Problem Set 4

MAS 622J/1.126J: Pattern Recognition and Analysis

Due: 5:00 p.m. on October 22

[Note: All instructions to plot data or write a program should be carried out using Matlab. In order to maintain a reasonable level of consistency and simplicity we ask that you do not use other software tools.]

If you collaborated with other members of the class, please write their names at the end of the assignment. Moreover, you will need to write and sign the following statement: "In preparing my solutions, I did not look at any old homeworks, copy anybody's answers or let them copy mine."

Problem 1: Expectation Maximization and Missing Data [20 points]

Consider data, $D = \left\{ \begin{pmatrix} 1 \\ 3 \end{pmatrix}, \begin{pmatrix} 4 \\ 5 \end{pmatrix}, \begin{pmatrix} 2 \\ * \end{pmatrix} \right\}$, sampled from a two-dimensional (separable) distribution, $p(x_1, x_2) = p_{x_1}(x_1)p_{x_2}(x_2)$, with

$$p_{x_1}(x_1) = \begin{cases} \frac{1}{\theta_1} e^{-x_1/\theta_1} & \text{if } x_1 \ge 0\\ 0 & \text{otherwise} \end{cases} \qquad p_{x_2}(x_2) = \begin{cases} \frac{1}{\theta_2} & \text{if } 0 \le x_2 \le \theta_2\\ 0 & \text{otherwise} \end{cases}$$

and a missing feature value, *.

a. What can you infer from θ_2 by looking at D?

Solution: θ_2 is necessarily greater or equal to 5 because $0 < x_{12} = 3 < 0$

$$x_{22} = 5 \le \theta_2.$$

b. Start with an initial estimate $\underline{\theta}^0 = \begin{pmatrix} 3 \\ 6 \end{pmatrix}$ and analytically calculate $Q(\underline{\theta}; \underline{\theta}^0)$. This is the *estimate* step in the EM algorithm.

Solution:

The probability of one sample can be determined by:

$$\begin{split} p(\mathbf{x}|\underline{\theta}) &= p(x_1, x_2|\underline{\theta}) = p_{x1}(x_1|\underline{\theta}) p_{x_2}(x_2|\underline{\theta}) \\ &= \left\{ \begin{array}{l} \left(\frac{1}{\theta_1} e^{-x_1/\theta_1}\right) \left(\frac{1}{\theta_2}\right) & \text{if } x_1 \geq 0 \text{ and } 0 \leq x_2 \leq \theta_2 \\ 0 & \text{otherwise} \end{array} \right. \\ &= \left\{ \begin{array}{l} \frac{1}{\theta_1 \theta_2} e^{-x_1/\theta_1} & \text{if } x_1 \geq 0 \text{ and } 0 \leq x_2 \leq \theta_2 \\ 0 & \text{otherwise} \end{array} \right. \end{split}$$

Formulating the log likelihood of the data and marginalizing over the possible values of the unknown data:

$$\begin{split} Q(\underline{\theta};\underline{\theta}^{0}) &= E_{x_{32}}[\ln p(\mathbf{x}_{g},\mathbf{x}_{b};\underline{\theta})|\underline{\theta}^{0};D_{g}] \\ &= \int_{-\infty}^{\infty} \left[\sum_{k=1}^{2} \ln p(\mathbf{x}_{k}|\underline{\theta}) + \ln p(\mathbf{x}_{3}|\underline{\theta}) \right] p(x_{32}|x_{31} = 2;\underline{\theta}^{0}) dx_{32} \\ &= \sum_{k=1}^{2} \ln p(\mathbf{x}_{k}|\underline{\theta}) + \int_{-\infty}^{\infty} \ln p\left(\left[\begin{array}{c} 2 \\ x_{32} \end{array} \right] |\underline{\theta} \right) \frac{p\left(\left[\begin{array}{c} 2 \\ x_{32} \end{array} \right] |\underline{\theta}^{0} \right)}{\int_{-\infty}^{\infty} p\left(\left[\begin{array}{c} 2 \\ x_{32} \end{array} \right] |\underline{\theta}^{0} \right) dx_{32}'} dx_{32} \end{split}$$

Note that x_{31} and x_{32} are independent with each other. Thus,

$$p\left(\begin{bmatrix}2\\x_{32}\end{bmatrix}|\underline{\theta}^{0}\right) = p_{x_{1}}\left(x_{31} = 2|\underline{\theta}^{0}\right)p_{x_{2}}(x_{32}|\underline{\theta}^{0}),$$

$$\int_{-\infty}^{\infty} p\left(\begin{bmatrix}2\\x'_{32}\end{bmatrix}|\underline{\theta}^{0}\right)dx'_{32} = p_{x_{1}}\left(x_{31} = 2|\underline{\theta}^{0}\right)\int_{-\infty}^{\infty} p_{x_{2}}(x'_{32}|\underline{\theta}^{0})dx'_{32} = p_{x_{1}}\left(x_{31} = 2|\underline{\theta}^{0}\right)$$

Therefore,

$$Q(\underline{\theta}; \underline{\theta}^0) = \sum_{k=1}^{2} \ln p(\mathbf{x}_k | \underline{\theta}) + \int_{-\infty}^{\infty} \ln p\left(\begin{bmatrix} 2 \\ x_{32} \end{bmatrix} | \underline{\theta}\right) p_{x_2}(x_{32} | \underline{\theta}^0) dx_{32}$$

Note that $p\left(\begin{bmatrix} 2 \\ x_{32} \end{bmatrix} | \underline{\theta}\right)$ is different than zero if $0 \leq x_{32} \leq \theta_2$, and $p_{x_2}(x_{32} | \underline{\theta}^0)$ is different than zero if $0 \leq x_{32} \leq \theta_2^0 = 6$.

Therefore, $p\left(\begin{bmatrix}2\\x_{32}\end{bmatrix}|\underline{\theta}\right)p_{x_2}(x_{32}|\underline{\theta}^0)$ is different than zero between $[0,\min(\theta_2,6)]$

$$\sum_{k=1}^{2} \ln p(\mathbf{x}_{k}|\underline{\theta}) = \ln \frac{1}{\theta_{1}\theta_{2}} e^{-1/\theta_{1}} + \ln \frac{1}{\theta_{1}\theta_{2}} e^{-4/\theta_{1}} = 2 \ln \frac{1}{\theta_{1}\theta_{2}} - \frac{5}{\theta_{1}}$$

$$\begin{split} \int_{-\infty}^{\infty} \ln p \left(\left[\begin{array}{c} 2 \\ x_{32} \end{array} \right] |\underline{\theta} \right) p_{x_{32}} \left(x_{32} | \underline{\theta}^0 \right) dx_{32} &= \int_{0}^{\min(\theta_2, 6)} \left(\ln \frac{1}{\theta_1 \theta_2} e^{-2/\theta_1} \right) \left(\frac{1}{6} \right) dx_{32} \\ &= \begin{cases} \theta_2 \frac{1}{6} \left(\ln \frac{1}{\theta_1 \theta_2} - \frac{2}{\theta_1} \right) & \text{if } \min(\theta_2, 6) = \theta_2, \quad \text{where } 5 \leq \theta_2 < 6 \\ \ln \frac{1}{\theta_1 \theta_2} - \frac{2}{\theta_1} & \text{if } \min(\theta_2, 6) = 6, \quad \text{where } \theta_2 \geq 6 \end{cases} \end{split}$$

Finally,

$$Q(\underline{\theta}; \underline{\theta}^{0}) = \sum_{k=1}^{2} \ln p(\mathbf{x}_{k}|\underline{\theta}) + \int_{-\infty}^{\infty} \ln p\left(\begin{bmatrix} 2 \\ x_{32} \end{bmatrix} |\underline{\theta}\right) p_{x_{32}}\left(x_{32}|\underline{\theta}^{0}\right) dx_{32}$$

$$= \begin{cases} 2 \ln \frac{1}{\theta_{1}\theta_{2}} - \frac{5}{\theta_{1}} + \frac{1}{6}\theta_{2}\left(\ln \frac{1}{\theta_{1}\theta_{2}} - \frac{2}{\theta_{1}}\right) & \text{where } 5 \leq \theta_{2} < 6 \\ 3 \ln \frac{1}{\theta_{1}\theta_{2}} - \frac{7}{\theta_{1}} & \text{where } \theta_{2} \geq 6 \end{cases}$$

c. Find the $\underline{\theta}$ that maximizes your $Q(\underline{\theta}; \underline{\theta}^0)$ – the maximization step of the EM algorithm.

Solution:

First,

$$\frac{\partial Q(\underline{\theta};\underline{\theta}^0)}{\partial \theta_2} = \begin{cases} -\frac{2}{\theta_2} - \frac{1}{6} \left(\frac{2}{\theta_1} - \ln \frac{1}{\theta_1 \theta_2} \right) - \frac{1}{6} & \text{where } 5 \le \theta_2 < 6 \\ -\frac{3}{\theta_2} & \text{where } \theta_2 \ge 6 \end{cases}$$

Thus, $Q(\underline{\theta}; \underline{\theta}^0) < 0$ where $\theta_2 \geq 5$. (Note $\frac{2}{\theta_1} - \ln \frac{1}{\theta_1 \theta_2} = \frac{1}{\theta_1} + (\frac{1}{\theta_1} - \ln \frac{1}{\theta_1}) + \ln \theta_2 > 0$ because $\frac{1}{\theta_1} > 0$, $(\frac{1}{\theta_1} - \ln \frac{1}{\theta_1}) \geq 1$ and $\ln \theta_2 > 0$ where $\theta_1 > 0$ and $\theta_2 \geq 5$.)

Second, $Q(\underline{\theta};\underline{\theta}^0)$ is continuous at $\theta_2 = 6$, because the two equations corresponding to the two intervals $(5 \le \theta_2 < 6, \theta_2 \ge 6)$ have the same value at $\theta_2 = 6$, which is $3 \ln \frac{1}{6\theta_1} - \frac{7}{\theta_1}$.

From these two facts, we can conclude that $Q(\underline{\theta}; \underline{\theta}^0)$ is a monotonically decreasing continuous function with respect to θ_2 . Thus, the maximum occurs when $\theta_2 = 5$.

To find the value of θ_1 that maximizes Q when $\theta_2 = 5$, $\frac{\partial Q(\underline{\theta};\underline{\theta}^0)}{\partial \theta_1} = -\frac{2}{\theta_1} + \frac{5}{\theta_1^2} - \frac{5}{6\theta_1} + \frac{10}{6\theta_1^2} = 0$. Thus, we find that $\underset{\theta_1}{\operatorname{arg max}} Q(\underline{\theta};\underline{\theta}^0) = \frac{40}{17}$ (i.e., $\theta_1 = \frac{40}{17}$ and $\theta_2 = 5$).

Problem 2: Baum-Welch Algorithm and Discrete HMMs [40 points]

Download the datasets from the course webpage. The datasets consist of training and testing sequences belonging to two classes. We assume the two HMMs for the two classes have the same configuration, i.e. the same number of states, zero transition probabilities and the number of output states.

Implement the Baum-Welch algorithm for training a discrete HMM. Train HMMs with one, three, and five states with transition probabilities in a strictly left-to-right configuration (see the figure below for a two-state HMM in left-to-right configuration). The visible output has four possible states 0, 1, 2 or 3. and

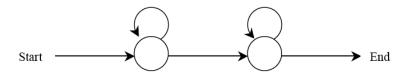


Figure 1: HMM in a left-to-right configuration

Repeat the following steps for each of the three HMM configurations with one, three, and five states:

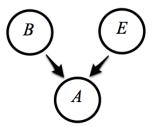
- a. Train two HMMs, one for each class of data. State very clearly the threshold you are using and the maximum number of iterations. List the output probabilities and state transition probabilities of each HMM.
- b. Implement the Viterbi algorithm to decode each test sequence using both HMMs. Show the log probability of each test sequence using each HMM.
- c. Show the confusion matrix of the test set.

Include a complete listing of your source code. **Solution:**

The MATLAB code for this problem can be downloaded from the course website. The code was developed by Ashish Kapoor (ash@media.mit.edu) ¹.

Problem 3: Bayes Networks [20 points]

Consider the following Bayesian Network where the three events are B (Burglary), E (Earthquake) and A (Alarm). Assume that the three nodes are binary nodes that can take on the value f (false) or t (true).



Here, B and E are marginally independent. That is, P(B,E) = P(B)P(E) (Note that this relation implies P(B|E) = P(B) or P(E|B) = P(E).) Now we want to show B and E are conditionally dependent given A. In other words, $P(B,E|A) \neq P(B|A)P(E|A)$.

¹Copyright 2002 by the Massachusetts Institute of Technology. All rights reserved. Permission to use, copy, or modify this software and its documentation for educational and research purposes only and without fee is hereby granted, provided that this copyright notice and the original authors' names appear on all copies and supporting documentation.

a. Prove that the relation $P(B, E|A) \neq P(B|A)P(E|A)$ implies $P(B|E, A) \neq P(B|A)$. Note that by the symmetry between B and E, this also means $P(E|B, A) \neq P(E|A)$.

Solution: From $P(B, E|A) \neq P(B|A)P(E|A)$, we get the relation

$$\frac{P(B,E|A)}{P(E|A)} \neq P(B|A)$$

Also, we know that $P(B|E,A) = \frac{P(B,E|A)}{P(E|A)}$ by Bayes' rule. Therefore, $P(B|E,A) \neq P(B|A)$

b. First, we consider a very simple case. Suppose A = B + E where the + sign means "Logical OR." This means that B and E are independent deterministic causes of A. Construct the CPT (conditional probability table) for $P(A \mid B, E)$.

c. For the obtained CPT above, prove that P(B = t|A = t) = P(B = t)/P(A = t). Which is greater between P(B = t|A = t) and P(B = t)? What is the meaning of this?

Solution: P(B = t | A = t) = P(A = t, B = t)/P(A = t) by Bayes'

rule.

Now for
$$P(A=t,B=t)$$
,
$$P(A=t,B=t) = \sum_{E=t,f} P(A=t,B=t,E) \text{ (by marginalization)}$$

$$= \sum_{E=t,f} P(B=t)P(E)P(A=t|B=t,E) \text{ (from the graphical model)}$$

$$= P(B=t)P(E=t)P(A=t|B=t,E=t) + P(B=t)P(E=f)P(A=t|B=t,E=f)$$

$$= P(B=t)P(E=t) + P(B=t)P(E=f)$$

$$= P(B=t) (P(E=t) + P(E=f))$$

$$= P(B=t)$$

Therefore, P(B = t|A = t) = P(B = t)/P(A = t)

Since
$$P(A = t) \le 1$$
, $P(B = t | A = t) \ge P(B = t)$

Note that P(B=t) means your belief on B=t (Burglary happened) without being given any information about E and A. However, P(B=t|A=t) means your belief on B=t (Burglary happened) when you know A=t (Alarm happened). Thus, $P(B=t|A=t) \geq P(B=t)$ means that your belief of B=t (Burglary happened) increases when you know A=t (Alarm happened).

d. For the obtained CPT above, prove that P(A = t) = P(B = t) + P(E = t) - P(B = t)P(E = t). If P(B = t) and P(E = t) are small, what happens to P(B = t|A = t), compared with P(B = t)?

Solution:

$$\begin{split} &P(A=t) \\ &= \sum_{B=t,f} \sum_{E=t,f} P(A=t,B,E) \\ &= \sum_{B=t,f} \sum_{E=t,f} P(B)P(E)P(A=t|B,E) \\ &= P(B=t)P(E=t)P(A=t|B=t,E=t) + P(B=t)P(E=f)P(A=t|B=t,E=f) \\ &+ P(B=f)P(E=t)P(A=t|B=f,E=t) + P(B=f)P(E=f)P(A=t|B=f,E=f) \\ &= P(B=t)P(E=t) + P(B=t)P(E=f) + P(B=f)P(E=t) \\ &= P(B=t)\left(P(E=t) + P(E=f)\right) + (1-P(B=t))P(E=t) \\ &= P(B=t) + P(E=t) - P(B=t)P(E=t) \end{split}$$

Thus, as P(B=t) and P(E=t) gets smaller, P(A=t) also gets smaller. Since P(B=t|A=t) = P(B=t)/P(A=t) (from (d)), P(B=t|A=t) becomes much greater than P(B=t).

e. For the obtained CPT above, prove that P(B=t|E=t,A=t)=P(B=t). Which is greater between P(B=t|E=t,A=t) and P(B=t|A=t)? What is the meaning of this? Hint: Note that the observation of Earthquake (E=t) alone is enough to explain the cause of Alarm (A=t). This occurrence is called "Explaining Away."

Solution:

$$\begin{split} &P(B=t|E=t,A=t)\\ &=\frac{P(B=t,E=t,A=t)}{P(E=t,A=t)}\\ &=\frac{P(B=t)P(E=t)P(A=t|B=t,E=t)}{P(E=t,A=t)} & \text{(from the graphical model)}\\ &=\frac{P(B=t)P(E=t)P(A=t|B=t,E=t)}{P(E=t)} & \text{(since }P(E=t,A=t)=P(E=t))\\ &=P(B=t) \end{split}$$

Note that in (c), we proved P(B = t, A = t) = P(B = t). In like manner, we can prove P(E = t, A = t) = P(E = t).

Since
$$P(B=t|E=t,A=t)=P(B=t)$$
 and $P(B=t|A=t)\geq P(B=t)$ (from (c)), we get

$$P(B = t|A = t) > P(B = t) = P(B = t|E = t, A = t)$$

Note that P(B=t|E=t,A=t) denotes your belief on B=t (Burglary happened) when you know both Earthquake and Alarm happened (E=t,A=t). Here the observation of Earthquake (E=t) alone is enough to explain the cause of Alarm (A=t). Therefore, since the fact that Earthquake happened can "explain away" the fact that Alarm happened, there is no increase in your belief on B=t (Burglary happened). Thus, P(B=t|E=t,A=t)=P(B=t).

f. We assume that P(B=t) = 0.15 and P(E=t) = 0.005. Calculate P(A=t), P(B=t|A=t), P(B=t|E=t, A=t), P(E=t|A=t), P(E=t|B=t, A=t) by hand.

Solution:

$$\begin{split} P(A=t) &= P(B=t) + P(E=t) - P(B=t)P(E=t) = 0.15 + \\ 0.005 - 0.15 \times 0.005 = 0.15425 \\ P(B=t|A=t) &= P(B=t)/P(A=t) = 0.15/0.15425 = 0.97245 \\ P(B=t|E=t,A=t) &= P(B=t) = 0.15 \\ P(E=t|A=t) &= P(E=t)/P(A=t) = 0.005/0.15425 = 0.0324 \\ P(E=t|B=t,A=t) &= P(E=t) = 0.005 \end{split}$$