

Problem Set 7

*Handed Out: November 18th, 2014**Due: December 3rd, 2014*

- Feel free to talk to other members of the class in doing the homework. I am more concerned that you learn how to solve the problem than that you demonstrate that you solved it entirely on your own. You should, however, write down your solution yourself. Please try to keep the solution brief and clear.
- Please use Piazza first if you have questions about the homework. Also feel free to send us e-mails and come to office hours.
- Please, no handwritten solutions. You will submit your solution manuscript as a single pdf file.
- The homework is due at **11:59 PM** on the due date. We will be using Compass for collecting the homework assignments. Please submit an electronic copy via Compass2g (<http://compass2g.illinois.edu>). Please do NOT hand in a hard copy of your write-up. Contact the TAs if you are face technical difficulties in submitting the assignment.
- **You cannot use the late submission credit hours for this problem set.**
- No code is needed for any of these problems. You can do the calculations however you please. You need to turn in only the report. Please name your report as `<NetID>-hw7.pdf`.

1. [EM Algorithm - 70 points]

Given a collection of documents $\{d_1, d_2, \dots, d_M\}$ where each document consists of words from vocabulary $\{w_1, w_2, \dots, w_V\}$, we want to cluster this collection into two categories, c_1 and c_2 . In this model, words and documents are observed variables and the category assignment of the words are latent variables. Each category has a multinomial distribution over the vocabulary and each document has a binomial distribution over the categories. Let us first introduce the model parameters:

- $P(d_i)$ is the probability of observing a particular document d_i .
- $P(c_k|d_i)$ is the probability that the document d_i has category c_k .
- $P(w_j|c_k)$ is the probability that word w_j appears in the category c_k .

Using these definitions, we can think about the generative process that resulted in the observed collection of documents as follows:

1. Pick a document d_i with probability $P(d_i)$. We generate M documents.
2. For each word position in the document d_i , pick a category c_k with probability $P(c_k|d_i)$.
3. For each word position, we generate a word from $\{w_1, \dots, w_V\}$ based on the category assignment c_k . That is, we generate a word w_j according to $P(w_j|c_k)$.

One way to estimate the parameters of this model is to use the EM algorithm. We are going to guide you through the steps of the EM algorithm for this model. Please use the notations defined above to answer following questions.

- (a) [10 points] What is the probability of observing a word in a document (in terms of the variables introduced above), $P(w_j, d_i)$?
- (b) [10 points] In the E-step, we estimate the posterior distribution of the latent variables given the current parameters. Derive $P(c_k|w_j, d_i)$.
- (c) [15 points] In the M-step, we maximize the expected complete data log-likelihood $E[LL]$ of the entire collection of documents. Derive $E[LL]$. (Please use $n(d_i, w_j)$ to denote the number of occurrences of w_j in document d_i . Note that it's possible that $n(d_i, w_j) = 0$ for some i and j .)
- (d) [20 points] Solve the optimization problem you formulated in (c) to derive the update rules for $P(d_i)$, $P(c_k|d_i)$ and $P(w_j|c_k)$.
- (e) [15 points] Examine the update rules and explain them in English. Also, describe in pseudocode how would you run the algorithm: initialization, iteration, and termination. What equations in the previous answers would you use at which steps in the algorithm?

2. [Tree Dependent Distributions - 30 points]

Note: In this problem, we will be looking at tree dependent distributions that will be covered in class soon. You may go through the lecture notes or wait for it to be taught in class before you attempt this problem. A brief introduction is given below.

A tree dependent distribution is a probability distribution over n variables, $\{x_1, \dots, x_n\}$ that can be represented as a tree built over n nodes corresponding to the variables. If there is a directed edge from variable x_i to variable x_j , then x_i is said to be the parent of x_j . Each directed edge $\langle x_i, x_j \rangle$ has a weight that indicates the conditional probability $\Pr(x_j | x_i)$. In addition, we also have probability $\Pr(x_r)$ associated with the root node x_r . While computing joint probabilities over tree-dependent distributions, we assume that a node is independent of all its non-descendants given its parent. For instance, in our example above, x_j is independent of all its non-descendants given x_i .

To learn a tree-dependent distribution, we need to learn three things: the structure of the tree, the probabilities on the edges of the tree, and the probabilities on the nodes. Assume that you have an algorithm to learn an *undirected* tree T with all required probabilities. To clarify, for all *undirected* edges $\langle x_i, x_j \rangle$, we have learned both probabilities, $\Pr(x_i | x_j)$ and $\Pr(x_j | x_i)$. (There exists such an algorithm and we will be covering that in class.) The only aspect missing is the directionality of edges to convert this undirected tree to a directed one.

However, it is okay to not learn the directionality of edges explicitly. In this problem, you would show that choosing any arbitrary node as the root and directing all edges away from it is sufficient, and that two directed trees obtained this way from the same underlying undirected tree T are equivalent.

- (a) [10 points] State exactly what is meant by the statement: “*The two directed trees obtained from T are equivalent.*”

- (b) **[20 points]** Show that no matter which node in T is chosen as the root for the “direction” stage, the resulting directed trees are all equivalent (based on your definition above).