

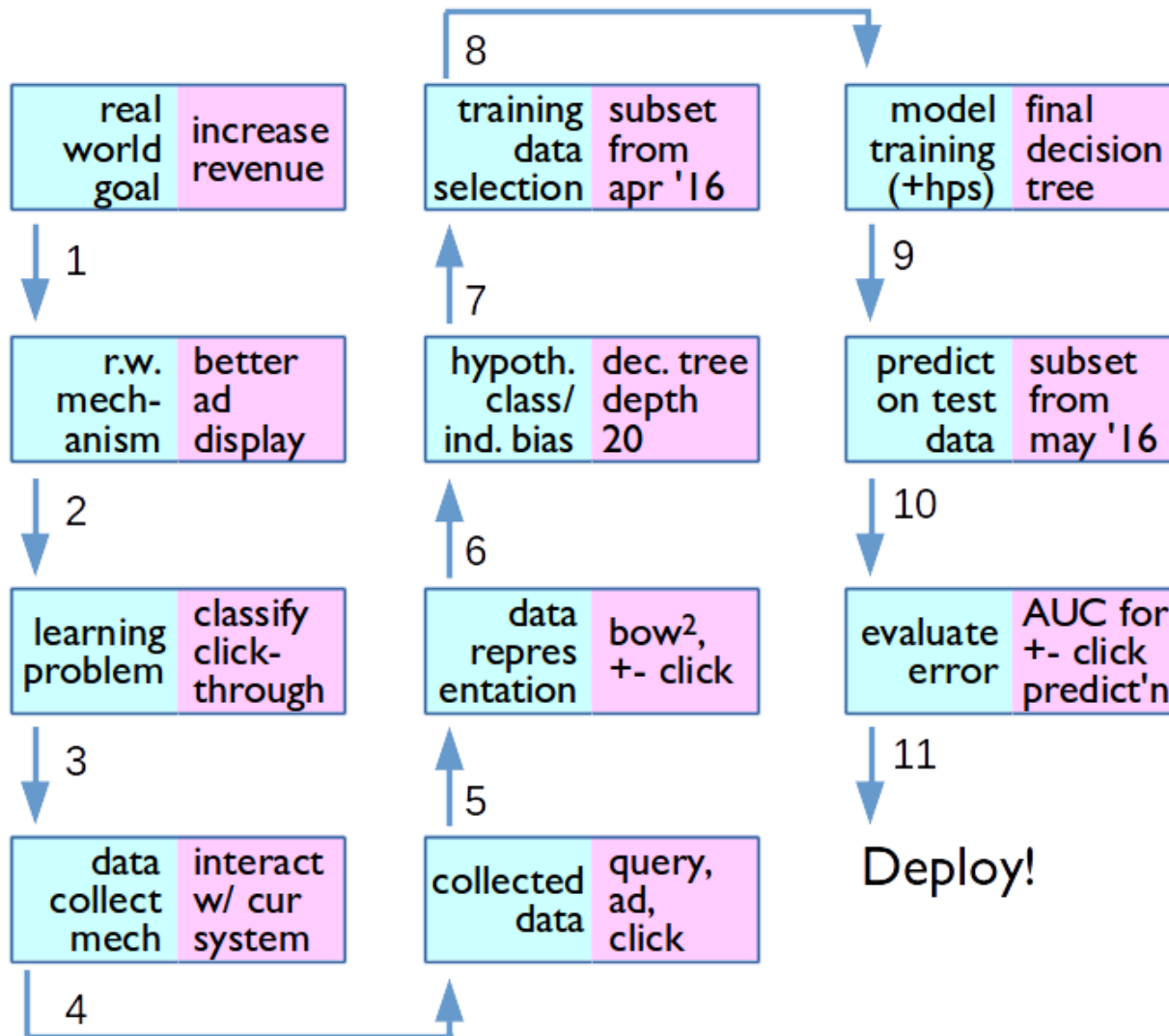
# Imbalanced Data and Reductions

CMSC 422

MARINE CARPUAT

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

# Typical Design Process for an ML Application



# Imbalanced data distributions

- Sometimes training examples are drawn from an imbalanced distribution
- This results in an imbalanced training set
  - “needle in a haystack” problems
  - E.g., find fraudulent transactions in credit card histories
- Why is this a big problem for the ML algorithms we know?

# Recall: Machine Learning as Function Approximation

## Problem setting

- Set of possible instances  $X$
- Unknown target function  $f: X \rightarrow Y$
- Set of function hypotheses  $H = \{h \mid h: X \rightarrow Y\}$

## Input

- Training examples  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$  of unknown target function  $f$

## Output

- Hypothesis  $h \in H$  that best approximates target function  $f$

# Recall: Loss Function

$l(y, f(x))$  where  $y$  is the truth and  $f(x)$  is the system's prediction

$$\text{e.g. } l(y, f(x)) = \begin{cases} 0 & \text{if } y = f(x) \\ 1 & \text{otherwise} \end{cases}$$

Captures our notion of what is important to learn

# Recall: Expected loss

- $f$  should make good predictions
  - as measured by loss  $l$
  - on **future** examples that are also drawn from  $D$
- Formally
  - $\varepsilon$ , the expected loss of  $f$  over  $D$  with respect to  $l$  should be small

$$\varepsilon \triangleq \mathbb{E}_{(x,y) \sim D} \{l(y, f(x))\} = \sum_{(x,y)} D(x, y) l(y, f(x))$$

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

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

Given a good algorithm for solving the binary classification problem, how can we solve the  $\alpha$ -weighted binary classification problem?

We define cost of misprediction as:  
 $\alpha > 1$  for  $y = +1$   
1 for  $y = -1$



# Solution: Train a binary classifier on an induced distribution

---

**Algorithm 11** SUBSAMPLEMAP( $\mathcal{D}^{\text{weighted}}, \alpha$ )

---

```
1: while true do  
2:    $(\mathbf{x}, y) \sim \mathcal{D}^{\text{weighted}}$            // draw an example from the weighted distribution  
3:    $u \sim$  uniform random variable in  $[0, 1]$   
4:   if  $y = +1$  or  $u < \frac{1}{\alpha}$  then  
5:     return  $(\mathbf{x}, y)$   
6:   end if  
7: end while
```

---

# Subsampling optimality

- **Theorem:** If the binary classifier achieves a binary error rate of  $\varepsilon$ , then the error rate of the  $\alpha$ -weighted classifier is  $\alpha \varepsilon$
- Let's prove it.  
(see also CIML 6.1)

# Strategies for inducing a new binary distribution

- Undersample the negative class
- Oversample the positive class

# Strategies for inducing a new binary distribution

- Undersample the negative class
  - More computationally efficient
- Oversample the positive class
  - Base binary classifier might do better with more training examples
  - Efficient implementations incorporate weight in algorithm, instead of explicitly duplicating data!

---

**Algorithm 1** DECISIONTREETRAIN(*data*, *remaining features*)

---

```
1: guess ← most frequent answer in data // default answer for this data
2: if the labels in data are unambiguous then
3:   return LEAF(guess) // base case: no need to split further
4: else if remaining features is empty then
5:   return LEAF(guess) // base case: cannot split further
6: else // we need to query more features
7:   for all  $f \in \textit{remaining features}$  do
8:     NO ← the subset of data on which  $f=no$ 
9:     YES ← the subset of data on which  $f=yes$ 
10:    score[f] ← # of majority vote answers in NO
11:                + # of majority vote answers in YES
// the accuracy we would get if we only queried on f
12:   end for
13:   f ← the feature with maximal score(f)
14:   NO ← the subset of data on which  $f=no$ 
15:   YES ← the subset of data on which  $f=yes$ 
16:   left ← DECISIONTREETRAIN(NO, remaining features \ {f})
17:   right ← DECISIONTREETRAIN(YES, remaining features \ {f})
18:   return NODE(f, left, right)
19: end if
```

---

# Reductions

- Idea is to re-use simple and efficient algorithms for binary classification to perform more complex tasks
- Works great in practice:
  - E.g., [Vowpal Wabbit](#)

# Learning with Imbalanced Data is an Example of Reduction

## 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?



# 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

# What you should know

- How can we take the standard binary classifier and adapt it to handle problems with
  - Imbalanced data distributions
  - Multiclass classification problems
- Algorithms & guarantees on error rate
- Fundamental ML concept: reduction

# 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  $k$  **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 =  $1 / (k + 1)$
  - Otherwise  $k / (k + 1)$ 
    - Efficiency =  $k / (k + 1) / k = 1 / (k + 1)$

# 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