

## Mixtures, EM, and Graphical Models

### Introduction

This homework assignment will have you work with EM for mixtures, PCA, and graphical models.

### Resources and Submission Instructions

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  $\text{\LaTeX}$  template, and start each problem on a new page.

Please submit the writeup PDF to the Gradescope assignment ‘HW5’. Remember to assign pages for each question. **You must include any plots in your writeup PDF.** . Please submit your  $\text{\LaTeX}$ file and code files to the Gradescope assignment ‘HW5 - Supplemental.’ The supplemental files will only be checked in special cases, e.g. honor code issues, etc. Your files should be named in the same way as we provide them in the repository, e.g. `hw5.pdf`, etc.

**Problem 1** (Expectation-Maximization for Gamma Mixture Models: Derivations, 10pts)

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

Let us suppose the following generative story for an observation  $x$ : first one of  $K$  classes is randomly selected, and then the features  $x$  are sampled according to this class. If

$$z \sim \text{Categorical}(\theta)$$

indicates the selected class, then  $x$  is sampled according to the class or “component” distribution corresponding to  $z$ . (Here,  $\theta$  is the mixing proportion over the  $K$  components:  $\sum_k \theta_k = 1$  and  $\theta_k > 0$ .) In this problem, we assume these component distributions are gamma distributions with shared shape parameter but different rate parameters:

$$x|z \sim \text{Gamma}(\alpha, \beta_k).$$

In an unsupervised setting, we are only given a set of observables as our training dataset:  $\mathcal{D} = \{x^{(n)}\}_{n=1}^N$ . The EM algorithm allows us to learn the underlying generative process (the parameters  $\theta$  and  $\{\beta_k\}$ ) despite not having the latent variables  $\{z^{(n)}\}$  corresponding to our training data.

1. **Intractability of the Data Likelihood.** We are generally interested in finding a set of parameters  $\beta_k$  that maximizes the likelihood of the observed data:

$$\log p(\{x^{(n)}\}_{n=1}^N; \theta, \{\beta_k\}_{k=1}^K).$$

Expand the data likelihood to include the necessary sums over observations  $x^{(n)}$  and to marginalize out the latents  $\mathbf{z}^{(n)}$ . Why is optimizing this likelihood directly intractable?

2. **Complete Data Log Likelihood.** The complete dataset  $\mathcal{D} = \{(x^{(n)}, \mathbf{z}^{(n)})\}_{n=1}^N$  includes latents  $\mathbf{z}^{(n)}$ . Write out the negative complete data log likelihood:

$$\mathcal{L}(\theta, \{\beta_k\}_{k=1}^K) = -\log p(\mathcal{D}; \theta, \{\beta_k\}_{k=1}^K).$$

Apply the power trick and simplify your expression using indicator elements  $z_k^{(n)}$ .<sup>a</sup> Notice 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 component variables  $\mathbf{z}^{(n)}$  conditioned on the observed data  $x^{(n)}$  with fixed parameters. That is:

$$\mathbf{q}^{(n)} = \begin{bmatrix} p(\mathbf{z}^{(n)} = \mathbf{C}_1 | x^{(n)}; \theta, \{\beta_k\}_{k=1}^K) \\ \vdots \\ p(\mathbf{z}^{(n)} = \mathbf{C}_K | x^{(n)}; \theta, \{\beta_k\}_{k=1}^K) \end{bmatrix}.$$

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. The main work is to find an expression for  $p(\mathbf{z}^{(n)} | x^{(n)}; \theta, \{\beta_k\}_{k=1}^K)$  for any choice of  $\mathbf{z}^{(n)}$ ; i.e., for any 1-hot encoded  $\mathbf{z}^{(n)}$ . With this, you can then construct the different components that make up the vector  $\mathbf{q}^{(n)}$ .

<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  $p(\mathbf{z}^{(n)}; \theta) = \prod_k \theta_k^{z_k^{(n)}}$ .

**Problem 1** (cont.)

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  $\{\beta_k\}_{k=1}^K$ .
- (a) Derive an expression for the expected complete data log likelihood using  $\mathbf{q}^{(n)}$ .
  - (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?
  - (c) Find an expression for  $\beta_k$  that maximizes the expected complete data log likelihood. Why does this optimal  $\beta_k$  make intuitive sense?
5. Suppose that this had been a classification problem. That is, you were provided the “true” components  $\mathbf{z}^{(n)}$  for each observation  $x^{(n)}$ , and you were going to perform the classification by inverting the provided generative model (i.e. now you’re predicting  $\mathbf{z}^{(n)}$  given  $x^{(n)}$ ). Could you reuse any of your derivations above to estimate the parameters of the model?

**Solution:** [Your solution here.](#)

**Problem 2** (Expectation-Maximization for Gamma Mixture Models: Coding, 15 pts)

In this problem, you will implement your EM derivations from Problem 1 and apply it to analyzing a synthetic example of the recovery time for patients following a surgical procedure, in hours. The doctors have noticed that some patients seem to recover at an expected rate, but sometimes the recovery takes a long time. They are keen to understand what is going on to improve their processes.

1. Plot the data. How would you describe the distribution? Based on what you see, why might a mixture model be an appropriate model?
2. Implement your solution from Problem 1 in `homework5.ipynb`. You do not need to include your code in your writeup.

Note that for numerical stability, we recommend using the log-probability directly; for example, you could use the `Gamma` class from `torch.distributions` and then use the `log_prob` and `logsumexp` methods.

3. Run your code for 1, 2, 3, and 4 mixture components. Plot the mixture models you find on top of the data distribution as well as the associated log likelihoods. How many mixtures does it seem that there are? How would you decide?
4. The doctors tell you that a normal recovery from the procedure is about 2-3 days, though sometimes patients recover a little faster. Does this match what you see in the data? Provide some hypotheses about what might be going on.
5. It's clear from the data that some patients take significantly longer than 2-3 days. Do you observe that there is evidence that these represent a different cluster, vs. a long tail from a single cluster? Why or why not?
6. The physician-scientists want to use this model to understand the characteristics of patients who have very long recoveries vs. those who do not. Is this mixture modeling approach appropriate for this task? Why or why not?
7. The physician-scientists develop a way of identifying someone's cluster based on a blood test—it seems that some patients in the longer group are ones that are at risk for clotting-related complications. The hospital operations staff want to use this model to help streamline operations. They plan to use the cluster of the patient to predict which patients will have a long length of stay. Is this plan sound? May there be some issues?

**Solution:** [Your solution here.](#)

### Problem 3 (PCA, 15 pts)

For this problem you will implement PCA from scratch on the first 6000 images of the MNIST dataset. Your job is to apply PCA on MNIST and discuss what kind of structure is found. Implement your solution in `homework5.ipynb` and attach the final plots below.

**You will receive no points for code not included below or for using third-party PCA implementations (i.e. `scikit-learn`).**

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  $k$ . Include this plot below.
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 these two plots below.

*Reminder: Center the data before performing PCA.*

3. Compute the reconstruction error on the dataset using the first 10 principal components. Then compute the reconstruction error when the reconstruction for each point is just the mean image of the dataset. How do these errors compare to the final objective loss achieved by using K-means on the dataset? Discuss any similarities and differences.

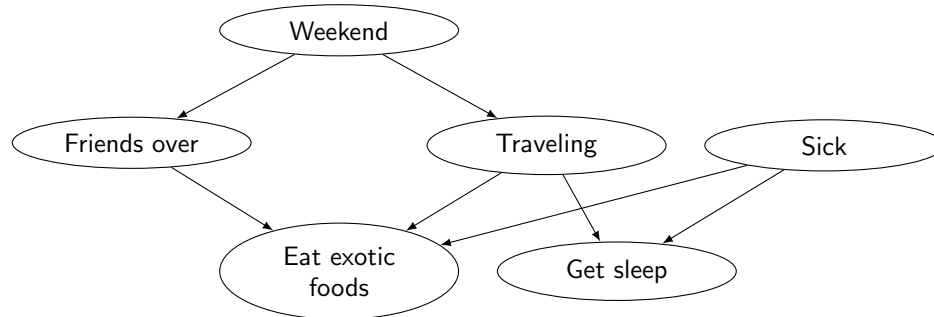
For consistency in grading, define the error function as the squared L2 norm of the difference between the true data and the reconstruction, averaged over all data points.

4. Suppose you took the original matrix of principle components that you found  $V$  and multiplied it on the right side by some rotation matrix  $R$  (i.e., you considered the matrix  $VR$ ). Would that change the quality of the reconstruction error in the last problem? The interpretation of the components? Why or why not?
5. Let's recall the zipcode application in Homework 3. A common application of PCA is to dimensionality reduction before running a classifier: You first project the data onto the first few PCA bases, and then you train a classifier from the projection to the output.
  - (a) First, how might this be advantageous to just applying the classification algorithm directly, from both a robustness and efficiency perspective?
  - (b) Second, recall from Homework 3 that adversaries can attack a classification algorithm by manipulating/perturbing the data; how could this approach help with such attacks?
6. You are collaborating with a penmanship analysis expert. They are able to identify the kind of pen used to make a mark by various characteristics such as the width of the line, its crispness, and the type (if any) of ink splatter. They have heard that your machine learning helped automate reading zip codes for the post office; they are wondering if you can help automate the manual process of classifying pen types.
  - (a) Does what the expert is describing correspond to some kind of hidden representation or latent variable? Describe why or why not.
  - (b) Do you think PCA will help the expert? Why or why not?

**Solution:** Your solution here.

**Problem 4** (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  $A \perp B \mid C$  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). For example, consider the following question and answer:

- *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 as we do not observe Weekend. Thus, the two are not independent.

**Actual Questions:**

1. Is Weekend  $\perp$  Get Sleep? If NO, give intuition for why.
2. Is Sick  $\perp$  Weekend? If NO, give intuition for why.
3. Is Sick  $\perp$  Friends over  $\mid$  Eat exotic foods? If NO, give intuition for why.
4. Is Friends over  $\perp$  Get Sleep? If NO, give intuition for why.
5. Is Friends over  $\perp$  Get Sleep  $\mid$  Traveling? If NO, give intuition for why.
6. Suppose the person stops traveling in ways that affect their sleep patterns. 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  $\text{\LaTeX}$ ).
7. 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:** Your solution here.



**Name:**  
**Collaborators and Resources:**