

`multiboost.org`: an implementation of
cost-sensitive multi-class/multi-label AdaBoost.MH

Róbert Busa-Fekete¹, Norman Casagrande², and Balázs Kégl¹

¹University of Paris-Sud / CNRS

²last.fm

ICML'10 open source workshop

June 25, 2010

reviewer 1: “using the program is marginally faster and much less instructive than writing your own version of these **simple** algorithms”

- Simple?
 - **yes**: **binary AdaBoost with decision stumps is simple**, but an average biologist will never sit down to implement it
 - **no**: implementing full **cost-sensitive multi-class/multi-label AdaBoost.MH**¹ with decision trees/products/Haar filters is **not simple** even for an average computer scientist
 - WEKA contains a **very suboptimal** multi-class AdaBoost implementation, other than that there is no widely used implementation that we know of

¹Schapire, R. E. and Singer, Y. Improved boosting algorithms using confidence-rated predictions. *Mach. Learn.*, 37(3):297–336, 1999.

INDICATORBASE(X, Y, W)

```

1   for j ← 1 to d           ▷ all (nominal) features
2       (αj, vj, uj) ← BESTINDICATOR(x(j), Y, W, I(j)) ▷ x(j) ≜ (x1(j), ..., xn(j))
3   j* ← argminj E(αj vj φj, uj, W)
4   return (αj*, vj*, φj*, uj*(·))

```

BESTINDICATOR(x, Y, W, I)

```

1   for ℓ ∈ I
2       for ℓ ← 1 to K
3           γℓ, ℓ+ ← γℓ, ℓ- ← 0
4           uℓ ← RANDOM(±1)
5   for i ← 1 to n for ℓ ← 1 to K
6       if wi, ℓ yi, ℓ > 0 then
7           γℓ, ℓ+ ← γℓ, ℓ+ + wi, ℓ yi, ℓ
8       else
9           γℓ, ℓ- ← γℓ, ℓ- - wi, ℓ yi, ℓ
10  α ← 0, v ← 0
11  while TRUE
12      αprev ← α, vprev ← v           ▷ save current optimal α and v
13      for ℓ ← 1 to K
14          vℓ ← sign ( ∑i ∈ I (γℓ, ℓ+ - γℓ, ℓ-) uℓ ) or vℓ ← ½ ln ( ∑i ∈ I (γℓ, ℓ+ I{ui > 0} + γℓ, ℓ- I{ui < 0}) / ∑i ∈ I (γℓ, ℓ- I{ui > 0} + γℓ, ℓ+ I{ui < 0}) )
15      α ← ½ ln ( ∑ℓ=1K ∑i ∈ I (γℓ, ℓ+ I{ui > 0} + γℓ, ℓ- I{ui < 0}) ) or α ← -1
16      if E(α v φu, W) ≥ E(αprev vprev φu, W) then
17          return (αprev, vprev, u)
18      αprev ← α, uprev ← u           ▷ save current optimal α and u
19      for ℓ ∈ I
20          uℓ ← sign ( ∑i=1K (γℓ, ℓ+ - γℓ, ℓ-) vℓ ) or uℓ ← ½ ln ( ∑ℓ=1K (γℓ, ℓ+ I{vℓ > 0} + γℓ, ℓ- I{vℓ < 0}) / ∑ℓ=1K (γℓ, ℓ- I{vℓ > 0} + γℓ, ℓ+ I{vℓ < 0}) )
21      α ← ½ ln ( ∑ℓ ∈ I ∑ℓ=1K (γℓ, ℓ+ I{vℓ > 0} + γℓ, ℓ- I{vℓ < 0}) ) or α ← -1
22      if E(α v φu, W) ≥ E(αprev v φuprev, W) then
23          return (αprev, v, uprev)

```

Cost sensitive multi-label/multi-class

- **Input** vector $\mathbf{x} \in \mathcal{X}^d$: $x^{(j)}$ is either **real-valued** or **nominal**
- **Label** vector $\mathbf{y} \in \{-1, 1\}^K$
- **Cost** (initial weight) vector $\mathbf{w}^{(1)} \in (\mathbb{R}^+ \cup \{0\})^K$
 - for example, **classical multi-class**

$$w_{i,\ell}^{(1)} = \begin{cases} \frac{1}{2n} & \text{if } \ell \text{ is the correct class of } \mathbf{x}_i \text{ (if } y_{i,\ell} = 1), \\ \frac{1}{2n(K-1)} & \text{otherwise (if } y_{i,\ell} = -1). \end{cases}$$

Cost sensitive multi-label/multi-class

- AdaBoost.MH learns a **vector-valued discriminant** function

$$\mathbf{f}(\mathbf{x}) : \mathcal{X}^d \rightarrow \mathbb{R}^K$$

by minimizing the **weighted Hamming loss**

$$R(\mathbf{f}^{(T)}, \mathbf{w}^{(1)}) = \sum_{i=1}^n \sum_{\ell=1}^K w_{i,\ell}^{(1)} \mathbb{I} \left\{ f^{(\ell)}(\mathbf{x}_i) y_{i,\ell} < 0 \right\}$$

- Strong learners
 - AdaBoost.MH²
 - FeatureBoost³
- Weak learners
 - **Decision stump** for real-valued features
 - **Selector** and **subset indicator**⁴ for nominal features
 - **Haar filter**⁵ for image input
 - Decision **tree** and decision **product**⁴ for combining simple base-classifiers
 - **Easy to add new base learners** without affecting the main boosting engine

²Schapire, R. E. and Singer, Y. Improved boosting algorithms using confidence-rated predictions. *Mach. Learn.*, 37(3):297–336, 1999.

³Bradley, J.K. and Schapire, R.E. FilterBoost: Regression and classification on large datasets. In *NIPS*, volume 20, 2008.

⁴Kégl, B. and Busa-Fekete, R. Boosting products of base classifiers. In *ICML*, volume 26, pp. 497–504, 2009.

⁵Viola, P. and Jones, M. Robust real-time face detection. *Int. J. of Comp. Vis.*, 57:137–154, 2004.

- Other features

- Bandit boosting⁶: [adaptive feature selection](#) to accelerate training
- Efficient [multi-platform C++](#) implementation
- Training and test data are input in [ARFF](#) format
- Support for [sparse data](#) and/or [label](#) matrix
- The [classifiers](#) are saved in [XML](#) format
- [Test/training error](#) and other iteration-wise statistics are saved in a [tab separated](#) data file

⁶Busa-Fekete, R. and Kégl, B. Boosting products of base classifiers. In *ICML*, volume 27, 2010.

Bechmark and challenge results

- **MIREX'05**

`www.music-ir.org/evaluation/mirex-results`

winner in genre classification track and runner up in the artist Identification track⁷

- **MNIST**

`yann.lecun.com/exdb/mnist`

Boosting decision product of stumps⁸ is the best reported no-domain-knowledge algorithm on MNIST after Hinton and Salakhutdinov's deep belief nets

- **Yahoo - Learning to Rank Challenge (ICML'10 workshop)**

`learningtorankchallenge.yahoo.com`

6th place in track 1 and 11th place in track 2. The difference between our calibrated AdaBoost.MH approach and the winners was less than 0.003 in both tracks.

⁷J. Bergstra et al. , "Aggregate features and AdaBoost for music classification," *Machine Learning Journal*, vol. 65, no. 2/3, pp. 473–484, 2006.

⁸Kégl, B. and Busa-Fekete, R. Boosting products of base classifiers. In *ICML*, volume 26, pp. 497–504, 2009.