

# Beyond Binary Classification: Reductions

CMSC 422

MARINE CARPUAT

[marine@cs.umd.edu](mailto:marine@cs.umd.edu)

# Topics

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:  $\mathbb{E}_{(x,y) \sim \mathcal{D}} [\alpha^{y=1} [f(x) \neq y]]$

## **Subsampling Optimality Theorem:**

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

# Multiclass classification

- Real world problems often have multiple classes (text, speech, image, biological sequences...)
- How can we perform multiclass classification?
  - Straightforward with decision trees or KNN
  - Can we use the perceptron algorithm?

# Today: Reductions for Multiclass Classification

## TASK: MULTICLASS CLASSIFICATION

*Given:*

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

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

## TASK: 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:  $\mathbb{E}_{(x,y) \sim \mathcal{D}} [f(x) \neq y]$

# How many classes can we handle in practice?

- In most tasks, number of classes  $K < 100$
- For much larger  $K$ 
  - we need to frame the problem differently
  - e.g, machine translation or automatic speech recognition

# 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( $\mathbf{D}^{multiclass}$ , BINARYTRAIN)

---

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

---

---

**Algorithm 13** ONEVERSUSALLTEST( $f_1, \dots, f_K, \hat{x}$ )

---

```
1:  $score \leftarrow \langle 0, 0, \dots, 0 \rangle$  // initialize  $K$ -many scores to zero
2: for  $i = 1$  to  $K$  do
3:    $y \leftarrow f_i(\hat{x})$ 
4:    $score_i \leftarrow score_i + y$ 
5: end for
6: return  $\operatorname{argmax}_k score_k$ 
```

---

# 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  $O(N^2)$  time to learn?

# Error bound

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

- If we have  $p$  **false positives** with the binary classifiers
- What is the probability that we will make an incorrect multiclass prediction?
  - If there is also a false negative: 1
    - Efficiency =  $1 / (p + 1)$
  - If there is no false negative:  $p / (p + 1)$ 
    - Efficiency =  $p / (p + 1) / p = 1 / (p + 1)$

# Error bound proof

- 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  $(K-1) \epsilon$**

## Reduction 2: AVA

- All versus all (aka all pairs)
- How many binary classifiers does this require?

---

**Algorithm 14** ALLVERSUSALLTRAIN( $\mathbf{D}^{multiclass}$ , BINARYTRAIN)

---

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

---

---

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

---

```
1:  $score \leftarrow \langle 0, 0, \dots, 0 \rangle$  // initialize  $K$ -many scores to zero
2: for  $i = 1$  to  $K-1$  do
3:   for  $j = i+1$  to  $K$  do
4:      $y \leftarrow f_{ij}(\hat{\mathbf{x}})$ 
5:      $score_i \leftarrow score_i + y$ 
6:      $score_j \leftarrow score_j - y$ 
7:   end for
8: end for
9: return  $\operatorname{argmax}_k score_k$ 
```

---



# 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  $O(N^2)$  time to learn?

# Error bound

- **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?

# Extensions

- 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

# Ranking

- 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  $d_i$  and  $d_j$ , the preference function outputs whether
  - $d_i$  should be preferred to  $d_j$
  - Or  $d_j$  should be preferred to  $d_i$
- 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)

---

```
1:  $\mathbf{D} \leftarrow []$ 
2: for  $n = 1$  to  $N$  do
3:   for all  $i, j = 1$  to  $M$  and  $i \neq j$  do
4:     if  $i$  is preferred to  $j$  on query  $n$  then
5:        $\mathbf{D} \leftarrow \mathbf{D} \oplus (\mathbf{x}_{nij}, +1)$ 
6:     else if  $j$  is preferred to  $i$  on query  $n$  then
7:        $\mathbf{D} \leftarrow \mathbf{D} \oplus (\mathbf{x}_{nij}, -1)$ 
8:     end if
9:   end for
10: end for
11: return BINARYTRAIN( $\mathbf{D}$ )
```

Features associated with comparing document  $j$  and document  $i$  for query  $n$

---

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

---

```
1:  $score \leftarrow \langle 0, 0, \dots, 0 \rangle$  // initialize  $M$ -many scores to zero
2: for all  $i, j = 1$  to  $M$  and  $i \neq j$  do
3:    $y \leftarrow f(\hat{\mathbf{x}}_{ij})$  // get predicted ranking of  $i$  and  $j$ 
4:    $score_i \leftarrow score_i + y$ 
5:    $score_j \leftarrow score_j - y$ 
6: end for
7: return ARGSORT( $score$ ) // 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  $f : \mathcal{X} \rightarrow \Sigma_M$  minimizing:

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

where  $\hat{\sigma} = f(\mathbf{x})$

# Example of cost functions

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

# Resulting Ranking Algorithms

---

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

---

```
1:  $\mathbf{D}^{bin} \leftarrow [ ]$ 
2: for all  $(\mathbf{x}, \sigma) \in \mathbf{D}^{rank}$  do
3:   for all  $u \neq v$  do
4:      $y \leftarrow \text{SIGN}(\sigma_v - \sigma_u)$  //  $y$  is +1 if  $u$  is preferred to  $v$ 
5:      $w \leftarrow \omega(\sigma_u, \sigma_v)$  //  $w$  is the cost of misclassification
6:      $\mathbf{D}^{bin} \leftarrow \mathbf{D}^{bin} \oplus (y, w, \mathbf{x}_{uv})$ 
7:   end for
8: end for
9: return BINARYTRAIN( $\mathbf{D}^{bin}$ )
```

---

---

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

---

```
1: if  $obj$  contains 0 or 1 elements then
2:   return  $obj$ 
3: else
4:    $p \leftarrow$  randomly chosen object in  $obj$  // 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:    $left \leftarrow$  RANKTEST( $f, \hat{x}, left$ ) // sort earlier elements
16:    $right \leftarrow$  RANKTEST( $f, \hat{x}, right$ ) // sort later elements
17:   return  $left \oplus \langle p \rangle \oplus right$ 
18: end if
```

---

- RankTest

- A probabilistic version of the quicksort algorithm

- Only  $O(M \log_2 M)$  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

\*for weighted binary classification and OVA only