

# AdaBoost Notes

Yang Gu@ Scopely

April 16, 2012

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>AdaBoost</b>	<b>1</b>
2.1	Scouting . . . . .	1
2.2	Drafting . . . . .	1
2.3	Weighting . . . . .	3
2.4	Pseudocode . . . . .	4
<b>3</b>	<b>Example</b>	<b>4</b>
3.1	Data Points . . . . .	5
3.2	Weak Classifiers . . . . .	5
3.3	Scout . . . . .	6
3.4	Initialization . . . . .	7
3.5	Training . . . . .	7
3.6	Results . . . . .	8

# 1 Introduction

AdaBoost or Adaptive Boosting, is an algorithm used for generating strong classifiers out of weak classifiers. For a given input  $x$  each expert classifier  $k_j$  can emit an opinion  $k_j(x) \in \{-1, 1\}$ , and the final decision of the committee  $K$  of experts is  $\text{sign}(C(x))$ , the sign the weighted sum of expert opinions, where

$$C(x) = \sum_j^L \alpha_j k_j(x),$$

where constants  $\alpha_j$  is the weight we assign to the opinion of the  $j$ th expert in the committee.

The process of AdaBoost consists three steps: 1) scouting prospective experts, 2) drafting them, and 3) assigning a weight to their contribution to the team [1].

## 2 AdaBoost

### 2.1 Scouting

Scouting is done by testing the classifiers in the pool using a training set  $T$  of  $N$  multidimensional data points  $x_i$ . For each point  $x_i$  we have a row  $i$  built in matrix with misses (with a 1) and hits (with a 0) of each classifier. Column  $j$  is reserved for the  $j$ -th classifier in the pool:

Table 1: Scout table of classifiers

	1	2	...	L
$x_1$	0	1	...	1
$x_2$	0	0	...	1
$\vdots$	$\vdots$	$\vdots$		$\vdots$
$x_N$	0	0	...	0

The main idea of AdaBoost is to proceed systematically by extracting one classifier from the pool in each of  $M$  iterations. At beginning, all elements ( $\alpha_j$ ) are assigned the same weight (e.g., 1). As the drafting progresses, the more difficult examples, that is, those where the committee still perform badly, are assigned larger and larger weights. The best "team players" are those which can provide new insights to the committee. Classifiers being drafted should complement each other in an optimal way.

### 2.2 Drafting

In each iteration, we need to rank all classifiers, so that we can select the current best out of the pool. At the  $m$ -th iteration we have already included  $m - 1$

classifiers in the committee and we want to draft the next one. The current linear combination of classifiers is

$$C_{m-1}(x) = \sum_j^{m-1} \alpha_j k_j(x)$$

and we want to extend it to

$$C_m(x) = C_{m-1}(x) + \alpha_m k_m(x).$$

At the first iteration ( $m = 1$ ,  $C_0$  is the zero function. We define the total cost, or total error, of the extended classifier as the exponential loss

$$E = \sum_{i=1}^N e^{-y_i C_m(x_i)} = \sum_{i=1}^N e^{-y_i (C_{m-1}(x_i) + \alpha_m k_m(x_i))}$$

where  $y_i \in \{-1, 1\}$ .

If the committee's opinion (the sign) is different (a miss) with the label, e.g.  $C_m(x_i) = \beta, y_i = -1$  where  $\beta > 0$  or vice versa,

$$-y_i C_m = -(-1 \cdot \beta) = \beta,$$

the expert pool is charged by a cost  $e^\beta$ .

If the committee's opinion (the sign) is the same (a hit) as the label, e.g.  $C_m(x_i) = -\beta, y_i = -1$  or vice versa,

$$-y_i C_m = -(-1 \cdot -\beta) = -\beta,$$

the expert pool is charged by a cost  $e^{-\beta}$ . Note that  $e^{-\beta} < e^\beta$ , so that misses are more heavily penalized than hits.

Since our intention is to draft  $k_m$  we rewrite the above expression as

$$E = \sum_{i=1}^N w_i^{(m)} e^{-y_i \alpha_m k_m(x_i)}, \quad (1)$$

where

$$w_i^{(m)} = e^{-y_i C_{m-1}(x_i)}, \quad (2)$$

In the first iteration  $w_i^{(1)} = 1$  for  $i = 1, 2, \dots, N$ . During later iterations, the vector  $w^{(m)}$  represents the weight assigned to each data point in the training set at iteration  $m$ . We can split the sum in Eq. 1 into two sums

$$E = \sum_{y_i = k_m(x_i)} w_i^{(m)} e^{-\alpha_m} + \sum_{y_i \neq k_m(x_i)} w_i^{(m)} e^{\alpha_m} \quad (3)$$

$$= W_c e^{-\alpha_m} + W_e e^{\alpha_m} \quad (4)$$

which means that the total cost is the weighted cost of all hits plus the weighted cost of all misses.

Minimizing  $E$  is equivalent to minimizing  $e^{\alpha_m} E$  for a fixed  $\alpha_m$ ,

$$e^{\alpha_m} E = W_c + W_e e^{2\alpha_m} \quad (5)$$

$$= (W_c + W_e) + W_e(e^{2\alpha_m} - 1) \quad (6)$$

$$= W + W_e(e^{2\alpha_m} - 1) \quad (7)$$

$W = W_c + W_e$  is the total sum of the weights of all data points, that is, a constant in the current iteration. The right hand side of the equation is minimized when at the  $m$ -th iteration we pick the classifier with the lowest total cost  $W_e$  (we know  $e^{2\alpha_m} > 1$ ). The next Drafter,  $k_m$ , should be the one with the lowest penalty given the current set of weights.

### 2.3 Weighting

Having picked the  $m$ -th member of the committee we need to determine its weight  $\alpha_m$ . From Eq 4 we see that (note  $de^x/dx = e^x$ )

$$\frac{dE}{d\alpha_m} = -W_c e^{-\alpha_m} + W_e e^{\alpha_m}$$

Setting the above equation to 0 and multiplying by  $e^{\alpha_m}$  we obtain

$$-W_c + W_e e^{2\alpha_m} = 0$$

The optimal  $\alpha_m$  is thus:

$$\alpha_m = \frac{1}{2} \ln \left( \frac{W_c}{W_e} \right) = \frac{1}{2} \ln \left( \frac{1 - e_m}{e_m} \right),$$

where  $e_m = W_e/W_m$ , is the percentage rate of error given the weights of the data points. Since the percentage is invariant if we normalize the weight, normalizing the weights of data at each iteration is irrelevant.

We proceed to get the recursive update equations:

$$\begin{aligned} e^{-\alpha_m} &= e^{-\frac{1}{2} \ln \left( \frac{1 - e_m}{e_m} \right)} \\ &= \left( e^{\ln \left( \frac{1 - e_m}{e_m} \right)} \right)^{-\frac{1}{2}} \\ &= \sqrt{\frac{e_m}{1 - e_m}} \end{aligned}$$

$$\begin{aligned} e^{\alpha_m} &= e^{\frac{1}{2} \ln \left( \frac{1 - e_m}{e_m} \right)} \\ &= \left( e^{\ln \left( \frac{1 - e_m}{e_m} \right)} \right)^{\frac{1}{2}} \\ &= \sqrt{\frac{1 - e_m}{e_m}} \end{aligned}$$

$$\begin{aligned}
w_i^{(m+1)} &= e^{-y_i C_m(x_i)} \\
&= e^{-y_i (C_{m-1}(x_i) + \alpha_m k_m(x_i))} \\
&= e^{-y_i C_{m-1}(x_i)} \cdot e^{-y_i \alpha_m k_m(x_i)} \\
&= w_i^{(m)} e^{-y_i \alpha_m k_m(x_i)}
\end{aligned}$$

$$w_i^{(m+1)} = \begin{cases} w_i^{(m)} e^{-\alpha_m} = w_i^{(m)} \sqrt{\frac{e_m}{1-e_m}}, & y_i = k_m(x_i) \quad \text{hit} \\ w_i^{(m)} e^{\alpha_m} = w_i^{(m)} \sqrt{\frac{1-e_m}{e_m}}, & y_i \neq k_m(x_i) \quad \text{miss} \end{cases} \quad (8)$$

## 2.4 Pseudocode

Given a training set  $T$  of data points  $x_i$  and their labels  $y_i$  in a two class problem, we assign initial weights  $w_j^{(1)} = 1$  to all data points  $x_i$ . We perform  $M$  iterations.

For  $m = 1$  to  $M$

1. Select and extract from the pool of classifiers the classifier  $k_m$  which minimizes

$$W_e = \sum_{y_i \neq k_m(x_i)} w_i^{(m)}$$

2. Set the weight  $\alpha_m$  of the classifier to

$$\alpha_m = \frac{1}{2} \ln \left( \frac{W_c}{W_e} \right)$$

3. Update the weights of the data points for the next iteration. If  $k_m(x_i)$  is a miss, set

$$w_i^{(m+1)} = w_i^{(m)} e^{\alpha_m} = w_i^{(m)} \sqrt{\frac{1-e_m}{e_m}}$$

otherwise

$$w_i^{(m+1)} = w_i^{(m)} e^{-\alpha_m} = w_i^{(m)} \sqrt{\frac{e_m}{1-e_m}}$$

## 3 Example

Define a two-class problem in the plane:

- class 1: points in circle  $x^2 + y^2 < 1$
- class 2: points not in circle and in box
  - $-2 \leq x \leq 2$  and
  - $-2 \leq y \leq 2$  and

$$-x^2 + y^2 > 1$$

which includes a circle of points inside a square. We will build a strong classifier out of a pool of  $L$  (10 or 100) randomly generated linear discriminates of the type  $\text{sign}(ax_1 + bx_2 + c)$ .

### 3.1 Data Points

Below python code generates 1000 data points for each class.

```
N = 1000
def unif(n):
    vec = []
    while len(vec) < n:
        vec.append(uniform(-1,1))
    return vec

def generate_class(label=1):
    X = []
    while len(X) < N:
        x,y = unif(2)
        if label ==1 and x**2 + y**2 >= 1: continue
        if label ==-1:
            x *= 2; y*= 2
            if x**2 + y**2 < 1: continue
        X.append([x,y])
    return X
```

### 3.2 Weak Classifiers

We randomly generate  $L$  linear discriminants using below code:

```
L = 10
def generate_classifiers():
    C = []
    while len(C) < L:
        a,b,c = unif(3)
        C.append([a,b,c])
    return C
```

To get predictions from the linear classifiers:

```
def is_prediction_hit(classifier, data, label):
    a,b,c = classifier
    x, y = data
    if linear(a,b,c,x,y) > 0:
        pred = 1
```

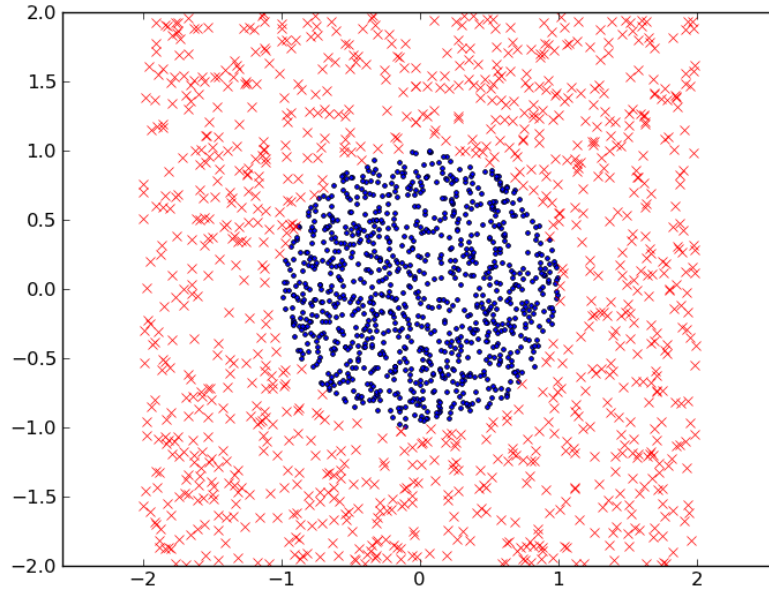


Figure 1: The data plot of the two-class problem. Blue dots are in class 1 and red crosses are in class 2.

```

else:
    pred = -1
    return pred == label

```

### 3.3 Scout

Assuming class 1 are labeled with 'True' and class 2 are labeled with 'False' in our data. We generate our scout table with below codes:

```

X1 = generate_class(1)
X2 = generate_class(-1)
C = generate_classifiers()
pred = fill_scout()

def fill_scout():
    pred = {}
    for j, classifier in enumerate(C):
        pred[j] = []
        for i, data in enumerate(X1):

```

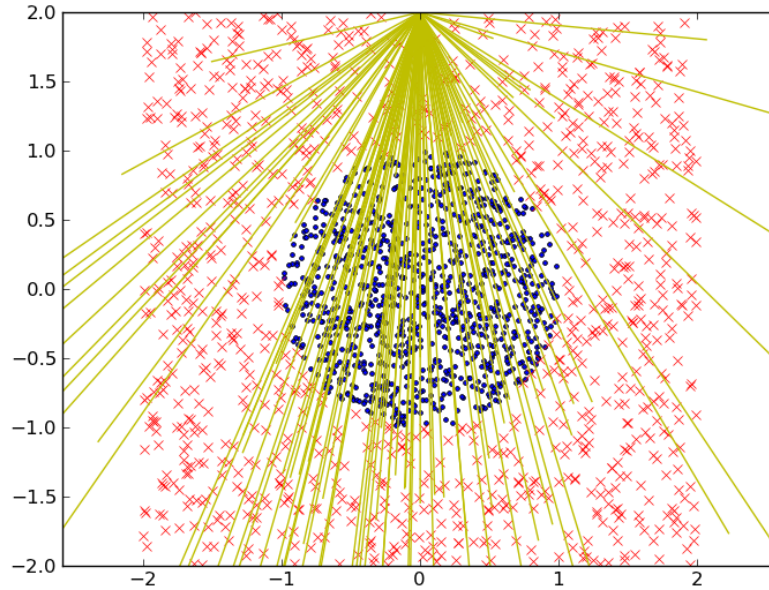


Figure 2: The data plot of the two-class problem. Blue dots are in class 1 and red crosses are in class 2. The yellow lines are linear classifiers.

```

        pred[j].append(is_prediction_hit(classifier, data, 1))
    for i,data in enumerate(X2):
        pred[j].append(is_prediction_hit(classifier, data, -1))
    return pred

```

### 3.4 Initialization

```

w = [1.0]*(2*N)
alpha = [0]*L
in_committee = set()

```

### 3.5 Training

```

def get_draftee(pred):
    min_We = 1e100
    min_W = 1e100
    min_classifier = -1
    for j in range(L):
        if j in in_committee: continue

```



```

We = 0
W = 0
for i, prediction in enumerate(pred[j]):
    W += w[i]
    if not prediction:
        We += w[i]
if We < min_We:
    min_classifier = j
    min_We = We
    min_W = W

em = float(min_We)/min_W
print min_We, min_W
alpha[min_classifier] = .5*log((1-em)/em)
in_committee.add(min_classifier)

for i in range(2*N):
    if not pred[min_classifier]:
        w[i] *= ((1-em)/em)**.5
    else:
        w[i] *= (em/(1-em))**.5

```

### 3.6 Results

Test code:

```

def test(X, y):
    cnt = 0
    for i, data in enumerate(X):
        pred = 0
        for j, classifier in enumerate(C):
            if is_prediction_hit(classifier, data, y):
                pred += alpha[j]
            else:
                pred += -alpha[j]
        if (pred > 0 and y > 0) or (pred < 0 and y < 0): cnt += 1
    return cnt

def testAll(X1, X2):
    cnt = 0
    cnt += test(X1, 1)
    cnt += test(X2, -1)
    print float(cnt)/2/N

T1 = generate_class(1)
T2 = generate_class(-1)

```

```
testAll(X1,X2)
testAll(T1,T2)
```

Training/Testing error : percentage of wrongly classified points

Table 2: Training and Testing error

number of classifiers	10	20	50	100
Train Error %	13.7	4.6	1.45	0
Test Error%	13.85	5.0	2.3	0

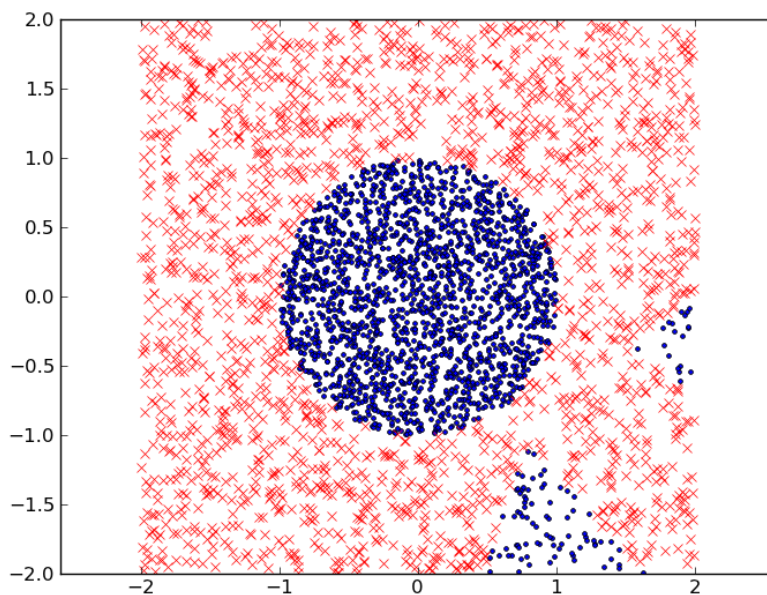


Figure 3: Predicted class labels in testing with 10 classifiers with some wrong labels. It is very clear that those error appears where we lack training examples.

## References

- [1] Baul Rojas, *AdaBoost and the Super Bowl of Classifiers: A tutorial introduction to Adaptive Boosting*, <http://www.inf.fu-berlin.de/inst/ag-ki/>

adaboost4.pdf