

CM146, Winter 2019  
Problem Set 5: Boosting, Unsupervised learning  
Due March 16, 2019, 11:59pm

## Submission instructions

- Submit your solutions electronically on the course Gradescope site as PDF files.
- If you plan to typeset your solutions, please use the LaTeX solution template. If you must submit scanned handwritten solutions, please use a black pen on blank white paper and a high-quality scanner app.

## 1 AdaBoost [5 pts]

In the lecture on ensemble methods, we said that in iteration  $t$ , AdaBoost is picking  $(h_t, \beta_t)$  that minimizes the objective:

$$\begin{aligned}(h_t^*(\mathbf{x}), \beta_t^*) &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} \sum_n w_t(n) e^{-y_n \beta_t h_t(\mathbf{x}_n)} \\ &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} (e^{\beta_t} - e^{-\beta_t}) \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)] \\ &\quad + e^{-\beta_t} \sum_n w_t(n)\end{aligned}$$

We define the weighted misclassification error at time  $t$ ,  $\epsilon_t$  to be  $\epsilon_t = \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)]$ . Also the weights are normalized so that  $\sum_n w_t(n) = 1$ .

- Take the derivative of the above objective function with respect to  $\beta_t$  and set it to zero to solve for  $\beta_t$  and obtain the update for  $\beta_t$ .
- Suppose the training set is linearly separable, and we use a hard-margin linear support vector machine (no slack) as a base classifier. In the first boosting iteration, what would the resulting  $\beta_1$  be?

## 2 K-means for single dimensional data [5 pts]

In this problem, we will work through K-means for a single dimensional data.

- Consider the case where  $K = 3$  and we have 4 data points  $x_1 = 1, x_2 = 2, x_3 = 5, x_4 = 7$ . What is the optimal clustering for this data? What is the corresponding value of the objective?

---

Parts of this assignment are adapted from course material by Jenna Wiens (UMich) and Tommi Jaakola (MIT).

- (b) One might be tempted to think that Lloyd's algorithm is guaranteed to converge to the global minimum when  $d = 1$ . Show that there exists a suboptimal cluster assignment (*i.e.*, initialization) for the data in the above part that Lloyd's algorithm will not be able to improve (to get full credit, you need to show the assignment, show why it is suboptimal *and* explain why it will not be improved).

### 3 Gaussian Mixture Models [8 pts]

We would like to cluster data  $\{x_1, \dots, x_N\}$ ,  $x_n \in \mathbb{R}^d$  using a Gaussian Mixture Model (GMM) with  $K$  mixture components. To do this, we need to estimate the parameters  $\theta$  of the GMM, *i.e.*, we need to set the values  $\theta = \{\omega_k, \mu_k, \Sigma_k\}_{k=1}^K$  where  $\omega_k$  is the mixture weight associated with mixture component  $k$ , and  $\mu_k$  and  $\Sigma_k$  denote the mean and the covariance matrix of the Gaussian distribution associated with mixture component  $k$ .

If we knew which cluster each sample  $x_n$  belongs to (we had complete data), we showed in the lecture on Clustering that the log likelihood  $l$  is what we have below and we can compute the maximum likelihood estimate (MLE) of all the parameters.

$$\begin{aligned} l(\theta) &= \sum_n \log p(\mathbf{x}_n, z_n) \\ &= \sum_k \sum_n \gamma_{nk} \log \omega_k + \sum_k \left\{ \sum_n \gamma_{nk} \log N(\mathbf{x}_n | \mu_k, \Sigma_k) \right\} \end{aligned} \quad (1)$$

Since we do not have complete data, we use the EM algorithm. The EM algorithm works by iterating between setting each  $\gamma_{nk}$  to the posterior probability  $p(z_n = k | \mathbf{x}_n)$  (step 1 on slide 26 of the lecture on Clustering) and then using  $\gamma_{nk}$  to find the value of  $\theta$  that maximizes  $l$  (step 2 on slide 26). We will now derive updates for one of the parameters, *i.e.*,  $\mu_j$  (the mean parameter associated with mixture component  $j$ ).

- (a) To maximize  $l$ , compute  $\nabla_{\mu_j} l(\theta)$ : the gradient of  $l(\theta)$  with respect to  $\mu_j$ .
- (b) Set the gradient to zero and solve for  $\mu_j$  to show that  $\mu_j = \frac{1}{\sum_n \gamma_{nj}} \sum_n \gamma_{nj} \mathbf{x}_n$ .
- (c) Suppose that we are fitting a GMM to data using  $K = 2$  components. We have  $N = 5$  samples in our training data with  $x_n, n \in \{1, \dots, N\}$  equal to:  $\{5, 15, 25, 30, 40\}$ .

We use the EM algorithm to find the maximum likelihood estimates for the model parameters, which are the mixing weights for the two components,  $\omega_1$  and  $\omega_2$ , and the means for the two components,  $\mu_1$  and  $\mu_2$ . The standard deviations for the two components are fixed at 1. Suppose that at the end of step 1 of iteration 5 in the EM algorithm, the soft assignment  $\gamma_{nk}$  for the five data items are as shown in Table 1.

$\gamma_1$	$\gamma_2$
0.2	0.8
0.2	0.8
0.8	0.2
0.9	0.1
0.9	0.1

Table 1: Entry in row  $n$  and column  $k$  of the table corresponds to  $\gamma_{nk}$

What are updated values for the parameters  $\omega_1$ ,  $\omega_2$ ,  $\mu_1$ , and  $\mu_2$  at the end of step 2 of the EM algorithm?