

# CAP 6610 Machine Learning, Spring 2022

## Homework 4

Due 4/20/2022 11:59PM

1. (15%) *Alternating optimization for PCA*. Recall that the principal component analysis (PCA) tries to find the projection matrix  $\Theta$  and embeddings  $\mathbf{y}_1, \dots, \mathbf{y}_n$  as the solution of

$$\underset{\Theta, \mathbf{y}_i}{\text{minimize}} \quad \sum_{i=1}^n \|\phi_i - \Theta \mathbf{y}_i\|^2.$$

We can define matrices  $\Phi = [\phi_1 \ \phi_2 \ \dots \ \phi_n]$  and  $\mathbf{Y} = [\mathbf{y}_1 \ \mathbf{y}_2 \ \dots \ \mathbf{y}_n]$  and rewrite the problem as

$$\underset{\Theta, \mathbf{Y}}{\text{minimize}} \quad \|\Phi - \Theta \mathbf{Y}\|^2. \quad (1)$$

Without imposing the constraint  $\Theta^\top \Theta = \mathbf{I}$ , derive an alternating optimization algorithm for (1).

2. (15%) *Naive Bayes GMM*. Consider a Gaussian mixture model in which the marginal distribution of the latent variable  $\mathbf{y}$  is  $\Pr(\mathbf{y} = \mathbf{e}_c) = \pi_c, c = 1, \dots, k$ , and the conditional distribution of  $\mathbf{x}$  given  $\mathbf{y}$  is  $\mathcal{N}(\mu_c, \Sigma_c)$  where each of the covariance matrices  $\Sigma_c$  is diagonal. This means given the latent variable  $\mathbf{y}$ , elements of the ambient variable  $\mathbf{x}$  are conditionally independent, resembling the naive Bayes model. Given a set of observations  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , derive the expectation-maximization algorithm for estimating the model parameters  $\pi_c, \mu_c$ , and  $\Sigma_c$  for  $c = 1, \dots, k$ .
3. (40%) *20 Newsgroup revisited*. Let us revisit the 20 Newsgroup data set <<http://qwone.com/~jason/20Newsgroups/>>, and apply some of the unsupervised methods by ignoring their labels. We will only consider the training data. You are required to code the algorithms by yourselves in the language of your choice.
  - (a) *LSI/PCA via orthogonal iteration*. Implement the orthogonal iteration algorithm that finds the PCA projection matrix  $\Theta$  of a data matrix  $\Phi$ . You are allowed to use a pre-existing function of QR. Apply tf-idf to the term-document matrix to obtain  $\Phi$  and feed it into your orthogonal iteration algorithm. Remember to use sparse matrix operations to avoid unnecessary memory/computational complexities. Set  $k = 2$  and let the algorithm run until  $\Theta$  doesn't change much. Then get  $\mathbf{Y} = \Theta^\top \Phi$ . Each column of  $\mathbf{Y}$  is a two-dimensional vector that you can plot on a plain. Plot all the documents on a two-dimensional plain, and use a different color for each point that belong to different news groups.
  - (b) *GMM via EM*. Implement the EM algorithm for the Gaussian mixture model (with different means and covariances for each Gaussian component). The data matrix  $\Phi$  is obtained from LSI with  $k_{\text{LSI}} = 100$  using the previous orthogonal iteration algorithm. Run the EM algorithm for GMM with  $k_{\text{GMM}} = 20$  until convergence. For each Gaussian component  $\mathcal{N}(\mu_c, \Sigma_c)$ , calculate  $\Theta \mu_c$  where  $\Theta$  is the PCA projection; the vector  $\Theta \mu_c$  should be element-wise nonnegative. For each cluster  $c$ , show the 10 terms that have the highest value in  $\Theta \mu_c$ . The index-term mapping can be found here <<http://qwone.com/~jason/20Newsgroups/vocabulary.txt>>. Does the result make sense?

4. (30%) *Mixture of multinomials.* Consider the latent variable model with the following probability distribution:

$$\Pr(y_i = c) = \frac{1}{k}, \quad \Pr(\mathbf{x}_i | y_i = c) \sim \text{Multi}(\mathbf{p}_c, L_i),$$

meaning that  $y_i$  is categorical, with  $k$  possible outcomes and equal probability;  $(\mathbf{x}_i | y_i = c)$  follows a multinomial distribution by drawing from  $\mathbf{p}_c$   $L_i$  times, i.e.,

$$p(\mathbf{x}_i | y_i = c) = \frac{L_i!}{\prod_{j=1}^d x_{ij}!} \prod_{j=1}^d p_{cj}^{x_{ij}}.$$

Given data samples  $\mathbf{x}_1, \dots, \mathbf{x}_n$ :

- (a) Write out the maximum likelihood formulation for estimating  $\mathbf{p}_1, \dots, \mathbf{p}_k$ . Simplify the objective function as much as possible.
- (b) Derive an expectation-maximization algorithm for approximately solving the aforementioned problem.
- (c) Implement this algorithm and try it on the 20 News Group data set with  $k = 20$  (on the raw word-count data, without tf-idf preprocessing). Show the top 10 words in each cluster.