

Sisteme de recunoaștere a formelor – Lab 10

Metoda AdaBoost

1. Obiective

În acest laborator vom studia un clasificator compus folosind metoda AdaBoost (Adaptive Boosting). Vom aplica acest clasificator pe o problemă de clasificare binară a unei mulțimi de puncte 2D.

2. Fundamente teoretice

AdaBoost, scurt pentru Adaptive Boosting, este o meta-metodă de învățare formulată de Yoav Freund și Robert Schapire în [1], care au câștigat premiul Gödel din 2003 pentru munca lor [2]. În acest laborator vom separa puncte 2D în două clase diferite, clasa punctelor fiind dată de culoarea lor.

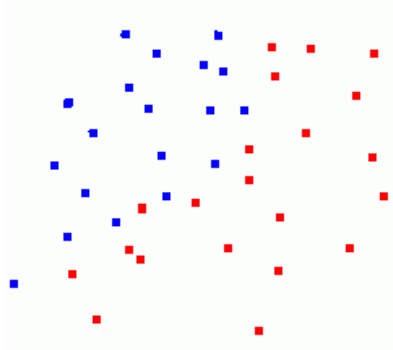


Figure 1 Exemplu de puncte din două clase

Ideea generală a algoritmului AdaBoost este construirea unui clasificator puternic $H_T(\mathbf{x})$ care este o combinație liniară a T clasificatori slabi denumiți h_t :

$$H_T(\mathbf{x}) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \right)$$

Fiecare clasificator slab returnează +1 sau -1 și este ponderat cu α_t . Rezultatul final este dat de semnul sumei. Vom utiliza un arbore de decizie cu un singur nivel (decision stump) ca clasificator slab. Acesta clasifică o instanță pe baza unei singure trăsături. Dacă valoarea trăsăturii este sub un prag se returnează -1 iar dacă este peste, se returnează +1.

Se dă setul de antrenare sub forma unei matrici de trăsături \mathbf{X} de dimensiune $n \times m$, care conține n exemple de antrenare, fiecare rând reprezentând un exemplu individual de dimensiune m . În cazul nostru, m este 2 și vom folosi ca trăsături coordonatele punctelor. Vectorul de etichete \mathbf{y} conține clasele punctelor. Vom considera punctele roșii ca având clasa +1 și punctele albastre clasa -1.

Metoda asociază o pondere fiecărui exemplu de antrenare. Vom stoca aceste ponderi în vectorul \mathbf{w} de dimensiune n . Inițial toate instanțele au aceeași pondere de $1/n$. Descriem în continuare algoritmul AdaBoost pentru a obține clasificatorul puternic H_T .

Algorithm AdaBoost

```
init  $w_i=1/n$ 
for  $t=1:T$ 
  //also returns the weighted training error  $\epsilon_t$ :
   $[h_t, \epsilon_t] = \text{findWeakLearner}(X, y, w)$ 
   $\alpha_t = 0.5 \ln\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$ 
   $s = 0$ 
  for  $i=1:n$ 
    //wrongly classified examples will have  $y_i h_t(X_i) < 0$ 
    //their weights will be higher in the next step
     $w_i \leftarrow w_i \cdot \exp(-\alpha_t y_i h_t(X_i))$ 
     $s += w_i$ 
  endfor
  //normalize the weights
  for  $i=1:n$ 
     $w_i \leftarrow w_i/s$ 
  endfor
endfor
//returns all the alpha values and the weak learners
return  $[\alpha, h]$ 
```

```
 $[h_t, \epsilon_t] = \text{findWeakLearner}(X, y, w)$ 
best_h = {}
best_err =  $\infty$ 
for  $j=1:X.\text{cols}$ 
  for threshold=0:img.size //cols or rows
    for class_label={-1,1}
       $e=0$ 
      for  $i=1:X.\text{rows}$ 
        if  $X(i, j) < \text{threshold}$ 
           $z_i = \text{class\_label}$ 
        else
           $z_i = -\text{class\_label}$ 
        endif
        if  $z_i y_i < 0$ 
           $e += w_i$ 
        endif
      endfor
      if  $e < \text{best\_err}$ 
         $\text{best\_err} = e$ 
         $\text{best\_h} = \{j, \text{threshold}, \text{class\_label}, e\}$ 
      endif
    endfor
  endfor
endfor
return best_h
```

Idee care stă la baza algoritmului este de a găsi cel mai bun clasificator slab pe baza erorii ponderate (cea mai mică valoare) și apoi să modificăm importanța exemplilor. Exemplele clasificate greșit vor avea o pondere mai mare, deci vor conta mai mult pentru selecția următorilor clasificatori slabi, iar cele corecte o pondere mai mică. Un exemplu de antrenare e clasificat greșit dacă expresia $y_i h_t(X_i)$ este negativă (eticheta de clasă și predicția noastră au semne diferite).

La pasul următor vom găsi un alt clasificator slab fiindcă acesta trebuie să clasifice corect instanțele cu pondere mai mare deoarece acestea au o contribuție mai mare în eroarea ponderată. Fiecare clasificator slab influențează scorul final. Această influență este ponderată în funcție de performanța lui pe setul de antrenare.

3. Detalii de implementare

Structură sugerată pentru un clasificator slab:

```
struct weaklearner{
    int feature_i;
    int threshold;
    int class_label;
    float error;
    int classify(Mat X){
        if (X.at<float>(feature_i)<threshold)
            return class_label;
        else
            return -class_label;
    }
};
```

Antetul funcției care returnează cel mai bun clasificator slab:

```
weaklearner findWeakLearner(Mat X, Mat y, Mat w)
```

Structură sugerată pentru clasificatorul puternic (MAXT este o constantă):

```
struct classifier{
    int T;
    float alphas[MAXT];
    weaklearner hs[MAXT];
    int classify(Mat X){
        return 0;
    }
};
```

Antetul funcției care colorează fundalul (să nu modificați imaginea originală):

```
void drawBoundary(Mat img, classifier clf)
```

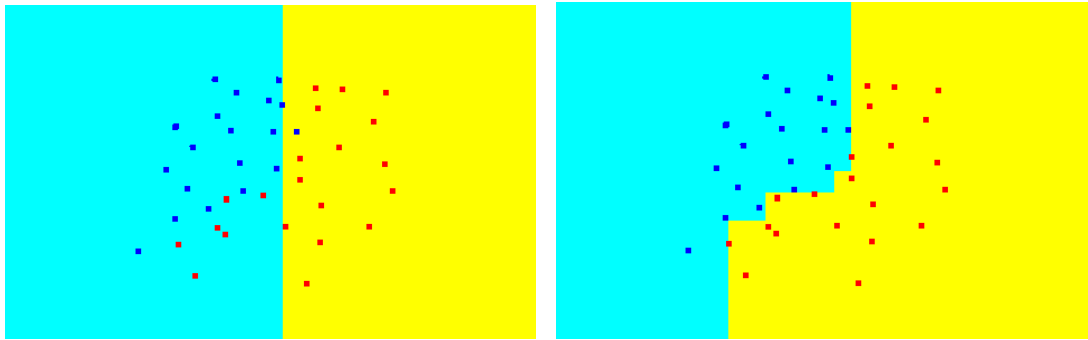


Figure 2. Rezultate obținute pe `points1` cu $T=1$ (stânga) și $T=13$ (dreapta)

4. Activitate practică

1. Citiți datele de antrenare dintr-o singură imagine (`points*.bmp`). Fiecare rând din matricea X trebuie să conțină rândul și coloana unui punct colorat din imagine. Clasa punctului se memorează în y considerând +1 pentru roșu și -1 pentru albastru.
2. Implementați clasificatorul slab.
3. Implementați funcția `findWeakLearner`.
4. Implementați funcția `drawBoundary` care colorează fundalul în funcție de apartenență. Fiecare pixel alb se colorează cu galben dacă face parte din clasa +1 și cu turcoaz pentru -1. Testați funcția cu un clasificator puternic format dintr-un singur clasificator slab.
5. Implementați algoritmul AdaBoost pentru a găsi un clasificator puternic format din T clasificatori slabi. Vizualizați bordura de decizie. Găsiți parametrul T pentru care obțineți eroare 0. Care sunt limitările metodei?

5. Referințe

- [1] Robert E. Schapire, The Boosting Approach to Machine Learning, An Overview, 2001
- [2] AdaBoost - <https://en.wikipedia.org/wiki/AdaBoost>