

# Homework 3

36-708, Spring 2021

**Due April 16 at 5PM EST**

Please attach all code to your homework. In an RMarkdown document for example, this can be done in one line (see [Yihui Xie's website](#) for how to do this).

## 1 Deriving boosting-like variants for square and logistic losses

- Data: Assume  $n$  labeled data points  $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R} \times \{-1, 1\}$ .
- Classifier: Denote by  $f = \sum_{t=1}^T \alpha_t h_t$  the linear combination of base classifiers  $h_t$  with weights  $\alpha_t$ .
- Loss function: Let the loss function to be used in the boosting derivation be  $L = \frac{1}{n} \sum_{i=1}^n \phi(-y_i f(x_i))$  for some function  $\phi$  (to be specified). The 0-1 loss is  $\phi(-u) = \mathbb{1}(u \leq 0)$  but we will use surrogates of this.
- Algorithm: *Boosting*, which we will derive from the point of view of coordinate descent on the loss function  $L$ . In particular at each step  $t$  of boosting, we find the best base classifier  $h_t$  and step size  $\alpha_t$  and update the function  $f$  with the new component  $\alpha_t h_t$ .

### 1.1 AdaBoost as coordinate descent

In this problem, we will use  $\phi_a(-u) := \exp(-u)$ , the exponential loss. We will rederive boosting as coordinate descent to recover AdaBoost.

(a) Write the objective function we wish to minimize:  $F(\alpha) := \frac{1}{n} \sum_{i=1}^n \exp\left(-y_i \sum_{j=1}^m h_j(x_i) \alpha_j\right)$ . Consider the “coordinate” of steepest descent:

$$\operatorname{argmin}_k \left. \frac{\partial F(\alpha + \eta e_k)}{\partial \eta} \right|_{\eta=0}.$$

Show that

$$\begin{aligned} \frac{\partial F(\alpha + \eta e_k)}{\partial \eta} &= \left( \prod_{j=1}^m Z_j \right) \sum_{i=1}^n -y_i h_k(x_i) w_i \\ &= 2\varepsilon_{t,k} - 1 \end{aligned}$$

where  $w_i$  and  $\prod_j Z_j$  are the weights and normalizing factor described in the lecture. Finally, conclude that the  $\operatorname{argmin}_k$  of the above is the weak learner with smallest weighted 0-1 error.

(b) Show that the step-size  $\eta$  is given by  $\frac{1}{2} \log\left(\frac{1-\varepsilon_{t,k}}{\varepsilon_{t,k}}\right)$ . That is, show that  $\frac{1}{2} \log\left(\frac{1-\varepsilon_{t,k}}{\varepsilon_{t,k}}\right)$  is the solution to

$$\frac{\partial F(\alpha + \eta e_k)}{\partial \eta} = 0.$$

(c) Write the pseudo-code for AdaBoost.

## 1.2 SquareBoost using the squared loss

We will repeat the previous problem but for the squared loss,  $\phi_s(-u) = (1 - u)^2 \mathbb{1}(u \leq 1)$ .

(a) Consider the objective function:

$$F(\alpha) = \frac{1}{n} \sum_{i=1}^n \left( 1 - y_i \sum_{j=1}^m h_j(x_i) \alpha_j \right)^2 \mathbb{1} \left( y_i \sum_{j=1}^m h_j(x_i) \alpha_j \leq 1 \right)$$

Using similar steps to before, show that the minimizer of the above loss is the weak learner with best 0-1 loss on the weighted data with weights given by

$$w_i := \frac{(1 - y_i f(x_i)) \mathbb{1} \left( y_i \sum_{j=1}^m h_j(x_i) \alpha_j \leq 1 \right)}{\sum_{i=1}^n (1 - y_i f(x_i)) \mathbb{1} \left( y_i \sum_{j=1}^m h_j(x_i) \alpha_j \leq 1 \right)}$$

(b) Unlike AdaBoost, the step size of SquareBoost cannot be computed in closed-form. How would you go about computing it?

(c) Write the pseudo-code for SquareBoost.

## 1.3 LogisticBoost using the logistic loss

We will repeat the previous problems but for the logistic loss,  $\phi_l(-u) = \log(1 + e^{-u})$ .

(a) Show that the minimizer of the logistic loss is given by the weak learner with smallest 0-1 loss on the weighted data where weights are given by

$$w_i = \frac{\text{logit}^{-1}(-y_i f(x_i))}{\sum_{i=1}^n \text{logit}^{-1}(-y_i f(x_i))}$$

(b) Show how finding the step-size is equivalent to training a logistic regression model with and offset and no intercept. Explain how you could compute this step size.

(c) Write the pseudo-code for LogisticBoost.

## 1.4 Comparison of AdaBoost, SquareBoost, and LogisticBoost

(a) How do the three previous boosting algorithms differ?

# 2 Implementing and evaluating boosting

In this problem, we will implement the previous three boosting algorithms and apply them to the spam dataset.

- Dataset: We will use the spam dataset available [here](#).
- Train and test split: Use a random 3:1 split.

(a) Implement AdaBoost, SquareBoost, and LogisticBoost with decision stumps as the base classifiers.

(b) Run the boosting algorithms for various values of the number of boosting rounds  $t$  (for some reasonable upper limit  $T$ ).

(c) Plot the train and test errors of all three algorithms as a function of  $t$ .

(d) Summarize your findings. In particular, comment on the following.

- How fast does the training error decrease?
- Does the training error reach zero? What happens to the testing error if you train longer?

**3 TBD**

**4 TBD**