# AdaBoost and K-Fold Cross-Validation on Hand-Written Digits

## KK Feng

# 1 Problem

## 1.1 Description

Classify grayscale images for hand-written digits.

## 1.2 Data

- **uspsdata.txt**: contains a matrix with one data point (= vector of length 256) per row. The 256-vector in each row represents a 16 by 16 image of a handwritten number.

- **uspscl.txt**: contains the corresponding class labels. The data contains two classes - the digits 5 and 6 - so the class labels are stored as -1 and +1, respectively.

## 1.3 Idea

- Adaptive Boosting algorithm with decision stumps as weak learners.

- K-Fold Cross-Validation to tune the number of weak learners.

# 2 Solution

## 2.1 Implementation

To train decision stumps, we implement following algorithm

---
**Algorithm 1** A simple training algorithm for decision stumps
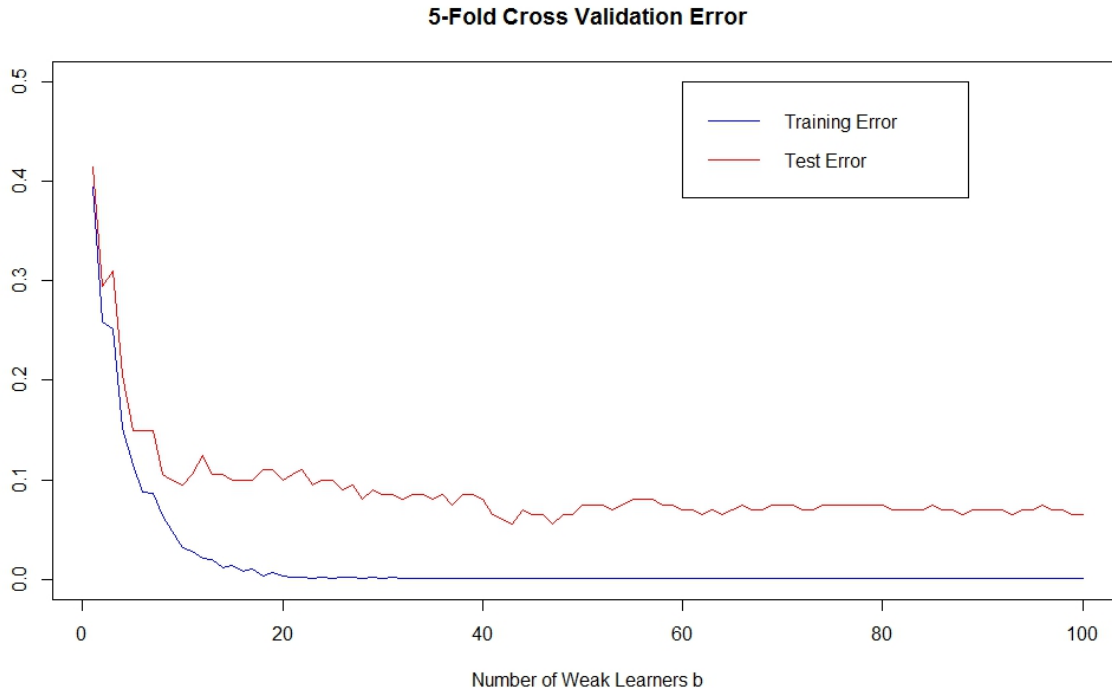
---
**Require:** Data $X = (x_1, \cdots, x_n)$ where $x_i \in \mathbb{R}^d$, weight $w$, label $y$
 1: **for** $j = 1 : d$ **do**
 2:    Sort samples $x_i$ in ascending order along dimension $j$
 3:    **for** $i = 1 : n$ **do**
 4:       Compute cumulative sums $cum_i^j = \sum_{k=1}^{i} w_k y_k$
 5:    **end for**
 6:    Threshold $\theta_j$ is obtained at the extrema of $cum_i^j$
 7:    Label $m_j$ is obtained from the sign of cumulative sum at extrema
 8:    Compute the error rate of classifier $(\theta_j, m_j)$ along dimension j
 9: **end for**
10: Find optimal $j^*, \theta^*$ in which the classifier $(\theta_j, m_j)$ gives the minimum error rate

---

(Reference: http://ais.informatik.uni-freiburg.de/teaching/ws09/robotics2/pdfs/rob2-10-adaboost.pdf )

## 2.2 Plot

Please find the plots for training error and test error as a function of b (number of weak learners) in the following. Note that the cross validation error is the average of errors of 5 folds.

**5-Fold Cross Validation Error**



From the plot, we can see that for the USPS data and using 5-fold cross validation, the training error reaches bottom of the curve when we use approximately 20 weak learners, and the training error curve become flat when number of weak learners go larger than 20. On the other hand, we need around 40 weak learners to ensure that we have the optimal test error. If number of weak learners go larger than 40, the test error will just have small oscillations around the optimal test error we get at 40 weak learners.

## 2.3 Code

```
1   ###### AdaBoost and K-Fold Cross-Validation on Hand-Written Digits ######
2
3   # Train: weak learner training routine
4   train <- function(X, w, y){
5       n <- dim(X)[1]
6       d <- dim(X)[2]
7       theta <- rep(0,d)
8       m <- rep(0,d)
9       error <- rep(0,d)
10      # find best stump classifier (theta_j, m_j) for each dimension j
11      for (j in 1:d){
12          x_order <- order(X[,j]) # get order of data along dimension j
13          x <- X[x_order,j]
14          cum <- rep(0,n) # compute cumulative sums cum_i^j = \sum_{k=1}^i w_ky_k
15          weighted_label <- w[x_order] * y[x_order]
16          cum[1] <- weighted_label[1]
17          for (i in 2:n){
18              cum[i] <- weighted_label[i] + cum[i-1]
19          }
20          index <- which.max(abs(cum))  # find theta_j, m_j
21          theta[j] <- x[index]
22          m[j] <- sign(cum[index])
```

2

```r
23      yy <- rep(-m[j], n)
24      yy[x > theta[j]] <- m[j]
25      error[j] <- (yy != y) %*% w    # compute error rate of classifier (theta_j, m_j)
26    }
27    j_star <- which.min(error)   # find optimal dimension j
28    pars <- list(j = j_star, theta = theta[j_star], m = m[j_star])
29    return(pars)
30  }
31
32  # Classify: evaluates the weak learner on X using the parametrization pars
33  classify <- function(X, pars){
34    label <- rep(-pars$m, dim(X)[1])
35    label[X[,pars$j] > pars$theta] <- pars$m
36    return(label)
37  }
38
39  # Agg_class: evaluates the boosting classifier ("aggregated classifier") on X.
40  agg_class <- function(X, alpha, allPars){
41    n <- dim(X)[1]
42    B <- length(alpha)
43    label_sum <- rep(0,n)
44    for (b in 1:B){
45      label_sum <- label_sum + alpha[b] * classify(X, allPars[[b]]) # sum up weighted labels
46    }
47    c_hat <- sign(label_sum)
48    return(c_hat)
49  }
50
51  # AdaBoost: implement the AdaBoost algorithm
52  AdaBoost <- function(X, y, B){
53    n <- length(y)
54    w <- rep(1/n, n) # Initialize weights
55    alpha <- rep(0,B) # Initialize alphas
56    allPars <- rep(list(list(0)), B)
57    for (b in 1:B) {
58      allPars[[b]] <- train(X, w, y)  # Train a weak learner c_b
59      index <- y != classify(X, allPars[[b]]) # misclassification index
60      error <- sum(w[index]) / sum(w)  # Compute error
61      alpha[b] <- log((1-error)/error) # Compute voting weights
62      w[index] <- w[index] * exp(alpha[b]) # Recompute weights
63    }
64    return(list(alpha = alpha, allPars = allPars)) # Return classifier
65  }
66
67  # Problem 1.3 Run algorithm on the USPS data and evaluate results using cross validation.
68  X <- read.table("uspsdata.txt")
69  y <- read.table("uspscl.txt")[,1]
70  n <- length(y)
71  B <- 100   # maximum number of weak learners
72  m <- 5   # 5-fold cross validation
73  train_error <- matrix(0, nrow = B, ncol = m)
74  test_error <- matrix(0, nrow = B, ncol = m)
75  for (i in 1:m){
76    # generate train data and test data for fold i
77    index <- round(n/m*(i-1)+1) : trunc(n/m*i)
78    data_train <- X[-index,]
79    y_train <- y[-index]
80    data_test <- X[index,]
81    y_test <- y[index]
82    # get AdaBoost classifer
83    AB <- AdaBoost(data_train, y_train, B)
84    alpha <- AB$alpha
85    allPars <- AB$allPars
86    for (b in 1:B){
87    # compute train and test error for fold i by AdaBoost with b weak learners
88      train_error[b,i] <- sum(y_train != agg_class(data_train, alpha[1:b], allPars[1:b]))/
              length(y_train)
```

3

```
89        test_error[b,i] <- sum(y_test != agg_class(data_test, alpha[1:b], allPars[1:b]))/length(
             y_test)
90      }
91   }
92   # compute cross validation error
93   cross_train_error <- rep(0, B)
94   cross_test_error <- rep(0, B)
95   for (b in 1:B)
96   {
97      cross_train_error[b] <- mean(train_error[b,])
98      cross_test_error[b] <- mean(test_error[b,])
99   }
100
101  # Plot the training error and the test error as a function of b.
102  plot(cross_train_error,type='l',ylim=c(0,0.5),col='blue',xlab='Number of Weak Learners b',
          ylab='', main='5-Fold Cross Validation Error')
103  lines(cross_test_error,col='red')
104  legend(60,0.5,c('Training Error','Test Error'),col=c('blue','red'),lty=1)
```