

CAP 6610 Machine Learning, Spring 2022

Homework 4 Solution

1. (15%) *Alternating optimization for PCA.* Recall that the principal component analysis (PCA) tries to find the projection matrix Θ 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).

Solution. Fixing Θ , minimizing $\|\phi_i - \Theta \mathbf{y}_i\|^2$ with respect to \mathbf{y}_i is a least squares problem with closed-form solution $(\Theta^\top \Theta)^{-1} \Theta^\top \phi_i$. Stacking them together for $i = 1, \dots, n$ as columns, we get

$$\mathbf{Y} \leftarrow (\Theta^\top \Theta)^{-1} \Theta^\top \Phi. \quad (2)$$

Now if we fix \mathbf{Y} , the problem is still a least squares problem if we transpose everything, thus

$$\Theta \leftarrow \Phi \mathbf{Y}^\top (\mathbf{Y} \mathbf{Y}^\top)^{-1}. \quad (3)$$

The resulting alternating optimization algorithm alternates between (2) and (3).

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 Σ_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 π_c, μ_c , and Σ_c for $c = 1, \dots, k$.

Solution. Let's start by fixing the model parameters $\pi_c, \mu_c, \Sigma_c, c = 1, \dots, k$, and compute the conditional mean $E[\mathbf{y}_i | \mathbf{x}_i]$

$$E[y_{ic} | \mathbf{x}_i] = \frac{\pi_c \det(2\pi \Sigma_c)^{-1/2} \exp(-\frac{1}{2}(\mathbf{x}_i - \mu_c)^\top \Sigma_c^{-1}(\mathbf{x}_i - \mu_c))}{\sum_j \pi_j \det(2\pi \Sigma_j)^{-1/2} \exp(-\frac{1}{2}(\mathbf{x}_i - \mu_j)^\top \Sigma_j^{-1}(\mathbf{x}_i - \mu_j))}. \quad (4)$$

Next we derive how to estimate the parameters by replacing \mathbf{y}_i with ψ_i in $p(\mathbf{x}_i, \mathbf{y}_i; \theta)$:

$$\sum_{i=1}^n E[-\log p(\mathbf{x}_i, \mathbf{y}_i)] = \sum_{i=1}^n \sum_{c=1}^k \psi_{ic} \left(\frac{1}{2}(\mathbf{x}_i - \mu_c)^\top \Sigma_c^{-1}(\mathbf{x}_i - \mu_c) + \frac{1}{2} \log \det(2\pi \Sigma_c) - \log(\pi_c) \right).$$

By taking the gradient equal to zero, we get

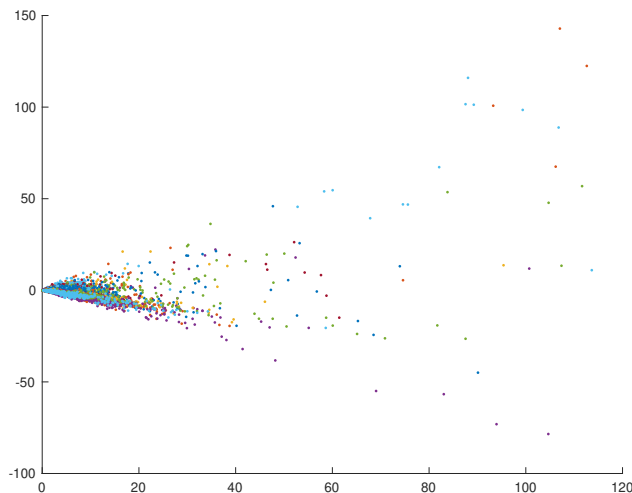
$$\begin{aligned} \pi_c &= \frac{\sum_{i=1}^n \psi_{ic}}{n}, & \mu_c &= \frac{1}{\sum_{i=1}^n \psi_{ic}} \sum_{i=1}^n \psi_{ic} \mathbf{x}_i, \\ \Sigma_c &= \frac{1}{n} \text{Diag} \left(\sum_{i=1}^n \psi_{ic} (\mathbf{x}_i - \mu_c)(\mathbf{x}_i - \mu_c)^\top \right), & c &= 1, \dots, k. \end{aligned} \quad (5)$$

The resulting EM algorithm alternates between (4) and (5). We see that the only difference between this and the general EM algorithm for GMM is that each Σ_c matrix takes the diagonal of the empirical sample covariance.

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 Θ of a data matrix Φ . You are allowed to use a pre-existing function of QR. Apply tf-idf to the term-document matrix to obtain Φ 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 Θ doesn't change much. Then get $Y = \Theta^T \Phi$. Each column of 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 Φ 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 Θ 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?

Solution. A sample MATLAB code is given.

- (a) The figure is given below. It looks like PCA embedding does not reveal much clustering structure from data.



- (b) The top-10 keywords for all 20 clusters seem to be the same, even though a k -means++ initialization is applied. This shows the limitation of EM and the difficulty of estimating the parameters of GMM.

Here's the top-10 words appearing in all clusters. They are all commonly used words and thus not very informative:

- as
- they
- was
- you
- we
- he
- are
- by
- will
- not

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.

Solution.

- (a) The likelihood of \mathbf{x}_i is

$$\begin{aligned} p(\mathbf{x}_i) &= \sum_{c=1}^k p(\mathbf{x}_i | y_i = c) \Pr(Y_i = c) \\ &= \sum_{c=1}^k \frac{L_i!}{\prod_{j=1}^m x_{ij}!} \prod_{j=1}^m p_{cj}^{x_{ij}} \times \frac{1}{k} \\ &= \frac{L_i!}{k \prod_{j=1}^m x_{ij}!} \sum_{c=1}^k \prod_{j=1}^m p_{cj}^{x_{ij}}. \end{aligned}$$

Its log-likelihood is

$$\log p(\mathbf{x}_i) = \log \sum_{c=1}^k \prod_{j=1}^m p_{cj}^{x_{ij}} + C_i,$$

where $C_i = \log \frac{L_i!}{k \prod_{j=1}^m x_{ij}!}$ is a constant that does not depend on the unknown parameter \mathbf{p}_c 's. Assuming the data samples are independent, the ML objective is

$$\begin{aligned} & \underset{\mathbf{p}_1, \dots, \mathbf{p}_k}{\text{maximize}} && \sum_{i=1}^n \log \sum_{c=1}^k \prod_{j=1}^m p_{cj}^{x_{ij}} \\ & \text{subject to} && \mathbf{p}_c \in \Delta^m, \quad c = 1, \dots, k, \end{aligned} \quad (6)$$

where Δ^d is the probability simplex defined as

$$\Delta^m = \{\mathbf{x} \in \mathbb{R}^m : \mathbf{x} \geq 0, \mathbf{1}^\top \mathbf{x} = 1\}.$$

(b) First we fix the model parameters $\mathbf{p}_c, c = 1, \dots, k$ and compute the conditional mean $\mathbb{E}[\mathbf{y}_i | \mathbf{x}_i]$

$$\begin{aligned} \psi_{ic} = \mathbb{E}[y_{ic} | \mathbf{x}_i] &= \frac{\frac{L_i!}{k \prod_{j=1}^m x_{ij}!} \prod_{j=1}^m p_{cj}^{x_{ij}}}{\frac{L_i!}{k \prod_{j=1}^m x_{ij}!} \sum_{s=1}^k \prod_{j=1}^m p_{sj}^{x_{ij}}} \\ &= \frac{\prod_{j=1}^m p_{cj}^{x_{ij}}}{\sum_{s=1}^k \prod_{j=1}^m p_{sj}^{x_{ij}}}. \end{aligned} \quad (7)$$

Next we maximize $\mathbb{E}[-\log p(\mathbf{y}, \mathbf{x})]$, which in this case is

$$\sum_{i=1}^n \sum_{c=1}^k \psi_{ic} \sum_{j=1}^m x_{ij} \log(p_{cj}).$$

To maximize it, we get

$$p_{cj} = \frac{\sum_{i=1}^n \psi_{ic} x_{ij}}{\sum_{i=1}^n \psi_{ic} L_i} \quad (8)$$

The EM algorithm for this model alternates between (7) and (8).

(c) An implementation is given. A practical issue is when calculating ψ_i 's using (7), it takes the products of many small numbers, which may cause underflows.

To calculate the numerator of (7), we can use the formula

$$\nu_{ic} = \exp \left(\sum_{j=1}^m x_{ij} \log(p_{cj}) \right).$$

Then

$$\psi_{ic} = \frac{\nu_{ic}}{\sum_{s=1}^k \nu_{is}}.$$

The result changes from different initializations, as expected when trying to minimize a nonconvex objective. Nevertheless, the solution seems to make more sense compared to the GMM model. To get more meaningful results, you might need to do some engineering first, for example by removing the 'stop words' (words that show up more often in all cases but do not carry much information, like 'that' or 'just').

topic 1	topic 2	topic 3	topic 4	topic 5	topic 6	topic 7	topic 8	topic 9	topic 10
more	am	hello	just	shipping	hello	subscribe	good	it	just
say	satan	testing	kidding	box	testing	quit	luck	only	kidding
need	just	just	good	manual	just	graphics	just	message	good
just	good	good	am	dennis	good	comp	am	test	am
good	luck	am	luck	included	am	just	satan	is	luck
am	kidding	luck	satan	just	luck	good	kidding	this	satan
luck	included	satan	included	good	satan	am	included	just	included
satan	shipping	kidding	shipping	am	kidding	luck	shipping	good	shipping
kidding	hello	included	hello	luck	included	satan	hello	am	hello
included	testing	shipping	testing	satan	shipping	kidding	testing	luck	testing
topic 11	topic 12	topic 13	topic 14	topic 15	topic 16	topic 17	topic 18	topic 19	topic 20
am	good	is	more	shipping	included	is	is	hello	included
satan	luck	this	say	box	shipping	this	this	testing	shipping
just	just	tesrt	need	manual	just	test	test	just	just
good	am	test	just	dennis	good	thanks	erme	good	good
luck	satan	thanks	good	included	am	erme	thanks	am	am
kidding	kidding	erme	am	just	luck	tesrt	tesrt	luck	luck
included	included	just	luck	good	satan	just	just	satan	satan
shipping	shipping	good	satan	am	kidding	good	good	kidding	kidding
hello	hello	am	kidding	luck	hello	am	am	included	hello
testing	testing	luck	included	satan	testing	luck	luck	shipping	testing