1 / (m + 1)

Otherwise m/(m+1)

Efficiency = [m / (m + 1)] / m = 1 / (m + 1)

MARYLAND

# **Error bound proof**

## 3. What is the worst case scenario?

False negative case: efficiency is (K-1)/K Larger than false positive efficiencies

There are K-many opportunities to get false negative, **overall error bound is (K-1)** ε





#### All versus all (aka all pairs)

# How many binary classifiers does this require?



Algorithm 14 ALLVERSUSALLTRAIN(D<sup>multiclass</sup>, BINARYTRAIN)

1: 
$$f_{ij} \leftarrow \emptyset, \forall 1 \leq i < j \leq K$$
  
2: for  $i = 1$  to  $K$ -1 do  
3:  $\mathbf{D}^{pos} \leftarrow \text{all } \mathbf{x} \in \mathbf{D}^{multiclass} \text{ labeled } i$   
4: for  $j = i+1$  to  $K$  do  
5:  $\mathbf{D}^{neg} \leftarrow \text{all } \mathbf{x} \in \mathbf{D}^{multiclass} \text{ labeled } j$   
6:  $\mathbf{D}^{bin} \leftarrow \{(\mathbf{x}, +1) : \mathbf{x} \in \mathbf{D}^{pos}\} \cup \{(\mathbf{x}, -1) : \mathbf{x} \in \mathbf{D}^{neg}\}$   
7:  $f_{ij} \leftarrow \mathbf{BINARYTRAIN}(\mathbf{D}^{bin})$   
8: end for  
9: end for  
10: return all  $f_{ij}$ s

#### Algorithm 15 ALLVERSUSALLTEST(all $f_{ij}$ , $\hat{x}$ )

1: $score \leftarrow \langle 0, 0, \dots, 0 \rangle$	// initialize K-many scores to zero	
<sup>2:</sup> for <i>i</i> = 1 to <i>K</i> -1 do		
$_{3:}$ for $j = i + 1$ to K do		
$_{4:} \qquad y \leftarrow f_{ij}(\hat{x})$		
5: $score_i \leftarrow score_i + y$		
6: $score_j \leftarrow score_j - y$		
7: end for		
8: end for		
9: <b>return</b> argmax <sub>k</sub> score <sub>k</sub>		

# **Time complexity**

- Suppose you have N training examples, in K classes. How long does it take to train an AVA classifier
  - if the base binary classifier takes O(N) time to learn?
  - if the base binary classifier takes O(N<sup>2</sup>) time to learn?





**Theorem:** Suppose that the average error of the K binary classifiers is  $\varepsilon$ , then the error rate of the AVA multiclass classifier is at most 2(K-1)  $\varepsilon$ 

Question: Does this mean that AVA is always worse than OVA?





- Divide and conquer
  - Organize classes into binary tree structures
- Use confidence to weight predictions of binary classifiers
  - Instead of using majority vote



## **Topics**

Given an arbitrary method for binary classification, how can we learn to make multiclass predictions?
OVA, AVA

Fundamental ML concept: reductions





#### Canonical example: web search

# Given all the documents on the web For a user query, retrieve relevant documents, ranked from most relevant to least relevant



# How can we reduce ranking to binary classification?



## **Preference function**

- Given a query q and documents d<sub>i</sub> and d<sub>j</sub>, the preference function outputs whether
  - d<sub>i</sub> should be preferred to d<sub>i</sub>
  - Or d<sub>i</sub> should be preferred to d<sub>i</sub>
- That's a binary classification problem!



# **Specifying the reduction from ranking to binary classification**

 How to train classifier that predicts preferences?

 How to turn the predicted preferences into a ranking?



**Algorithm 16** NAIVERANKTRAIN(*RankingData*, BINARYTRAIN)



#### **Algorithm 17** NAIVERANKTEST( $f, \hat{x}$ )

- 1: Score  $\leftarrow \langle 0, 0, \ldots, 0 \rangle$
- 2: for all i, j = 1 to M and  $i \neq j$  do
- 3:  $y \leftarrow f(\hat{x}_{ij})$
- $_{4:} \quad score_i \leftarrow score_i + y$
- 5:  $score_j \leftarrow score_j y$
- 6: end for
- 7: **return Argsort**(*score*)

// initialize *M*-many scores to zero

// get predicted ranking of i and j

// return queries sorted by score

Naïve approach

## Works well for bipartite problems

#### "is this document relevant or not?"

# Not ideal for full ranking problems,

#### because

Binary preference problems are not all equally important

Separates preference function and sorting



#### Improving on naïve approach

#### TASK: $\omega$ -RANKING

Given:

- 1. An input space  $\mathcal{X}$
- 2. An unknown distribution  $\mathcal{D}$  over  $\mathcal{X} \times \Sigma_M$

*Compute:* A function  $f : \mathcal{X} \to \Sigma_M$  minimizing:

$$\mathbb{E}_{(\boldsymbol{x},\sigma)\sim\mathcal{D}}\left[\sum_{u\neq v} [\sigma_u < \sigma_v] \left[\hat{\sigma}_v < \hat{\sigma}_u\right] \,\omega(\sigma_u, \sigma_v)\right]$$
(5.7)  
where  $\hat{\sigma} = f(\boldsymbol{x})$ 

MARYLAND

#### **Example of cost functions**

# $\omega(i,j) = \begin{cases} 1 & \text{if } \min\{i,j\} \le K \text{ and } i \ne j \\ 0 & \text{otherwise} \end{cases}$



# **Resulting Ranking Algorithms**

#### Algorithm 18 RANKTRAIN( $D^{rank}$ , $\omega$ , BINARYTRAIN)

- $\mathbf{D}^{bin} \leftarrow [ ]$
- 2: for all  $(x, \sigma) \in \mathbf{D}^{rank}$  do
- $_{3:}$  for all  $u \neq v$  do
- $_{4:} \qquad y \leftarrow \mathbf{SIGN}(\sigma_{v} \sigma_{u})$

5: 
$$w \leftarrow \omega(\sigma_u, \sigma_v)$$

6: 
$$\mathbf{D}^{bin} \leftarrow \mathbf{D}^{bin} \oplus (y, w, x_{uv})$$

- 7: end for
- 8: end for
- 9: return BINARYTRAIN $(\mathbf{D}^{bin})$

// y is +1 if u is prefered to v // w is the cost of misclassification



#### **Algorithm 19 RANKTEST**(f, $\hat{x}$ , obj)

- <sup>11</sup> if *obj* contains 0 or 1 elements then
- 2: return *obj*

#### 3: **else**

4:	$p \leftarrow$ randomly chosen object in $o$	<i>bj</i> // pick pivot	
5:	$left \leftarrow []$	// elements that seem smaller than $p$	
6:	$right \leftarrow []$	// elements that seem larger than $p$	
7:	for all $u \in obj \setminus \{p\}$ do		
8:	$\hat{y} \leftarrow f(x_{up})$	// what is the probability that $u$ precedes $p$	
9:	if uniform random variable $<$	$\hat{y}$ then	
10:	$left \leftarrow left \oplus u$		
11:	else		
12:	$right \leftarrow right \oplus u$		
13:	end if		
14:	end for		
15:	$left \leftarrow \text{RankTest}(f, \hat{x}, left)$	// sort earlier elements	
16:	$right \leftarrow RankTest(f, \hat{x}, right)$	// sort later elements	
17:	<b>return</b> <i>left</i> $\oplus \langle p \rangle \oplus right$		
18: end if			

## RankTest

A probabilistic version of the quicksort algorithm

Only O(Mlog<sub>2</sub>M) calls to f in expectation

Better error bound than naïve algorithm (see CIML for theorem)



# What you should know

- What are reductions and why they are useful
- Implement, analyze and prove error bounds of algorithms for
  - Weighted binary classification
  - Multiclass classification (OVA, AVA)
- Understand algorithms for
  - $\omega$  –ranking



Furong Huang 3251 A.V. Williams, College Park, MD 20740 301.405.8010 / furongh@cs.umd.edu