

Aggressive Double Sampling for Reducing Multi-class Classification to Binary Classification

Bikash Joshi (PhD Student)
AMA team, LIG

Supervised By:

Prof. Massih-Reza Amini and Dr. Franck Iutzeler

March 20, 2017

Outline

- 1 Introduction
- 2 Multiclass to Binary Reduction
- 3 Double-Sampled Multiclass to Binary Reduction
- 4 Experimental Results
- 5 Conclusion

Outline

- 1 Introduction
- 2 Multiclass to Binary Reduction
- 3 Double-Sampled Multiclass to Binary Reduction
- 4 Experimental Results
- 5 Conclusion

Multiclass Classification: Introduction

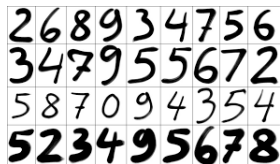


Figure : Digit Classification

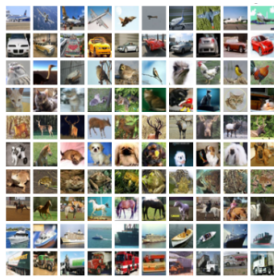


Figure : Image Classification

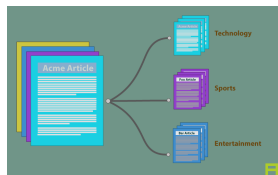


Figure : Text Classification

- Finite set of categories ($K > 2$)
- Popular applications: image and text classification.

Multiclass classification: Related Work

① Combined approaches based on binary classification:

▶ One-Vs-Rest

- ★ One binary problem for each class
- ★ K binary problems
- ★ $O(K \times d)$

▶ One-Vs-One

- ★ One binary problem for each pair of classes
- ★ $O(K^2 \times d)$

② Uncombined Approaches

- ▶ for example: multiclass SVM, MLP
- ▶ One scoring function per class

③ Logarithmic Time Algorithms

- ▶ For example: logTree, Recall-Tree
- ▶ Each leaf node represents a class
- ▶ $O(\log K)$

Multiclass classification : Challenges

- The number of classes, K , in new emerging multiclass problems, for example in text and image classification, may reach 10^5 to 10^6 categories.
- For example:

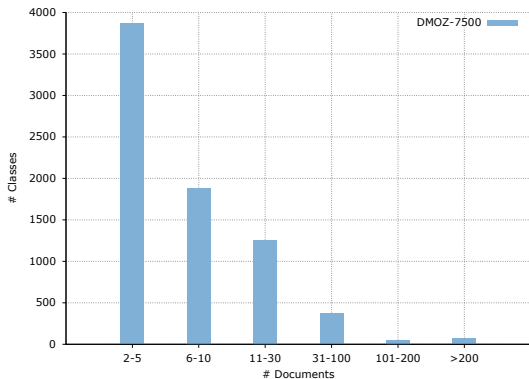


The screenshot shows the DMOZ website interface. At the top, there is a green header with the 'dmoz' logo and 'AOL.' branding. Below the header is a navigation bar with links for 'Follow @dmoz', 'about dmoz', 'dmoz blog', 'suggest URL', 'help', 'link', and 'editor login'. A search bar is present with the text 'Search advanced'. The main content area is a grid of category links, each with a sub-link. The categories include: Arts (Movies, Television, Music...), Business (Jobs, Real Estate, Investing...), Computers (Internet, Software, Hardware...), Games (Video Games, RPGs, Gambling...), Health (Fitness, Medicine, Alternative...), Home (Family, Consumers, Cooking...), Kids and Teens (Arts, School Time, Teen Life...), News (Media, Newspapers, Weather...), Recreation (Travel, Food, Outdoors, Humor...), Reference (Maps, Education, Libraries...), Regional (US, Canada, UK, Europe...), Science (Biology, Psychology, Physics...), Shopping (Clothing, Food, Gifts...), Society (People, Religion, Issues...), Sports (Baseball, Soccer, Basketball...), and World (Català, Český, Dansk, Deutsch, Español, Esperanto, Français, Galego, Hrvatski, Italiano, Lietuvių, Magyar, Nederlands, Norsk, Polski, Português, Română, Slovensky, Suomi, Svenska, Türkçe, Български, Ελληνικά, Русский, Українська, العربية, עברית, עברית, 日本語, 簡體中文, 繁體中文, ...). At the bottom, there is a 'Become an Editor' button and a copyright notice: 'Copyright © 1998-2015 AOL, Inc.' and '4,001,665 sites - 90,684 editors - over 1,026,659 categories'.

- ▶ 4×10^6 sites
- ▶ 10^6 categories
- ▶ 10^5 editors
- ▶ Imbalanced nature of hierarchies

Multiclass classification : Challenges

- Class imbalance problem
- Majority of classes have few representative examples
- Long tailed distribution



Text Classification:

Task: Automatic classification of an example text to one of fixed set of categories.

Feature Representation:

- Bag of Words:
 - ▶ From training corpus extract vocabulary.
 - ▶ Represent each terms as 0 or 1
 - ▶ Highly sparse
- Document-class joint feature representation:
 - ▶ Inspired by learning to rank
 - ▶ Similarity features between an example and class of examples
 - ▶ For example:

$$\sum_{t \in y \cap x} 1$$

Where,

$x \rightarrow$ One document

$y \rightarrow$ Class of documents

Outline

- 1 Introduction
- 2 Multiclass to Binary Reduction**
- 3 Double-Sampled Multiclass to Binary Reduction
- 4 Experimental Results
- 5 Conclusion

Motivation of our work

- Baselines: Model complexity increases with classes(K) and feature dimension (d).
- Algorithm that scales well for large scale data
- Does not suffer from class imbalance problem
- Less complex model
- Competitive with the state of the art approaches

Framework

- $\mathcal{X} \subseteq \mathbb{R}^d$: Input Space
- $\mathcal{Y} = 1, \dots, K$: Output Space
- $S = (x_i^{y_i})_{i=1}^m$: Training set of i.i.d. pairs
- $G = g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$: Class of predictors

Instantaneous Loss

$$e(g, x^y) = \frac{1}{K-1} \sum_{y' \in \mathcal{Y} \setminus y} \mathbb{1}_{g(x^y) \leq g(x^{y'})} \quad (1)$$

- $\mathbb{1}_\pi$ is the indicator function (Value is 0 or 1)
- Average number of classes that get greater scoring by g than true class
- Ranking loss used in Multiclass-SVM ^a

^aWeston et. al. (1998)

Framework

Empirical Loss

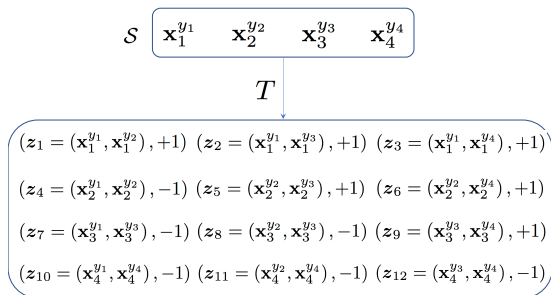
Empirical error of $g \in G$ over S is:

$$L_m(g, S) = \frac{1}{m(K-1)} \sum_{i=1}^m \sum_{y' \in Y \setminus y_i} \mathbb{1}_{g(x_i^{y_i}) \leq g(x_i^{y'})} \quad (2)$$

$$= \frac{1}{m(K-1)} \sum_{i=1}^m \sum_{y' \in Y \setminus y_i} \mathbb{1}_{\underbrace{h(x_i^{y_i}, x_i^{y'})}_{g(x_i^{y_i}) - g(x_i^{y'})} \leq 0} \quad (3)$$

- Resembles to binary-classification-loss based risk
- Selection of a hypothesis in G minimizing risk over S is equivalent to search a hypothesis in H minimizing risk over $T(S)$ of size $m \times (K-1)$

Multiclass to binary reduction example



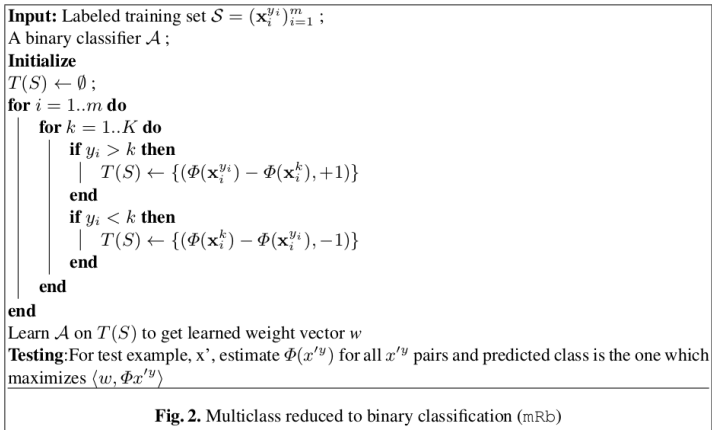
- We consider the following transformation

$$T(S) = \left(\left\{ \begin{array}{ll} (\mathbf{z}_j = (\mathbf{x}_i^k, \mathbf{x}_i^{y_i}), \tilde{y}_j = -1) & \text{if } k < y_i \\ (\mathbf{z}_j = (\mathbf{x}_i^{y_i}, \mathbf{x}_i^k), \tilde{y}_j = +1) & \text{elsewhere} \end{array} \right\}_{j \doteq (i-1)(K-1)+k} \right)$$

- $|T(S)| = m \times (K - 1)$

Multiclass to binary reduction algorithm

[Bikash et al. 2015]



Improvements and New challenges

Improvements:

- One parameter vector for all classes.
- Low-dimensional feature space.
- Overcome class imbalance.

New Challenges:

- Number of transformations huge for larger K
- Large computational overhead
- Large memory requirement

Outline

- 1 Introduction
- 2 Multiclass to Binary Reduction
- 3 Double-Sampled Multiclass to Binary Reduction**
- 4 Experimental Results
- 5 Conclusion

Aggressive double sampling

- 1 Drawing uniformly μ examples per class, in order to form practical set S_{μ} :
 - ▶ Reduce redundancy in examples
 - ▶ Emphasizing rare classes
- 2 For each example \mathbf{x}^y in S_{μ} , drawing uniformly κ adversarial classes in $\mathcal{Y} \setminus \{y\}$.
 - ▶ Reduces time complexity
 - ▶ Low memory requirement

Double Sampled Multi to Binary Reduction

```
Input: Labeled training set  $\mathcal{S} = (\mathbf{x}_i^{y_i})_{i=1}^m$   
Initialize  
 $T_\kappa(\mathcal{S}_\mu) \leftarrow \emptyset$   
 $\mathcal{S}_\mu \leftarrow \emptyset$   
for  $k = 1..K$  do  
  Draw randomly a set  $\underline{\mathcal{S}}$  of  $\mu$  examples of class  $k$  from  $\mathcal{S} \triangleright \mu \ll |\mathcal{S}_k|$   
   $\mathcal{S}_\mu \leftarrow \mathcal{S}_\mu \cup \underline{\mathcal{S}}$   
end  
forall the  $\mathbf{x}^y \in \mathcal{S}_\mu$  do  
  Draw uniformly a set  $\underline{\mathcal{K}}$  of  $\kappa$  classes from  $\mathcal{Y} \setminus \{y\} \triangleright \kappa \ll |K|$   
  forall the  $k \in \underline{\mathcal{K}}$  do  
    if  $k < y$  then  
       $T_\kappa(\mathcal{S}_\mu) \leftarrow T_\kappa(\mathcal{S}_\mu) \cup (\mathbf{z} = (\phi(\mathbf{x}^k), \phi(\mathbf{x}^y)) \quad , \tilde{y} = -1)$   
    end  
    else  
       $T_\kappa(\mathcal{S}_\mu) \leftarrow T_\kappa(\mathcal{S}_\mu) \cup (\mathbf{z} = (\phi(\mathbf{x}^y), \phi(\mathbf{x}^k)) \quad , \tilde{y} = +1)$   
    end  
  end  
end  
return  $T_\kappa(\mathcal{S}_\mu)$ 
```

Outline

- 1 Introduction
- 2 Multiclass to Binary Reduction
- 3 Double-Sampled Multiclass to Binary Reduction
- 4 Experimental Results**
- 5 Conclusion

Experimental Setup

Datasets:

- Application: Text Classification
- DMOZ and Wikipedia datasets. (LSHTC challenge)
- Pre-processed with stop word removal and stemming.
- Random samples of 1000, 2000, 3000, 4000, 5000, 7500, 10000, 20000.

Comparison:

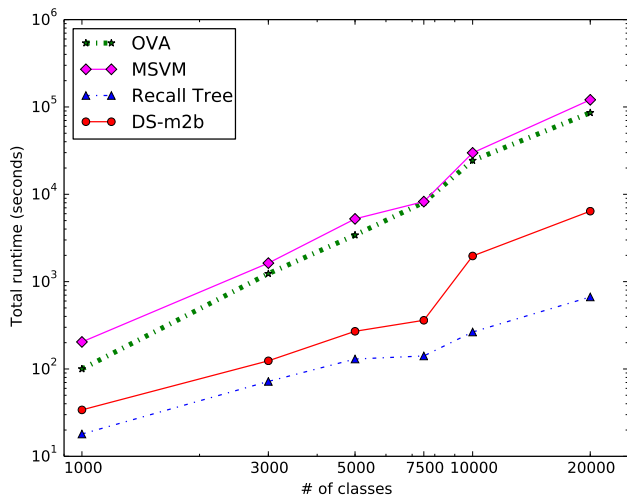
- DS-m2b: Proposed double sampled multiclass to binary algorithm
- OVA: One-Vs-All algorithm
- M-SVM: Crammar-Singer implementation of multiclass SVM
- Recall Tree: Hierarchical One-Vs-Some algorithm

Feature representation $\Phi(x^y)$

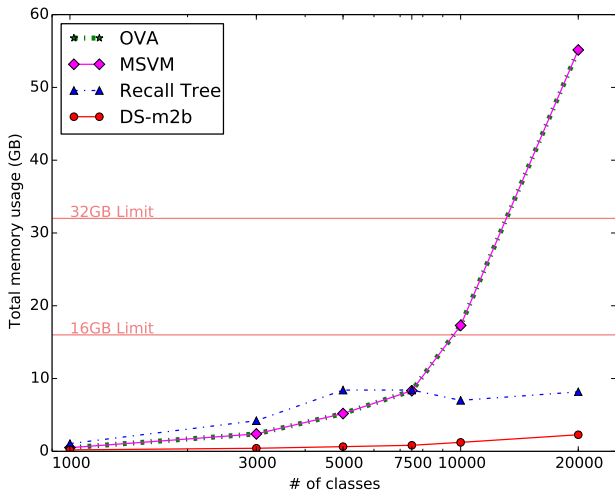
Features	
1. $\sum_{t \in Y \cap X} \ln(1 + y_t)$	2. $\sum_{t \in Y \cap X} \ln(1 + \frac{l_S}{S_t})$
3. $\sum_{t \in Y \cap X} l_t$	4. $\sum_{t \in Y \cap X} \ln(1 + \frac{y_t}{ y })$
5. $\sum_{t \in Y \cap X} \ln(1 + \frac{y_t}{ y } \cdot l_t)$	6. $\sum_{t \in Y \cap X} \ln(1 + \frac{y_t}{ y } \cdot \frac{l_S}{S_t})$
7. $\sum_{t \in Y \cap X} 1$	8. $\sum_{t \in Y \cap X} \frac{y_t}{ y } \cdot l_t$
9. BM25	10. $d(x^y, \text{centroid}(y))$

- x_t : number of occurrences of terme t in document x ,
- \mathcal{V} : Number of distinct terms in \mathcal{S} ,
- $y_t = \sum_{x \in Y} x_t$, $|y| = \sum_{t \in \mathcal{V}} y_t$, $S_t = \sum_{x \in \mathcal{S}} x_t$, $l_S = \sum_{t \in \mathcal{V}} S_t$.
- l_t : idf of the term t ,

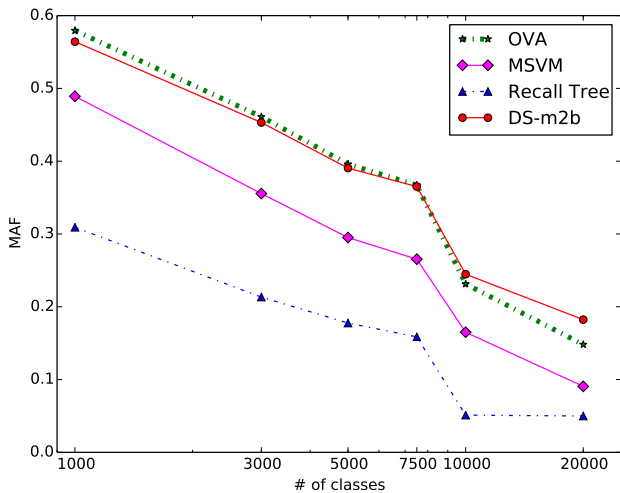
Results: Runtime Comparison



Results: Memory Comparison



Results: Prediction Performance Comparison



Outline

- 1 Introduction
- 2 Multiclass to Binary Reduction
- 3 Double-Sampled Multiclass to Binary Reduction
- 4 Experimental Results
- 5 Conclusion**

Conclusion:

- Multiclass to binary reduction to handle large-class scenario and overcome class imbalance problem.
- Use of double sampling to further improve computational complexity and memory usage.

Questions?