

Imbalanced Data and Reductions

CMSC 422

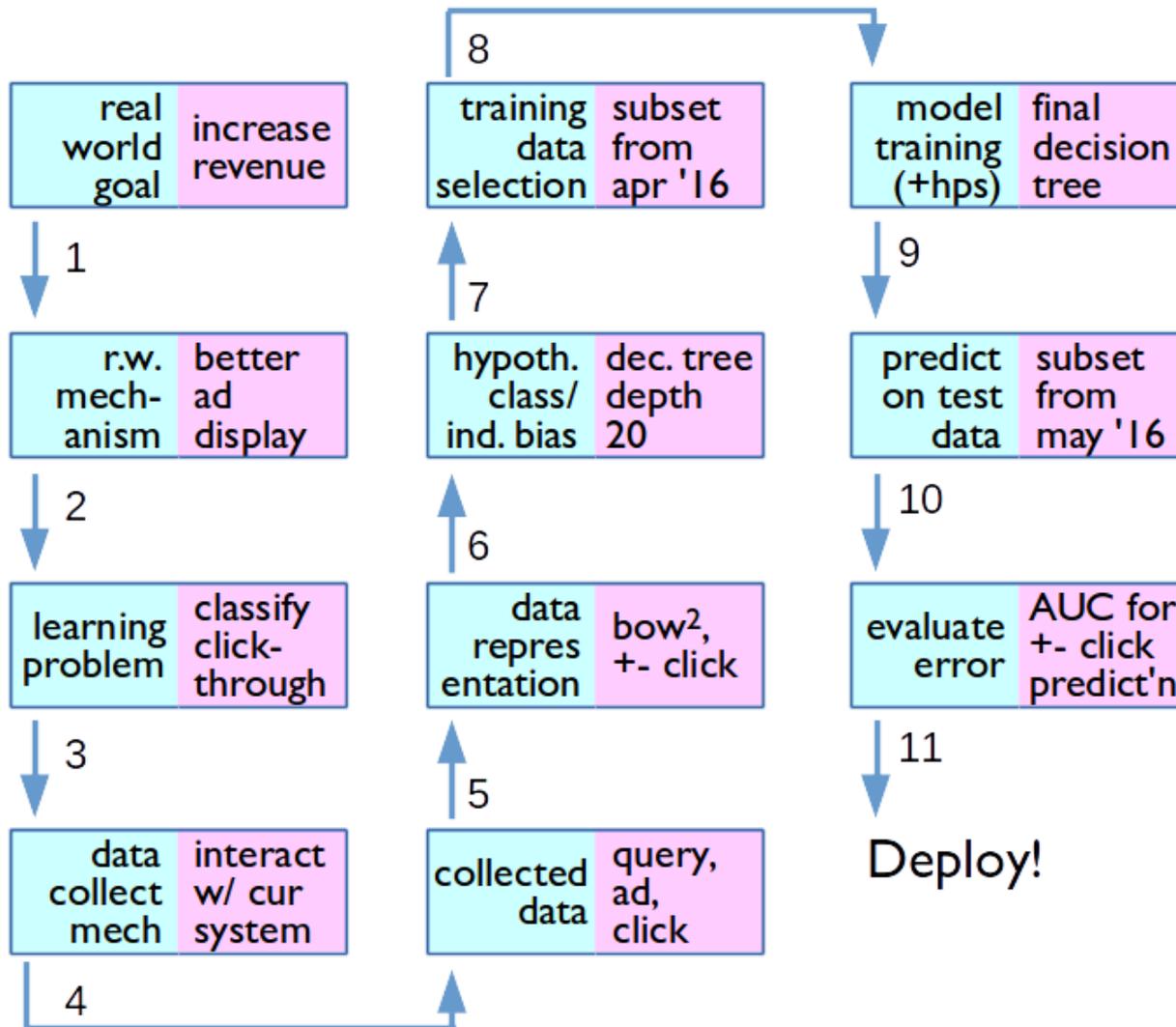
MARINE CARPUAT

marine@cs.umd.edu

Today's topics

- Using standard binary classifiers to solve other problems
 - Weighted classification
 - Multiclass classification
- Algorithms & guarantees on error rate
- Fundamental ML concept: reduction

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: Expected loss

- f should make good predictions
 - as measured by loss l
 - on **future** examples that are also drawn from D
- Formally
 - ε , 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: α -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 α -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 ε , then the error rate of the α -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  $\leftarrow$  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  $f \in$  remaining features do
8:     NO  $\leftarrow$  the subset of data on which  $f=no$ 
9:     YES  $\leftarrow$  the subset of data on which  $f=yes$ 
10:    score[ $f$ ]  $\leftarrow$  # 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:    $f \leftarrow$  the feature with maximal score( $f$ )
14:   NO  $\leftarrow$  the subset of data on which  $f=no$ 
15:   YES  $\leftarrow$  the subset of data on which  $f=yes$ 
16:   left  $\leftarrow$  DECISIONTREETRAIN(NO, remaining features  $\setminus$  { $f$ })
17:   right  $\leftarrow$  DECISIONTREETRAIN(YES, remaining features  $\setminus$  { $f$ })
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: α -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 ε , then the error rate of the α -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

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?

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?

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