

CS182, Spring 2022  
Homework 3  
(Due Tuesday, May. 24 at 11:59pm (CST))

1. [15 points] Given a Bayesian network (Fig. 1) with five discrete variables  $\{F, A, S, H, N\}$ , where  $\{F, A, S, H, N\}$  are boolean variables. Suppose that  $\{F, A, H, N\}$  are observed variables and  $\{S\}$  is a latent variable. Now we implement EM algorithm for this model.

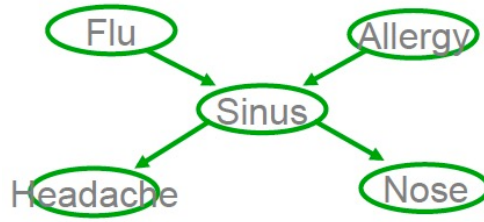


Figure 1: The Bayesian network with five discrete variables  $\{F, A, S, H, N\}$ .

- (a) Derive the E-step. [5 points]
- (b) Derive the M-step. [5 points]
- (c) Guess the solution of parameter estimation in the M-step, according to the MLE solution with all variables being observed. [5 points]

**Solution:**

- (a) In E-step, we need to compute  $P(S|F, A, H, N, \theta)$   
Since  $\{F, A, S, H, N\}$  are boolean variables,

$$P(S_k = 0 | f_k, a_k, h_k, n_k, \theta) = \frac{P(S_k = 0, f_k, a_k, h_k, n_k | \theta)}{P(S_k = 0, f_k, a_k, h_k, n_k | \theta) + P(S_k = 1, f_k, a_k, h_k, n_k | \theta)}$$

$$P(S_k = 1 | f_k, a_k, h_k, n_k, \theta) = \frac{P(S_k = 1, f_k, a_k, h_k, n_k | \theta)}{P(S_k = 0, f_k, a_k, h_k, n_k | \theta) + P(S_k = 1, f_k, a_k, h_k, n_k | \theta)}$$

$$\begin{aligned} P(S_k = 1 | f_k, a_k, h_k, n_k, \theta) &= E[s_k] \\ &= \frac{P(S_k = 1, f_k, a_k, h_k, n_k | \theta)}{P(S_k = 0, f_k, a_k, h_k, n_k | \theta) + P(S_k = 1, f_k, a_k, h_k, n_k | \theta)} \end{aligned}$$

- (b) In M-step, we need to compute the  $\theta'$  which maximize

$$\begin{aligned} &E_{P(S|F,A,H,N,\theta)} \log[P(F, A, S, H, N | \theta')] \\ &= \sum_{k=1}^K \sum_{i=0}^1 P(S_k = i | f_k, a_k, h_k, n_k, \theta) [\log P(f_k) + \log P(a_k) + \log P(s_k | f_k, a_k) + \log P(h_k | s_k) + \log P(n_k | s_k)] \end{aligned}$$

- (c) MLE solution:

$$\theta_f = \frac{\sum_{k=1}^K \delta(f_k = 1)}{K}$$

$$\begin{aligned}
\theta_a &= \frac{\sum_{k=1}^K \delta(a_k = 1)}{K} \\
\theta_{s|f,a} &= \frac{\sum_{k=1}^K \delta(s_k = s, f_k = f, a_k = a)}{\sum_{k=1}^K \delta(f_k = f, a_k = a)} \\
\theta_{h|s} &= \frac{\sum_{k=1}^K \delta(h_k = 1, s_k = s)}{\sum_{k=1}^K \delta(s_k = s)} \\
\theta_{n|s} &= \frac{\sum_{k=1}^K \delta(n_k = 1, s_k = s)}{\sum_{k=1}^K \delta(s_k = s)}
\end{aligned}$$

M-step solution:

$$\begin{aligned}
\theta_f &= \frac{\sum_{k=1}^K \delta(f_k = 1)}{K} \\
\theta_a &= \frac{\sum_{k=1}^K \delta(a_k = 1)}{K} \\
\theta_{s|f,a} &= \frac{\sum_{k=1}^K \delta(f_k = f, a_k = a) P(s_k = s)}{\sum_{k=1}^K \delta(f_k = f, a_k = a)} \\
\theta_{h|s} &= \frac{\sum_{k=1}^K \delta(h_k = 1) P(s_k = s)}{\sum_{k=1}^K P(s_k = s)} \\
\theta_{n|s} &= \frac{\sum_{k=1}^K \delta(n_k = 1) P(s_k = s)}{\sum_{k=1}^K P(s_k = s)}
\end{aligned}$$

2. [20 points] Suppose two data points ( $x_1 = 0, x_2 = 1$ ) are generated from two Gaussian mixture model (A and B). The parameter of the two Gaussian model are unknown. We want to use EM to guess parameters of the two Gaussian models. For simplicity, the priors are set to equal, which means  $P(a) = P(b) = \frac{1}{2}$ . EM can be divided into following steps:

- Randomly choose  $(\mu_a, \sigma_a^2)$  and  $(\mu_b, \sigma_b^2)$ .
- For each point  $x_i$ , calculate  $P(a|x_i)$  and  $P(b|x_i)$ .
- Adjust  $(\mu_a, \sigma_a^2)$  and  $(\mu_b, \sigma_b^2)$ .
- Repeat 2 and 3 until convergence.

Suppose we randomly choose parameters as:  $(\mu_a, \sigma_a^2) = (0, 1)$  and  $(\mu_b, \sigma_b^2) = (1, 1)$

- E-step: calculate  $P(a|x_i)$  and  $P(b|x_i)$ ,  $i = 1, 2$ . [10 points]

- M-step: Adjust  $(\mu_a, \sigma_a^2)$  and  $(\mu_b, \sigma_b^2)$  with following formula. [10 points]

$$\mu_a = \frac{a_1 x_1 + a_2 x_2}{a_1 + a_2}, \quad \sigma_a^2 = \frac{a_1 (x_1 - \mu_a)^2 + a_2 (x_2 - \mu_a)^2}{a_1 + a_2}$$

where  $a_1 = P(a|x_1)$  and  $a_2 = P(a|x_2)$

**Solution:**

(1)

$$\begin{aligned} P(a|x_i) &= P(a | x_i) \\ &= \frac{P(x_i | a) P(a)}{P(x_i)} \\ &= \frac{1}{2} \frac{P(x_i | a) P(a)}{P(x_i | a) P(a) + P(x_i | b) P(b)} \\ &= \frac{P(x_i | a)}{P(x_i | a) + P(x_i | b)} \end{aligned}$$

similarly,

$$P(b|x_i) = \frac{P(x_i | b)}{P(x_i | a) + P(x_i | b)}$$

By Gaussian mixture model, we know

$$\begin{aligned} P(x_i | a) &= \frac{1}{\sigma_a \sqrt{2\pi}} e^{-\frac{(x_i - \mu_a)^2}{2\sigma_a^2}} \\ P(x_i | b) &= \frac{1}{\sigma_b \sqrt{2\pi}} e^{-\frac{(x_i - \mu_b)^2}{2\sigma_b^2}} \end{aligned}$$

Bringing in the values gives,

$$\begin{aligned} P(a | x_1) &= \frac{1}{1 + e^{-\frac{1}{2}}} \\ P(b | x_1) &= \frac{e^{-\frac{1}{2}}}{1 + e^{-\frac{1}{2}}} \\ P(a | x_2) &= \frac{e^{-\frac{1}{2}}}{1 + e^{-\frac{1}{2}}} \\ P(b | x_2) &= \frac{1}{1 + e^{-\frac{1}{2}}} \end{aligned}$$

- $a_1 = P(a|x_1) = \frac{1}{1+e^{-\frac{1}{2}}}$ ,  $a_2 = P(a | x_2) = \frac{e^{-\frac{1}{2}}}{1+e^{-\frac{1}{2}}}$

$$\begin{aligned} \mu_a &= \frac{a_1 x_1 + a_2 x_2}{a_1 + a_2} \\ &= \frac{a_2}{a_1 + a_2} \\ &= a_2 \\ &= \frac{e^{-\frac{1}{2}}}{1 + e^{-\frac{1}{2}}} \end{aligned}$$

$$\begin{aligned}
\sigma_a^2 &= \frac{a_1(x_1 - \mu_a)^2 + a_2(x_2 - \mu_a)^2}{a_1 + a_2} \\
&= a_1\mu_a^2 + a_2(1 - \mu_a)^2 \\
&= \frac{e^{-\frac{1}{2}}}{(1 + e^{-\frac{1}{2}})^2}
\end{aligned}$$

$$b_1 = P(b \mid x_1) = \frac{e^{-\frac{1}{2}}}{1 + e^{-\frac{1}{2}}}, \quad b_2 = P(b \mid x_2) = \frac{1}{1 + e^{-\frac{1}{2}}}$$

$$\begin{aligned}
\mu_b &= \frac{b_1x_1 + b_2x_2}{b_1 + b_2} \\
&= \frac{b_2}{b_1 + b_2} \\
&= b_2 \\
&= \frac{1}{1 + e^{-\frac{1}{2}}}
\end{aligned}$$

$$\begin{aligned}
\sigma_b^2 &= \frac{b_1(x_1 - \mu_b)^2 + b_2(x_2 - \mu_b)^2}{b_1 + b_2} \\
&= b_1\mu_b^2 + b_2(1 - \mu_b)^2 \\
&= \frac{e^{-\frac{1}{2}}}{(1 + e^{-\frac{1}{2}})^2}
\end{aligned}$$

Table 1: The training data in (a).

$i$	$x_{i1}$	$x_{i2}$	$y_i$
1	1.5	0.5	1
2	2.5	1.5	1
3	3.5	3.5	1
4	6.5	5.5	1
5	7.5	10.5	1
6	1.5	2.5	-1
7	3.5	1.5	-1
8	5.5	5.5	-1
9	7.5	8.5	-1
10	1.5	10.5	-1

3. [30 points] Suppose that we are interested in learning a classifier, such that at any turn of a game we can pose a question, like "should I attack this ant hill now?", and get an answer. That is, we want to build a classifier which we can feed some features on the current game state, and get the output "attack" or "don't attack". There are many possible ways to define what the action "attack" means, but for now let's define it as sending all friendly ants that can see the ant hill under consideration towards it.

Let's recall the AdaBoost algorithm described in class. Its input is a dataset  $\{(x_i, y_i)\}_{i=1}^n$ , with  $x_i$  being the  $i$ -th sample, and  $y_i \in \{-1, 1\}$  denoting the  $i$ -th label,  $i = 1, 2, \dots, n$ . The features might be composed of a count of the number of friendly ants that can see the ant hill under consideration, and a count of the number of enemy ants these friendly ants can see. For example, if there were 10 friendly ants that could see a particular ant hill, and 5 enemy ants that the friendly ants could see, we would have:

$$x_1 = \begin{bmatrix} 10 \\ 5 \end{bmatrix}.$$

The label of the example  $x_1$  is  $y_1 = 1$ , once the friendly ants were successful in razing the enemy ant hill, and  $y_1 = 0$  otherwise. We could generate such examples by running a greedy bot (or any other opponent bot) against a bot that we make periodically try to attack an enemy ant hill. Each time this bot tries the attack, we record (say, after 20 turns or some other significant amount of time) whether the attack was successful or not.

- (a) Let  $\epsilon_t$  denote the error of a weak classifier  $h_t$ :

$$\epsilon_t = \sum_{i=1}^n D_t(i) \mathbb{1}(y_i \neq h_t(x_i)).$$

In the simple "attack" / "don't attack" scenario, suppose that we have implemented the following six weak classifiers:

$$\begin{aligned} h^{(1)}(x_i) &= 2 * \mathbb{1}(x_{i1} \geq 2) - 1, & h^{(4)}(x_i) &= 2 * \mathbb{1}(x_{i2} \leq 2) - 1, \\ h^{(2)}(x_i) &= 2 * \mathbb{1}(x_{i1} \geq 6) - 1, & h^{(5)}(x_i) &= 2 * \mathbb{1}(x_{i2} \leq 6) - 1, \\ h^{(3)}(x_i) &= 2 * \mathbb{1}(x_{i1} \geq 10) - 1, & h^{(6)}(x_i) &= 2 * \mathbb{1}(x_{i2} \leq 10) - 1. \end{aligned}$$

Given ten training data points ( $n = 10$ ) as shown in Table 1, please show that what is the minimum value of  $\epsilon_1$  and which of  $h^{(1)}, \dots, h^{(6)}$  achieve this value? Note that there may be multiple classifiers that all have the same  $\epsilon_1$ . You should list all classifiers that achieve the minimum  $\epsilon_1$  value. [6 points]

- (b) For all the questions in the remainder of this section, let  $h_1$  denote  $h^{(1)}$  chosen in the first round of boosting. (That is,  $h^{(1)}$  was the classifier that achieved the minimum  $\epsilon_1$ .)

- (1) What is the value of  $\alpha_1$  (the weight of this first classifier  $h_1$ )? [2 points]

- (2) What should  $Z_t$  be in order to make sure the distribution  $D_{t+1}$  is normalized correctly? That is, derive the formula of  $Z_t$  in terms of  $\epsilon_t$  that will ensure  $\sum_{i=1}^n D_{t+1}(i) = 1$ . Please also derive the formula of  $\alpha_t$  in terms of  $\epsilon_t$ . [6 points]

- (3) Which points will increase in significance in the second round of boosting? That is, for which points will we have  $D_1(i) < D_2(i)$ ? What are the values of  $D_2$  for these points? [5 points]
- (4) In the second round of boosting, the weights on the points will be different, and thus the error  $\epsilon_2$  will also be different. Which of  $h^{(1)}, \dots, h^{(6)}$  will minimize  $\epsilon_2$ ? (Which classifier will be selected as the second weak classifier  $h_2$ ?) What is its value of  $\epsilon_2$ ? [6 points]
- (5) What will the average error of the final classifier  $H$  be, if we stop after these two rounds of boosting? That is, if  $H(x) = \text{sign}(\alpha_1 h_1(x) + \alpha_2 h_2(x))$ , what will the training error  $\epsilon = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i \neq h(x_i))$  be? Is this more, less, or the same as the error we would get, if we just used one of the weak classifiers instead of this final classifier  $H$ ? [5 points]

**Solution:**

- (a) By six different weak classifiers, the  $\epsilon_1$  values of  $h^{(1)}, \dots, h^{(6)}$  is 4/10, 4/10, 5/10, 4/10, 4/10, 5/10 which equals 0.4, 0.4, 0.5, 0.4, 0.4, 0.5 the minimum one is 0.4, the 1st, 2nd, 4th, 5th classifiers achieve the minimum.

(b)

- (1) Using formula

$$\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$$

Then,

$$\begin{aligned} \alpha_1 &= \frac{1}{2} \ln \left( \frac{1 - \epsilon_1}{\epsilon_1} \right) \\ &= \frac{1}{2} \ln \left( \frac{3}{2} \right) \\ &\approx 0.2027 \end{aligned}$$

- (2) Since  $D_{t+1}(i) = \frac{D_t(i)}{Z_t} e^{\{-\alpha_t y_i h_t(x_i)\}}$

$$\begin{aligned} D_{t+1}(i) &= \frac{D_t(i)}{Z_t} e^{\{-\alpha_t y_i h_t(x_i)\}} \\ 1 &= \sum_{i=1}^n D_{t+1}(i) = \sum_{i=1}^n \frac{D_t(i)}{Z_t} e^{\{-\alpha_t y_i h_t(x_i)\}} \\ Z_t &= \sum_{i=1}^n D_t(i) e^{\{-\alpha_t y_i h_t(x_i)\}} \end{aligned}$$

Since

$$\begin{aligned} D_{t+1}(i) &= \frac{D_t(i)}{Z_t} e^{\{-\alpha_t\}} \text{ if } y_i = h_t(x_i) \\ D_{t+1}(i) &= \frac{D_t(i)}{Z_t} e^{\{\alpha_t\}} \text{ if } y_i \neq h_t(x_i) \\ Z_t &= \sum_{i=1}^n D_t(i) e^{\{-\alpha_t y_i h_t(x_i)\}} \\ &= \sum_{y_i = h_t(x_i)} D_t(i) e^{\{-\alpha_t\}} + \sum_{y_i \neq h_t(x_i)} D_t(i) e^{\{\alpha_t\}} \\ &= (1 - \epsilon_t) e^{\{-\alpha_t\}} + \epsilon_t e^{\{\alpha_t\}} \end{aligned}$$

To minimize the training error,

$$\alpha_t = \underset{\alpha}{\operatorname{argmin}} [(1 - \epsilon_t) e^{\{-\alpha_t\}} + \epsilon_t e^{\{\alpha_t\}}]$$

By  $\frac{\partial Z_t}{\partial \alpha_t} = 0$ , we know  $\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$ . Then,  $Z_t = 2\sqrt{\epsilon_t(1 - \epsilon_t)}$

- (3) The points which misclassifies will increase in weight in  $D_2$ . They are 1st, 7th, 8th, 9th ( $i = 1, 7, 8, 9$ ) points. For these points,

$$\begin{aligned} D_2(i) &= \frac{D_1(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t} \\ &= \frac{\frac{1}{10} \exp\{\alpha_1\}}{4 * \exp\{\alpha_1\} + 6 * \exp\{-\alpha_1\}} \\ &= \frac{1}{8} \end{aligned}$$

which means  $D_2(1) = D_2(7) = D_2(8) = D_2(9) = \frac{1}{8}$

- (4) By (3), similarly, for  $i = 2, 3, 4, 5, 6, 10$ ,  $D_2(i) = \frac{1}{12}$ .  
the error values of  $\epsilon_2$  for  $h^{(1)}, \dots, h^{(6)}$  is  $\frac{4}{8}, \frac{2}{8} + \frac{2}{12}, \frac{1}{8} + \frac{4}{12}, \frac{1}{8} + \frac{3}{12}, \frac{2}{8} + \frac{2}{12}, \frac{3}{8} + \frac{2}{12}$  which equals  $\frac{1}{2}, \frac{5}{12}, \frac{11}{24}, \frac{3}{8}, \frac{5}{12}, \frac{13}{24}$   
Thus,  $h_4$  will minimize  $\epsilon_2$ , the value is  $\frac{3}{8}$ .

- (5) By (4) and (2),  $\alpha_2 = \frac{1}{2} \ln \left( \frac{1-\epsilon_2}{\epsilon_2} \right) = \frac{1}{2} \ln \frac{7}{5}$

Thus,

$$\begin{aligned} H(x) &= \text{sign}(\alpha_1 h_1(x) + \alpha_2 h_4(x)) \\ &= \text{sign}(h_1(x) \ln \frac{3}{2} + h_4(x) \ln \frac{7}{5}) \end{aligned}$$

Since  $\alpha_1 \neq \alpha_2$  and  $h_t(x_i) = 1$  or  $-1$ ,  $H(x)$  equals to one of  $h_1(x), h_4(x)$  ( $H(x) = \text{sign}(\alpha_1 h_1(x) + \alpha_2 h_4(x))$ )  
sign of it is only about the large coefficient one  $h_t(x)$

Since  $\alpha_1 > \alpha_2$ , we know  $H(x) = h_1(x)$ , which means  $\epsilon = \epsilon_1 = 0.4$

4. [15 points] Given valid kernels  $k_1(\mathbf{x}, \mathbf{x}')$  and  $k_2(\mathbf{x}, \mathbf{x}')$ , please verify the following new kernels will also be valid:
- (a)  $k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$ , where  $f(\cdot)$  is any function. [5 points]
  - (b)  $k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}'))$ , where  $q(\cdot)$  is a polynomial with nonnegative coefficients. [5 points]
  - (c)  $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{A} \mathbf{x}'$ , where  $\mathbf{A}$  is a symmetric positive semi-definite matrix. [5 points]

**Solution:**

- (a) Since  $k_1(x, x')$  can be written as  $\langle \phi(x), \phi(x') \rangle$ ,

$$k(x, x') = f(x)\phi(x)^T \phi(x')f(x') = \langle \phi(x)f(x)^T, \phi(x')f(x') \rangle$$

it can be written as inner product of feature vectors, so it is valid kernel.

- (b) First, product of valid kernels is valid kernel, sum of valid kernels is valid kernel.  
So,

$$q(x) = \sum_{i=1}^n a_n x^n$$

$$k(x, x') = \sum_{i=1}^n a_n (k_1(x, x'))^n$$

Thus, it is a valid kernel.

- (c) Eigendecomposition of  $A$  is  $Q\Lambda Q^T$ , and  $\Lambda$  is diagonal matrix. So,

$$\begin{aligned} k(x, x') &= x^T A x' \\ &= x^T Q \Lambda Q^T x \\ &= (\Lambda^{0.5} Q^T x)^T (\Lambda^{0.5} Q^T x') \\ &= \langle (\Lambda^{0.5} Q^T x), (\Lambda^{0.5} Q^T x') \rangle \end{aligned}$$

it can be written as inner product of feature vectors, so it is valid kernel.



5. [20 points] We have learned that when solving a SVM problem, we need to first construct Lagrangian function  $L(w, b, \alpha)$  and set partial derivative to zero. By using KKT conditions, we can get the dual problem of SVM :

$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j < x_i, x_j > - \sum_{i=1}^n \alpha_i,$$

$$s.t. \alpha_i \geq 0, \quad \sum_{i=1}^n \alpha_i y_i = 0$$

Now we use a simple example to better understand how SVM works. We consider the separating hyperplane being  $wx + b = 0$ . Suppose we have three data points:  $x_1 = (2, -1)^T, x_2 = (2, -3)^T, x_3 = (4, -1)^T$ , the corresponding labels are:  $y_1 = -1, y_2 = -1, y_3 = 1$ . Use SVM to find the values of  $w^* = (w_1^*, w_2^*)$ ,  $b^*$  and give the separating hyperplane. Please show your calculation process.

**Solution:**

First,  $y_1 = -1, y_2 = -1, y_3 = 1$ , so  $\alpha_1 + \alpha_2 - \alpha_3 = 0$ .

Then by  $\min_{\alpha} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j < x_i, x_j > - \sum_{i=1}^n \alpha_i$ ,

we know

$$\alpha = \min_{\alpha} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j < x_i, x_j > - \sum_{i=1}^n \alpha_i$$

$$= \min_{\alpha} 5\alpha_2^2 + 2\alpha_3^2 - 2\alpha_3$$

Thus,  $\alpha_2 = 0$  and  $\alpha_3 = \frac{1}{2}, \alpha_1 = \frac{1}{2}$ .

$$w^* = \sum_i \alpha_i y_i x_i$$

$$= (1, 0)$$

$$y_i = w^* x_i + b \Rightarrow b^* = -3$$

the separating hyperplane:  $0 = x_1 - 3$