

## CS480/680: Introduction to Machine Learning

## Homework 3

Due: 11:59 pm, March 18, 2025, submit on LEARN and Crowdmark.

NAME  
student number

Submit your writeup in pdf and all source code in a zip file (with proper documentation). Write a script for each programming exercise so that the TA can easily run and verify your results. Make sure your code runs!  
[Text in square brackets are hints that can be ignored.]

## Exercise 1: Adaboost (5 pts)

In this exercise we will interpret Adaboost as minimizing the exponential loss:

$$\min_{\alpha} \frac{1}{n} \sum_{i=1}^n \exp \left[ -y_i \sum_t \alpha_t h_t(\mathbf{x}_i) \right], \quad (1)$$

where  $h_t$  are the so-called weak learners, and the aggregated classifier

$$h_{\alpha}(\mathbf{x}) := \sum_t \alpha_t h_t(\mathbf{x}). \quad (2)$$

Note that we assume  $y_i \in \{\pm 1\}$  in this exercise.

Let us introduce the uniform distribution  $\mathbf{p}_1 = \frac{1}{n} \mathbf{1}$  over our training set  $\{\mathbf{x}_i, y_i\}_{i=1}^n$ , and rewrite (1) as:

$$\min_{\alpha} \mathbb{E}_{\mathbf{p}_1} \exp \left[ -Y \sum_t \alpha_t h_t(\mathbf{X}) \right], \quad (3)$$

where  $(\mathbf{X}, Y) \sim \mathbf{p}_1$ , i.e., with probability  $p_{i1} = \frac{1}{n}$ ,  $\mathbf{X} = \mathbf{x}_i, Y = y_i$ .

- (1 pt) Let  $q > 0$  be an arbitrary function (or vector). By normalization, i.e.,  $q \leftarrow q / \int q$  or  $q \leftarrow q / q^{\top} \mathbf{1}$  we obtain a density function (or probability mass function). Find probability density (vector)  $\mathbf{p}_2$  below so that

$$\mathbb{E}_{\mathbf{p}_1} \exp \left[ -Y \sum_{t=1}^2 \alpha_t h_t(\mathbf{X}) \right] = Z_1 \cdot \mathbb{E}_{\mathbf{p}_2} \exp \left[ -Y \sum_{t=2}^2 \alpha_t h_t(\mathbf{X}) \right], \quad (4)$$

as well as the formula for  $Z_1$  (a positive constant).

Ans:

- (1 pt) Apply the previous exercise repeatedly with probability densities (vectors)  $\mathbf{p}_t$  so that

$$\mathbb{E}_{\mathbf{p}_1} \exp \left[ -Y \sum_{t=1}^T \alpha_t h_t(\mathbf{X}) \right] = \prod_{t=1}^T Z_t, \quad (5)$$

where each  $Z_t$  is a positive constant. Explain what is  $Z_t$  for each  $t$ .

Ans:

3. (1 pt) Prove the following bound on the training error:

$$\mathbb{E}_{\mathbf{p}_1} [\underbrace{Y h_{\alpha}(\mathbf{X})}_{\hat{Y}} \leq 0] = \mathbb{E}_{\mathbf{p}_1} [Y \sum_{t=1}^T \alpha_t h_t(\mathbf{X}) \leq 0] \leq \prod_{t=1}^T Z_t, \quad (6)$$

where each  $Z_t$  is given in the previous exercise. Recall that  $h_{\alpha}(\mathbf{X}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{X})$  is the aggregated classifier.

Ans:

4. (1 pt) Assuming in the  $t$ -th iteration we have found  $h_t$ . We now aim to find its coefficient  $\alpha_t$  by considering the following (convex) minimization problem:

$$\min_{\alpha_t} \mathbb{E}_{\mathbf{p}_t} \exp[-Y \alpha_t h_t(\mathbf{X})], \quad (7)$$

Suppose there is indeed a minimizer, then it must satisfy (by setting derivative to 0):

$$0 = \mathbb{E}_{\mathbf{p}_t} \{Y h_t(\mathbf{X}) \cdot \exp[-Y \alpha_t h_t(\mathbf{X})]\}. \quad (8)$$

From the above result deduce that

$$0 = \mathbb{E}_{\mathbf{p}_{t+1}} [Y h_t(\mathbf{X})], \quad (9)$$

and show that for (deterministic) weak classifiers  $h_t \in \{\pm 1\}$ :

$$\mathbb{E}_{\mathbf{p}_{t+1}} [h_t(\mathbf{X}) \neq Y] = \frac{1}{2}, \quad (10)$$

namely that in the next iteration  $t + 1$ , the previous classifier  $h_t$  has error exactly  $\frac{1}{2}$ .

Ans:

5. (1 pt) In (8) above we obtained a nonlinear equation of  $\alpha_t$ . Although it is possible to find  $\alpha_t$  through numerical root finding algorithms, we prefer to derive a closed-form solution. Assuming  $h_t \in [-1, 1]$  is given, we can apply the bound (the so-called **Jensen's inequality**)

$$\exp(-\alpha u) \leq \frac{1+u}{2} \exp(-\alpha) + \frac{1-u}{2} \exp(\alpha) \quad (11)$$

to (7) first and then derive the optimal  $\alpha_t$ .

[This is essentially the coefficient  $\log \frac{1}{\beta_t}$  that we saw in class, up to some trivial changes.]

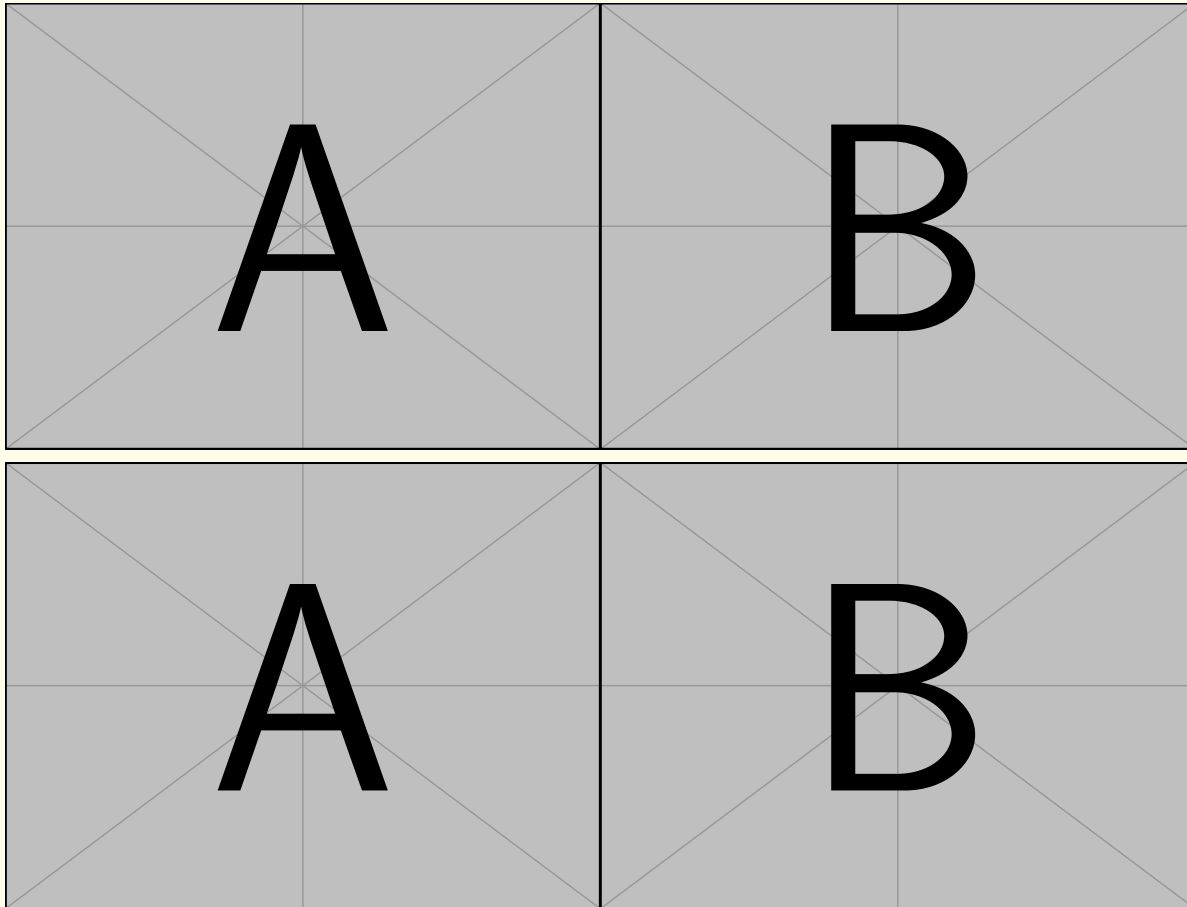
Ans:

(Remark) We have not talked much about choosing weak classifiers  $h_t$ . Here is the catch: we could simply pretend we enumerate **all** weak classifiers (infinitely many!) in our final aggregate  $h_\alpha$ . All we need to figure out is the weight  $\alpha_t$  that we assign to each weak classifier  $h_t$ , and a zero  $\alpha$  means the corresponding weak classifier is effectively discarded. The Adaboost algorithm starts with  $\alpha \equiv 0$  and only changes one  $\alpha$  into nonzero in each iteration. Thus, after  $T$  iterations, we have at most  $T$  nonzero  $\alpha$ 's. On a high level, this is very similar to kernels where the dual problem only involves  $n$  nonzero Lagrangian multiplier  $\alpha$ 's ( $n$  being the size of the training set). See, we do not need to fear infinite dimensions!

## Exercise 2: Vision Transformers (10 pts)

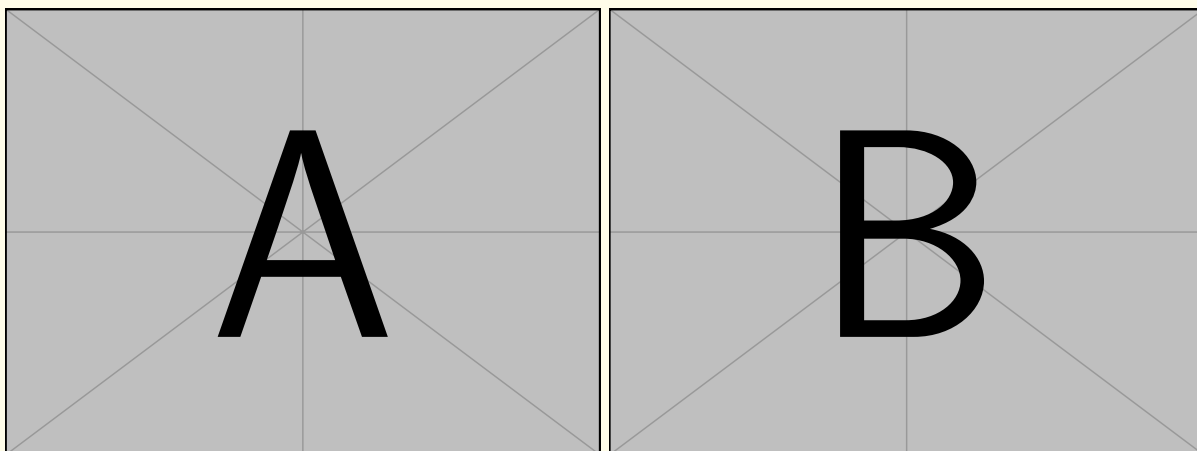
Please follow the instructions of this [ipynb file](#).

- (1+3+2 = 6 pts) Complete the missing coding parts in the provided [ipynb file](#).
- (1 pt) Visualization of patches:



3. (1 pt) The test accuracy I obtained on MNIST is: [xxx%](#)

4. (2 pts) Training / Validation accuracy vs. epoch:



### Exercise 3: Generative Adversarial Networks (5 pts)

Let us consider the game between the generator  $q(\mathbf{x})$  (the implicit density of  $T_{\theta}(Z)$ ) and the discriminator  $S(\mathbf{x})$ :

$$\inf_q \sup_S \int_{\mathbf{x}} S(\mathbf{x})p(\mathbf{x})d\mathbf{x} + \int_{\mathbf{x}} \log(1 - \exp(S(\mathbf{x})))q(\mathbf{x})d\mathbf{x} + \log 4. \quad (12)$$

We remind that  $q$  is a probability density (so is  $p$  which is given) while  $S$  is any real-valued function.

1. (1 pt) Fix an arbitrary generator  $q$  and find the resulting optimal discriminator  $S$ .

Ans:

2. (1 pt) Plug the optimal discriminator  $S$  above back to (12) and find the optimal generator  $q$ .

Ans:

Now we swap the order of the two players:

$$\sup_S \inf_q \int_{\mathbf{x}} S(\mathbf{x})p(\mathbf{x})d\mathbf{x} + \int_{\mathbf{x}} \log(1 - \exp(S(\mathbf{x})))q(\mathbf{x})d\mathbf{x} + \log 4. \quad (13)$$

3. (1 pt) Fix an arbitrary discriminator  $S$  and find an optimal generator  $q$ .

Ans:

4. (1 pt) Plug the optimal generator  $q$  above back to (13) and find the optimal discriminator  $S$ .

[Hint: average  $\leq$  max.]

Ans:

5. (1 pt) Let  $f : \mathbb{R}_+ \rightarrow \mathbb{R}$  be a convex function. We see in class that the  $f$ -divergence admits the following variational form:

$$-\mathbb{D}_f(\mathbf{p}||\mathbf{q}) = \inf_{S:\mathbb{R}^d \rightarrow [0,1]} -\mathbb{E}_{\mathbf{X} \sim \mathbf{p}}[S(\mathbf{X})] + \mathbb{E}_{\mathbf{X} \sim \mathbf{q}}[f^*(S(\mathbf{X}))], \quad (14)$$

where for simplicity we have restricted the range of  $S$  to  $[0, 1]$ . Now consider the following distribution  $(\mathbf{X}, Y) \sim \mathcal{D}$  where  $Y = \pm 1$  with equal probability while

$$[\mathbf{X} | Y = 1] \sim \mathbf{p} \text{ and } [\mathbf{X} | Y = -1] \sim \mathbf{q}. \quad (15)$$

We claim that

$$-\frac{1}{2}\mathbb{D}_f(\mathbf{p}||\mathbf{q}) = \inf_{S:\mathbb{R}^d \rightarrow [0,1]} \mathbb{E}_{(\mathbf{X}, Y) \sim \mathcal{D}}[\ell(Y, S(\mathbf{X}))]. \quad (16)$$

Express the binary loss function  $\ell$  in terms of  $f$ . Thus, given an  $f$ -divergence, we may rewrite it as a binary classification problem! Conversely, given any *proper* loss function  $\ell$ , we may reverse the argument and induce an  $f$ -divergence from the binary loss  $\ell$ .

Ans: