

6.867: Recitation Handout (Week 10)

November 10, 2016

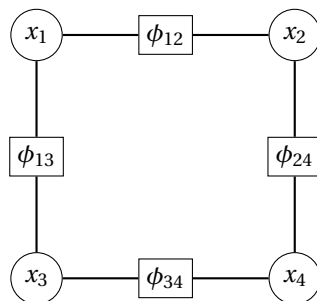
Contents

1	Learning graphical models using co-occurrence counts	2
2	EM for word counts	3
3	Buying widgets	4
4	Missing data	4
5	Mixed-up mixture	5
6	More mixture	6
7	NB: Nota bene	7
8	Parameterized CPT	8
9	What's up, doc?	9

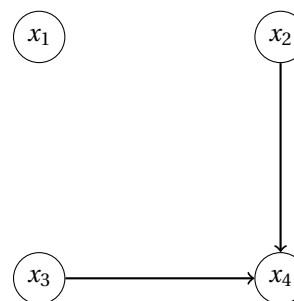
1 Learning graphical models using co-occurrence counts

We wanted to estimate the parameters of several different graphical models from medical data. For privacy reasons, we weren't allowed access to the actual records but only statistics computed from the records. We requested all pairwise counts involving four binary (0/1) variables x_1, \dots, x_4 . In other words, we have co-occurrence counts $\hat{n}_{ij}(x_i, x_j)$ for $i \neq j$ and $x_i \in \{0, 1\}, x_j \in \{0, 1\}$.

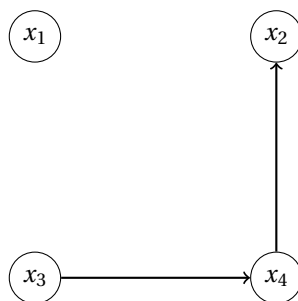
- (a) Consider the graphical models shown below. Check which of them we could estimate based on the pairwise counts. By estimation we mean finding the same maximum likelihood estimates of the parameters as we would if we had access to the full records. In answering this question, make no additional assumptions about the models above and beyond the graph structure.



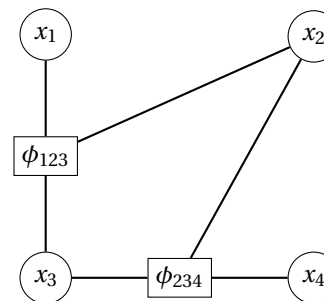
(a) **a** ()



(b) **b** ()



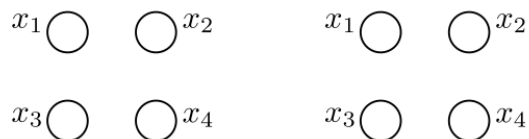
(c) **c** ()



(d) **d** ()

- (b) Provide a brief justification for your answer to model (b) above
- (c) We also considered building a Bayesian network model from expert knowledge. Based on our discussions with a prominent physician, we were able to extract the following properties concerning the four binary variables.
- x_1 and x_2 are independent
 - x_1 and x_3 are independent
 - x_3 and x_4 are independent

Draw two NOT equivalent Bayesian network models consistent with these properties avoiding (if possible) any additional assumptions not already implied by the above properties.



2 EM for word counts

Consider a probabilistic model of the following form:

$$p(w_2, c | w_1) = q(c | w_1)q(w_2 | c)$$

Here w_1 and w_2 are words, drawn from some set of possible words \mathcal{V} . c is a word class, which can take any value in the set $\{1, 2, \dots, k\}$ for some integer k . $q(c | w_1)$ and $q(w_2 | c)$ are the parameters of the model. We can interpret this a model where: (1) a class c is generated by word w_1 ; (2) w_2 is then generated by the class c chosen in step 1.

Under this model, we can derive

$$p(w_2 | w_1) = \sum_{c=1}^k q(c | w_1)q(w_2 | c)$$

This will be a model of the conditional probability of seeing the word w_2 given that the previous word in a sentence was w_1 .

- Write down an expression for $p(c | w_1, w_2)$ as a function of the q parameters.
- Say for each pair of words w_1, w_2 , $\text{count}(w_1, w_2)$ is the number of times w_1 is followed by w_2 in our training data. We are going to derive an EM algorithm for optimization of the following log likelihood function:

$$L(\theta) = \sum_{w_1, w_2} \text{count}(w_1, w_2) \log p(w_2 | w_1)$$

For given parameter values q , define $\text{count}(c | w_1)$ to be expected number of times that w_1 generates class c , and define $\text{count}(w_2 | c)$ to be the expected number of times w_2 is generated by class c . (Here expectation is taken with respect to the distribution defined by the q parameters.) State how $\text{count}(c | w_1)$ and $\text{count}(w_2 | c)$ can be calculated as a function of the q parameters:

- Now describe how the q parameters are recalculated in the EM algorithm, based on the $\text{count}(c | w_1)$ and $\text{count}(w_2 | c)$ counts derived in the previous part.
- Say we initialize the parameters to be $q(c | w_1) = 1/k$ for all w_1, c , and $q(w_2 | c) = 1/|\mathcal{V}|$ for all w_2, c . Given these initial parameter values, what parameter values will the EM algorithm converge to?

3 Buying widgets

You need widgets as components in a large system you are building. You have bought a large number of them, and you find that there is considerable variability in their longevity (length of time before failure), energy consumption, and reliability (correctness of output). You've studied the data and can find no independence relations among any of the variables.

- (a) Using three variables (L, E, and R), draw a Bayesian network model of their relationships.

Now, you discover that the widgets are actually manufactured by three different suppliers, and you have reason to believe that longevity, energy consumption, and reliability are independent given the manufacturer. You are unable to tell which widget came from which manufacturer, but you can observe L, E, and R.

- (b) Draw a new Bayesian network to describe this situation.
- (c) How would you learn the parameters for each model?
- (d) Which learned model would assign higher likelihood to the data? Why?
- (e) Why might you prefer each of these models?

4 Missing data

We'll start with a very simple problem, in which single attribute of a single data set is missing. There are two attributes, A and B, and this is our data set, \mathcal{D} :

i	A	B
1	1	1
2	1	1
3	0	0
4	0	0
5	0	0
6	0	H ***missing **
7	0	1
8	1	0

Assume the data is *missing completely at random* (MCAR): that is, that the fact that it is missing is independent of its value.

Our goal is to estimate $\Pr(A, B)$ from this data. We'd really like to find the maximum-likelihood parameter values, if we can. The likelihood is

$$\mathcal{L}(\theta) = \log \Pr(\mathcal{D}; \theta) = \log (\Pr(\mathcal{D}, H = 0; \theta) + \Pr(\mathcal{D}, H = 1; \theta)) \quad .$$

- (a) Kim is lazy and decides to ignore $x^{(6)}$ all together, and estimate the parameters:

$$\hat{\theta}^1 = \begin{pmatrix} 3/7 & 1/7 \\ 1/7 & 2/7 \end{pmatrix} = \begin{pmatrix} .429 & .143 \\ .143 & .285 \end{pmatrix}$$

What is $\mathcal{L}(\hat{\theta}^1)$?

- (b) Jan thinks we should let H be the ‘best’ value it could have, that is to make the log likelihood as large as possible, and so tries setting $H = 0$ and then $H = 1$ and computes the log likelihood of the complete data in both cases. What value gives the highest complete-data log likelihood? What is the likelihood value?

- (c) Evelyn thinks this is all unprincipled messing around and says we should optimize the thing we want to optimize! That is,

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta) .$$

Evelyn also thinks we can just use the code for gradient descent that we already built in 6.867 to do this job.

Is Evelyn right?

- (d) Ariel was paying close attention in lecture and thinks this problem is an example of estimation in the presence of a latent variable and that we should use EM.

Let’s start with the guess

$$\theta_0 = \begin{pmatrix} .25 & .25 \\ .25 & .25 \end{pmatrix}$$

What is the formula for the E step in this problem? What is the numerical result in this particular case?

- (e) Ariel’s roommate Angel joins in the EM game and computes the M step, to get θ_1 . What is the numerical value in this case, and why?
- (f) Will EM always find a solution that maximizes \mathcal{L} ?

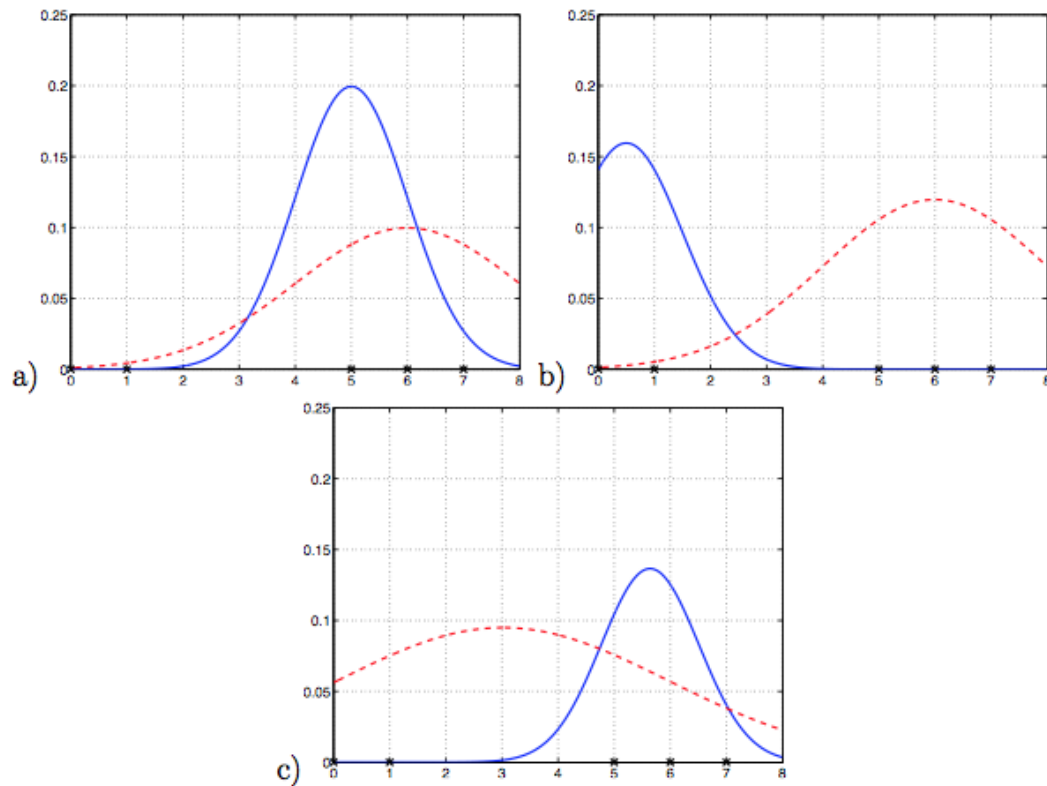
5 Mixed-up mixture

Here we are estimating a mixture of two Gaussians via the EM algorithm. The mixture distribution over x is given by

$$P(x; \theta) = P(1)N(x; \mu_1, \sigma_1^2) + P(2)N(x; \mu_2, \sigma_2^2)$$

Any student in this class could solve this estimation problem easily. Well, one student, devious as they were, scrambled the order of figures illustrating EM updates. They may have also slipped in a figure that does not belong. Your task is to extract the figures of successive updates and explain why your ordering makes sense from the point of view of how the EM algorithm works. All the figures plot $P(1)N(x; \mu_1, \sigma_1^2)$ as a function of x with a solid line and $P(2)N(x; \mu_2, \sigma_2^2)$ with a dashed line.

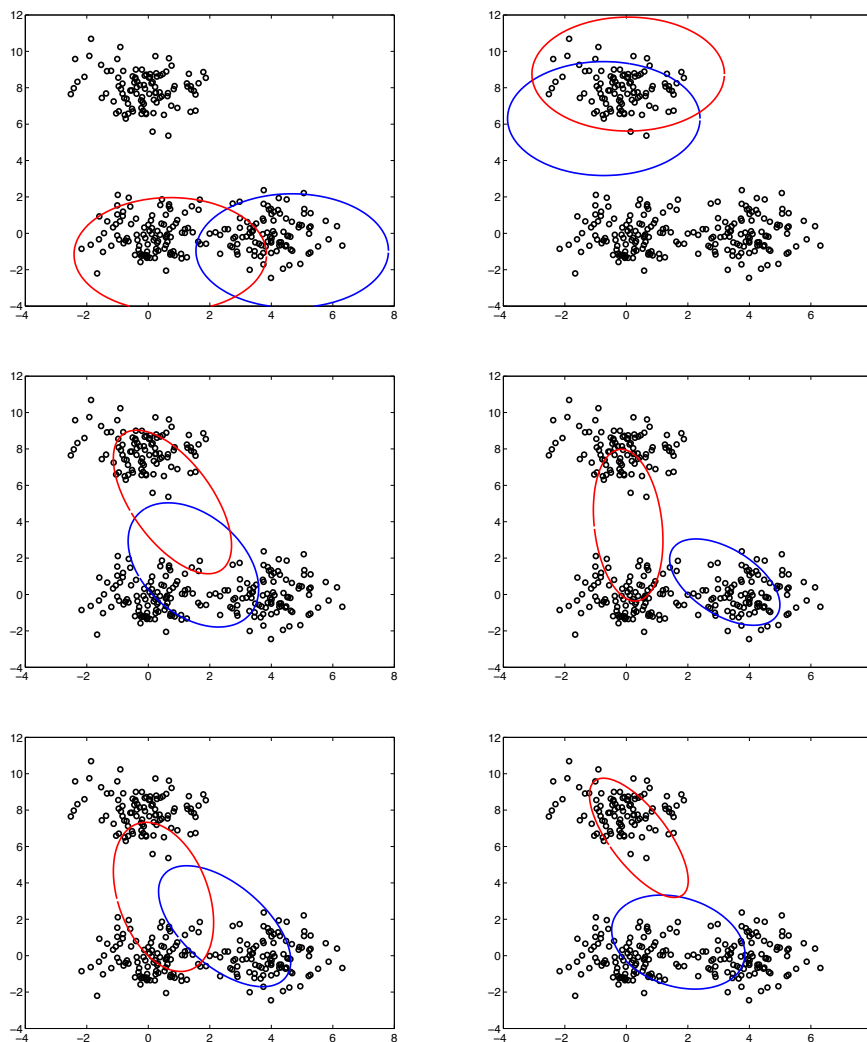
- (a) (True/False) In the mixture model, we can identify the most likely T posterior assignment, i.e., j that maximizes $P(j | x)$, by comparing the values of $P(1)N(x; \mu_1, \sigma_1^2)$ and $P(1)N(x; \mu_2, \sigma_2^2)$



- (b) Assign two figures to the correct steps in the EM algorithm.
- Step 0: () initial mixture distribution
 - Step 1: () after one EM-iteration
- (c) Briefly explain how the mixture you chose for “step 1” follows from the mixture you have in “step 0”.

6 More mixture

We estimated a two Gaussians mixture model based on two-dimensional data shown in the figure below. The mixture was initialized randomly in two different ways and run for three iterations based on each initialization. However, the figures got mixed up (yes, again!). Please draw an arrow from one figure to another to indicate how they follow from each other (you should draw only four arrows).

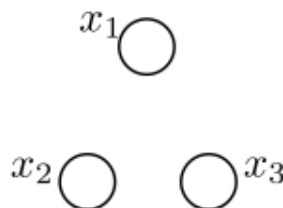


7 NB: Nota bene

Consider a Naive-Bayes model over three binary (0/1) features x_1 , x_2 , and x_3 and two classes $y = 0, 1$. By marginalizing over the class variable y , we obtain a two component mixture distribution

$$P(x_1, x_2, x_3; \theta) = \sum_{y=0,1} q(y) \prod_{j=1}^3 q_j(x_j | y)$$

- (a) If we wanted to represent this mixture distribution as a graphical model over x_1 , x_2 , and x_3 , what would the graph look like? Draw the necessary edges or arrows in the figure below.



- (b) We will use the EM algorithm to learn the parameters of the Naive Bayes model from feature observations (x_1, x_2, x_3) . We initialized the parameters as follows:

$q(y = 1) = 1/2$, the class conditional feature distributions are randomized but with the constraint that $q_j(x_j = 1 \mid y = 0) = q_j(x_j = 1 \mid y = 1)$, for all j .

What can you say about $P(y = 1 \mid x_1, x_2, x_3; \theta)$ as a function of (x_1, x_2, x_3) given our initialization?

- (c) Let $\hat{n}_j(x_j)$ denote the number of times the j^{th} feature has value x_j in the data. Given our initialization above, what are the expected counts $\hat{n}_j(x_j, y)$ we would evaluate in the first step of the EM-algorithm?
- (d) Check all the statements below that are consistent with running the EM algorithm with our initialization.
- ☐ The parameters $q(y)$ would remain as initialized.
 - ☐ None of the model parameters would change after the first M-step.
 - ☐ The property $q_j(x_j = 1 \mid y = 0) = q_j(x_j = 1 \mid y = 1)$, for all j , holds throughout the EM iterations

8 Parameterized CPT

Consider the Bayesian network $X \rightarrow Y$.

Instead of the usual CPT's we have decided to use the following parametric model with two parameters for this problem:

$$\Pr(X = 1) = \alpha$$

$$\Pr(Y = 0 \mid X = 0) = \theta$$

$$\Pr(Y = 1 \mid X = 0) = 1 - \theta$$

$$\Pr(Y = 0 \mid X = 1) = 1 - \theta$$

$$\Pr(Y = 1 \mid X = 1) = \theta$$

Given the following data \mathcal{D} :

X	Y	#occ
0	0	3
0	1	5
1	0	6
1	1	2

- Give an expression for $\Pr(\mathcal{D}|\alpha, \theta)$.
- For what value of α is it maximized?
- For what value of θ is it maximized?
- Give a numerical example of a probability distribution that can't be well modeled with this parameterization.

9 What's up, doc?

Consider the simple topic model in which:

- Words are drawn from a vocabulary \mathcal{W} ;
- Each document d is a sequence of words $w_1^d, \dots, w_{n_d}^d$;
- The distribution of words depends on the topic; for a particular topic $z = \{1 \dots k\}$, the word distribution is a multinomial: $\Pr(w | z, \theta) = \theta_{w|z}$, $\sum_{w \in \mathcal{W}} \theta_{w|z} = 1$ for all z ;
- The distribution of topics is a multinomial: $\Pr(z | \theta) = \theta_z$, $\sum_{z=1}^k \theta_z = 1$;
- The document can be viewed as being drawn from a mixture

$$\Pr(d | \theta) = \sum_{z=1}^k \Pr(z | \theta) \Pr(d | z, \theta) = \sum_{z=1}^k \theta_z \prod_{w \in d} \theta_{w|z} ,$$

in which there is a single topic per document.

- Draw the graphical model corresponding to this model. Here as it is inconvenient to write down all z^j and w_i^j in a graphical model, we therefore introduce a notation called *plate*. For more reference about *plate*, please look at Section 8.1.1 in Bishop book.
- What is the complete log likelihood, $\Pr(D, z | \theta)$, for a set of documents $D = \{d^1 \dots d^n\}$?
- Describe the E and M steps for an EM algorithm that estimates the θ values given D . Give these expressions explicitly.