

## Homework 5: EM with Mixtures, PCA, and Graphical Models

This homework assignment will have you work with EM for mixtures, PCA, and graphical models. We encourage you to read sections 9.4 and 8.2.5 of the course textbook.

Please type your solutions after the corresponding problems using this L<sup>A</sup>T<sub>E</sub>X template, and start each problem on a new page.

Please submit the **writup PDF to the Gradescope assignment ‘HW5’**. Remember to assign pages for each question.

Please submit your **L<sup>A</sup>T<sub>E</sub>X file and code files to the Gradescope assignment ‘HW5 - Supplemental’**.

**Problem 1** (Expectation-Maximization for Categorical-Geometric Mixture Models, 25pts)

In this problem we will explore expectation-maximization for a Categorical-Geometric Mixture model.

Specifically, we assume that there are a set of  $K$  parameters  $p_k \in [0, 1]$ . To generate each observation  $n$ , we first choose which of those  $K$  parameters we will use based on the (unknown) overall mixing proportion over the components  $\theta \in [0, 1]^K$ , where  $\sum_{k=1}^K \theta_k = 1$ . Let the (latent)  $\mathbf{z}_n$  indicate which of the  $K$  components we use to generate observation  $n$ . Next we sample the observation  $\mathbf{x}_n$  from a geometric distribution with parameter  $p_{\mathbf{z}_n}$ . This process can be written as:

$$\begin{aligned}\mathbf{z}_n &\sim \text{Categorical}(\theta) \\ \mathbf{x}_n &\sim \text{Geometric}(p_{\mathbf{z}_n})\end{aligned}$$

We encode observation  $n$ 's latent component-assignment  $\mathbf{z}_n \in \{0, 1\}^K$  as a one-hot vector. Element indicator variables  $z_{nk}$  equal 1 if  $\mathbf{z}_n$  was generated using component  $k$ .

A geometric distribution corresponds to the number of trials needed to get to the first success, if success occurs with probability  $p$ . Its PMF is given by  $p(x_n|p_k) = (1 - p_k)^{x_n-1}p_k$

1. **Intractability of the Data Likelihood** We are generally interested in finding a set of parameters  $p_k$  that maximize the data likelihood  $\log p(\{\mathbf{x}_n\}_{n=1}^N | \{p_k\}_{k=1}^K)$ . Expand the data likelihood to include the necessary sums over observations  $\mathbf{x}_n$  and to marginalize out (via more sums) the latents  $\mathbf{z}_n$ . Why is optimizing this likelihood directly intractable?
2. **Complete-Data Log Likelihood** The complete data  $D = \{(\mathbf{x}_n, \mathbf{z}_n)\}_{n=1}^N$  includes latents  $\mathbf{z}_n$ . Write out the complete-data negative log likelihood. Apply the “power trick”<sup>a</sup> and simplify your expression using indicator elements  $z_{nk}$ .

$$\mathcal{L}(\theta, \{p_k\}_{k=1}^K) = -\ln p(D | \theta, \{p_k\}_{k=1}^K).$$

Note that optimizing this loss is now computationally tractable if we know  $\mathbf{z}_n$ .

3. **Expectation Step** Our next step is to introduce a mathematical expression for  $\mathbf{q}_n$ , the posterior over the hidden topic variables  $\mathbf{z}_n$  conditioned on the observed data  $\mathbf{x}_n$  with fixed parameters, i.e  $\mathbf{q}_n = p(\mathbf{z}_n | \mathbf{x}_n; \theta, \{p_k\}_{k=1}^K)$ .
  - **Part 3.A** Write down and simplify the expression for  $\mathbf{q}_n$ . Note that because the  $\mathbf{q}_n$  represents the posterior over the hidden categorical variables  $\mathbf{z}_n$ , the components of vector  $\mathbf{q}_n$  must sum to 1.
  - **Part 3.B** Give an algorithm for calculating the expression for  $\mathbf{q}_n$  found in Part 3.A for all  $n$ , given the observed data  $\{\mathbf{x}_n\}_{n=1}^N$  and settings of the parameters  $\theta$  and  $\{p_k\}_{k=1}^K$ .
4. **Maximization Step** Using the  $\mathbf{q}_n$  estimates from the Expectation Step, derive an update for maximizing the expected complete data log likelihood in terms of  $\theta$  and  $\{p_k\}_{k=1}^K$ .
  - **Part 4.A** Derive an expression for the expected complete-data log likelihood using  $\mathbf{q}_n$ .
  - **Part 4.B** Find an expression for  $\theta$  that maximizes this expected complete-data log likelihood. You may find it helpful to use Lagrange multipliers in order to enforce the constraint  $\sum \theta_k = 1$ . Why does this optimal  $\theta$  make intuitive sense?
  - **Part 4.C** Find an expression for the  $\{p_k\}_{k=1}^K$  that maximize the expected complete-data log likelihood. Why does this optimal  $\{p_k\}_{k=1}^K$  make intuitive sense?
5. Suppose that this had been a classification problem. That is, you were provided the “true” categories  $\mathbf{z}_n$  for each observation  $\mathbf{x}_n$ , and you were going to perform the classification by inverting the provided generative model (i.e. now you’re predicting  $z$  given  $x$ ). Could you reuse any of your inference derivations above?

(Continued on next page.)

<sup>a</sup>The “power trick” is used when terms in a PDF are raised to the power of indicator components of a one-hot vector. For example, it allows us to rewrite  $p(\mathbf{z}_n; \theta) = \theta_k^{z_{nk}}$ .

**Problem 1** (cont.)

6. Finally, implement your solution (see `T5_P1.py` for starter code). You are responsible for implementing the `expected_loglikelihood`, `e_step` and `m_step` functions. Test it out with data given 10 samples from 3 components with  $p_1 = .1$ ,  $p_2 = .5$ , and  $p_3 = .9$ . How does it perform? What if you increase the number of samples to 1000 from each of the components? What if you change  $p_2 = .2$ ? Hypothesize reasons for the differences in performance when you make these changes. You may need to record five to ten trials (random restarts) in order to observe meaningful insights.

**Solution**

**Problem 2** (PCA, 15 pts)

For this problem you will implement PCA from scratch. Using `numpy` to call SVDs is fine, but don't use a third-party machine learning implementation like `scikit-learn`.

We return to the MNIST data set from T4. You have been given representations of 6000 MNIST images, each of which are  $28 \times 28$  greyscale handwritten digits. Your job is to apply PCA on MNIST, and discuss what kind of structure is found.

The given code in `T5_P2.py` loads the images into your environment. File `T5_P2_Autograder.py` contains a test case to check your cumulative proportion of variance.

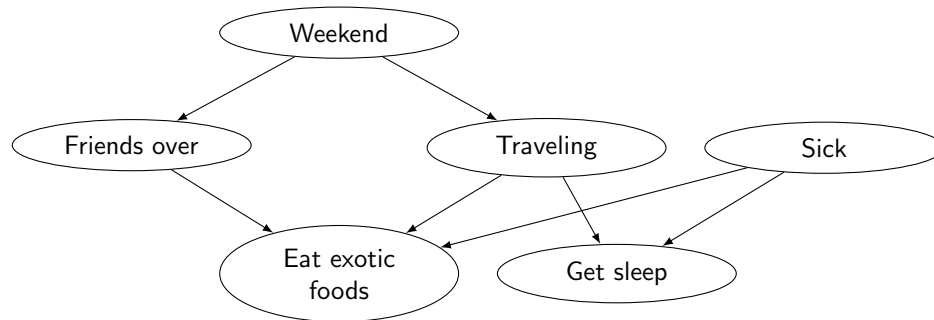
1. Compute the PCA. Plot the eigenvalues corresponding to the most significant 500 components in order from most significant to least. Make another plot that describes the cumulative proportion of variance explained by the first  $k$  most significant components for values of  $k$  from 1 through 500. How much variance is explained by the first 500 components? Describe how the cumulative proportion of variance explained changes with  $k$ .
2. Plot the mean image of the dataset and plot an image corresponding to each of the first 10 principle components. How do the principle component images compare to the cluster centers from K-means? Discuss any similarities and differences. Include all 11 plots in your PDF submission.
3. Compute the reconstruction error on the data set using the mean image of the dataset. Then compute the reconstruction error using the first 10 principal components. How do these errors compare to the final objective loss achieved by using K-means on the dataset? Discuss any similarities and differences.

*Include your plots in your PDF. There may be several plots for this problem, so feel free to take up multiple pages.*

**Solution**

**Problem 3** (Bayesian Networks, 10 pts)

In this problem we explore the conditional independence properties of a Bayesian Network. Consider the following Bayesian network representing a fictitious person's activities. Each random variable is binary (true/false).



The random variables are:

- **Weekend:** Is it the weekend?
- **Friends over:** Does the person have friends over?
- **Traveling:** Is the person traveling?
- **Sick:** Is the person sick?
- **Eat exotic foods:** Is the person eating exotic foods?
- **Get Sleep:** Is the person getting sleep?

For the following questions,  $A \perp B$  means that events A and B are independent and  $A \perp B|C$  means that events A and B are independent conditioned on C.

**Use the concept of d-separation** to answer the questions and show your work (i.e., state what the blocking path(s) is/are and what nodes block the path; or explain why each path is not blocked).

*Example Question:* Is Friends over  $\perp$  Traveling? If NO, give intuition for why.

*Example Answer:* NO. The path from Friends over – Weekend – Traveling is not blocked following the d-separation rules. Thus, the two are not independent. Intuitively, this makes sense as if say we knew that the person was traveling, it would make it more likely to be the weekend. This would then make it more likely for the person to have friends over.

**Actual Questions:**

1. Is Sick  $\perp$  Weekend? If NO, give intuition for why.
2. Is Sick  $\perp$  Friends over | Eat exotic foods? If NO, give intuition for why.
3. Is Friends over  $\perp$  Get Sleep? If NO, give intuition for why.
4. Is Friends over  $\perp$  Get Sleep | Traveling? If NO, give intuition for why.
5. Suppose the person stops traveling in ways that affect their sleep patterns (as various famous people have done). Travel still affects whether they eat exotic foods. Draw the modified network. (Feel free to reference the handout file for the commands for displaying the new network in L<sup>A</sup>T<sub>E</sub>X).
6. For this modified network, is Friends over  $\perp$  Get Sleep? If NO, give an intuition why. If YES, describe what observations (if any) would cause them to no longer be independent.

## Solution

**Name**

**Collaborators and Resources**

Whom did you work with, and did you use any resources beyond cs181-textbook and your notes?

**Calibration**

Approximately how long did this homework take you to complete (in hours)?