

# EECS545 Machine Learning

## Homework #5

revision: 1

2023. 3.

Last Update: 4/1 23:30PM (rev1)

## Contents

<b>1 [25 points] K-means and GMM for image compression</b>	<b>2</b>
1.1 K-means Clustering . . . . .	2
1.2 Gaussian Mixtures . . . . .	4
<b>2 [20 points] Expectation Maximization for GDA with missing labels</b>	<b>6</b>
2.1 Lower bound [3 points] . . . . .	7
2.2 E-step [2 points] . . . . .	7
2.3 M-step for $\mu_k$ [6 points] . . . . .	7
2.4 M-step for $\phi$ [6 points] . . . . .	9
2.5 M-step for $\Sigma_k$ [3 points] . . . . .	10
<b>3 [20 points] PCA and eigenfaces</b>	<b>11</b>
<b>4 [10 points] Independent Component Analysis</b>	<b>15</b>
<b>5 [25 points] Conditional Variational Autoencoders</b>	<b>16</b>

# 1 [25 points] K-means and GMM for image compression

In this problem, we will apply the K-means algorithm and Gaussian Mixture Models (GMM) to lossy image compression, by reducing the number of colors used in an image.

You are given two image files `mandrill-large.tiff` and `mandrill-small.tiff` of different resolution. Here, `mandrill-large.tiff` is a  $512 \times 512$  image of a mandrill represented in 24-bit color. This means that, for each of the  $512 \times 512 = 262,144$  pixels in the image, there are three 8-bit numbers (each ranging from 0 to 255) that represent the red, green, and blue intensity values for that pixel. The file `mandrill-small.tiff` contains a  $128 \times 128$  version of `mandrill-large.tiff`. The straightforward representation of this image therefore takes about  $262,144 \times 3 = 786,432$  bytes (a byte being 8 bits).

## 1.1 K-means Clustering

To compress the image, we will use K-means to reduce the image to  $K = 16$  colors. More specifically, each pixel in the image is considered as a point in the three-dimensional  $(r, g, b)$ -space:  $[0, 255]^3$ . To compress the image, we will cluster these points in color-space into 16 clusters, and replace each pixel with the closest cluster centroid.

We use the Euclidean distance  $\|x - y\|_2$  to measure the distance between two points  $x$  and  $y$ . In the pixel (RGB) space, the distance (or the *pixel error*) between  $(r, g, b)$  and  $(r', g', b')$  can be written as:  $\sqrt{(r - r')^2 + (g - g')^2 + (b - b')^2}$ .

Now, follow the instructions below.

- (a) (6 pts) **(Autograder)** Start working on `kmeans_gmm.ipynb` and `kmeans.py`, and submit them when you are done implementing.

Treating each pixel's  $(r, g, b)$  values as an element of  $\mathbb{R}^3$ , implement and run K-means with 16 clusters on the pixel data from `mandrill-small.tiff` image, running 50 update steps. We provided the initial centroids as `initial_centroids` in the starter code to ensure a deterministic result.

You will need to implement (a general version of) K-means algorithm in `kmeans.train_kmeans()`, which will be graded with the provided sample data and some random data. In order to get full points, your implementation should be efficient and fast enough (otherwise you will get only partial points).

Hint: You may use `sklearn.metrics.pairwise_distances` function to compute the distance between centroids and data points, although it would be not difficult to implement this function (in a vectorized version) on your own.

Solution.

Example output (we do not ask submitting nor grade this, as grading will be done via Autograder):

```
Iteration  0: mean error = 27.78
Iteration  1: mean error = 22.08
Iteration  2: mean error = 21.43
...
Iteration 47: mean error = 19.49
Iteration 48: mean error = 19.49
Iteration 49: mean error = 19.49
```

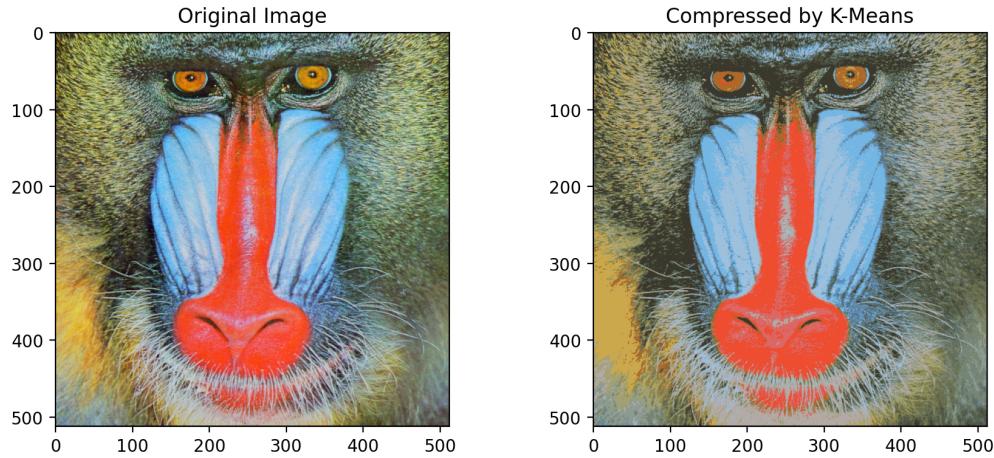
It will converge around to mean pixel error of 19.50 by around Iteration 30.

- (b) (3 pts) After training on the train image `mandrill-small.tiff`, read the test image `mandrill-large.tiff`, and replace each pixel's  $(r, g, b)$  values with the value of the closest cluster centroid.

Use the notebook's plotting code to display the original and compressed images side-by-side, and attach the plots to the **write-up**. (Note: you should have reasonable image quality/resolution to make the difference discernable). Also, measure and write down the mean pixel error between the original and compressed image.

Solution.

Pixel Error = 14.64576



- (c) (1 pts) If we represent the image with these reduced 16 colors, by (approximately) what factor have we compressed the image (in terms of bits used to represent pixels) in terms of the data size? Include an explanation of why.

Solution.

**Answer:**

The original image uses 24 bits to represent each pixel. The compressed image uses 16 clusters, which requires  $\log_2(16) = 4$  bits per pixel. So the compression factor is  $24/4 = 6$ .

## 1.2 Gaussian Mixtures

Now let's repeat the same process with Gaussian mixtures (with *full* covariances), with  $K = 5$  instead of K-means.

- (d) (10 pts) (**Autograder**) Work on the notebook `kmeans_gmm.ipynb` to implement the EM algorithm for GMM. You will need to implement `gmm.train_gmm()` to train a GMM model, which will be graded with the provided sample data and some random data. In order to get full points, your implementation should be efficient and fast enough (otherwise you will get only partial points).

[Hint 1: You may use `scipy.stats.multivariate_normal()` to compute the *log*-likelihood of the data.]

[Hint 2: You may use `scipy.special.logsumexp()` when computing  $\gamma(z_{nk})$ . You would need trick this because division by small probabilities can become computationally unstable when the likelihood values are too small. In practice, it is recommended to represent (possibly small) probabilities  $\mathcal{N}(\mathbf{x}^{(n)} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  with log-likelihood. Note that  $\mathcal{N}(\mathbf{x}^{(n)} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \exp[\log \mathcal{N}(\mathbf{x}^{(n)} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]$ .]

[Note: Do not use (and you don't need use) any other `scipy` or `scikit-learn` APIs.]

- (e) (3 pts) Train Gaussian mixtures with  $K = 5$  on the training data `mandrill-small.tif`. Use the provided initial mean and covariance matrices, and prior distribution of latent cluster (`initial_mu`, `initial_sigma`, and `initial_pi`) to ensure deterministic output.

Report in the **write-up**: the log-likelihood of the training data after running 50 EM steps. In addition, report the GMM model parameters  $\{(\pi_k, \boldsymbol{\mu}_k) : k = 1, \dots, 5\}$ . You do not need to write down  $\boldsymbol{\Sigma}_k$ . You can either write down the values, or simply attach the figure of visualization plots whose legend shows  $\pi_k$  and  $\boldsymbol{\mu}_k$  values.

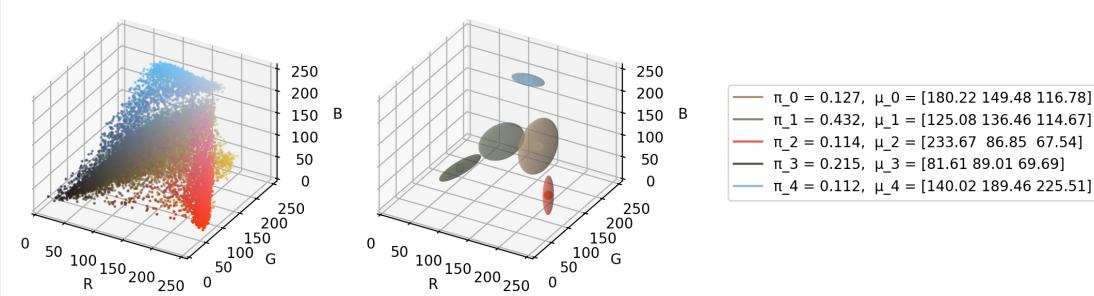
Solution.

Example Output:

```
Iteration  0: log-likelihood=-256674.70
Iteration  1: log-likelihood=-234014.08
Iteration  2: log-likelihood=-231557.90
Iteration  3: log-likelihood=-230375.66
...
Iteration 46: log-likelihood=-228103.70
Iteration 47: log-likelihood=-228103.60
Iteration 48: log-likelihood=-228103.51
Iteration 49: log-likelihood=-228103.43
```

The log-likelihood after 50 steps is around  $-228103.35$ . (Note that the log-likelihood printed above in the last step,  $-228103.43$  is the LL value right before M-step).

Model parameters and plot of the five mixture components:  $\boldsymbol{\pi} = (0.127, 0.432, 0.114, 0.215, 0.112)$  and  $\boldsymbol{\mu} \approx (180, 149, 116)(125, 136, 114)(233, 86, 67)(81, 89, 69)(140, 189, 225)$ .



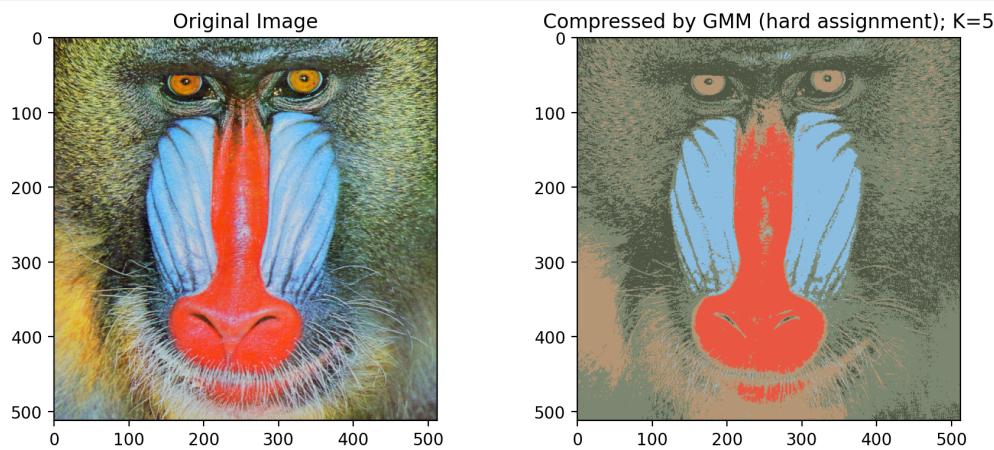
- (f) (2 pts) After training on the train image `mandrill-small.tiff`, read the test image `mandrill-large.tiff` and replace each pixel's  $(r, g, b)$  values with the value of latent cluster mean, where we use the **MAP** (Maximum A Posteriori) estimation for the latent cluster-assignment variable for each pixel.

Use the notebook's plotting code to display the original and compressed images side-by-side, and attach the plots to the **write-up**. (Note: you should have reasonable image quality/resolution to make the difference discernable). Also, measure and write down the mean pixel error between the original and compressed image.

Solution.

Mean pixel error: 16.3850

Compressed image:



## 2 [20 points] Expectation Maximization for GDA with missing labels

In this problem, we will work on using the EM algorithm for Gaussian Discrimination Analysis (GDA) with missing labels.

Suppose that you are given dataset where some portion of the data is labeled and the other portion is unlabeled. We want to learn a generative model over this partially-labeled dataset.<sup>1</sup> In particular, suppose there are  $l$  examples with labels and  $u$  examples without labels, i.e.,  $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(l)}, y^{(l)}), \mathbf{x}^{(l+1)}, \dots, \mathbf{x}^{(l+u)}\}$ .

We also make have the following assumptions:

- The data is real-valued and  $M$ -dimensional, i.e.,  $\mathbf{x} \in \mathbb{R}^M$
- The label  $y$  can take one of  $\{0, 1\}$  (i.e., binary classification problem).
- We model the data following the same assumption as in Gaussian Discrimination Analysis, i.e.,

$$P(\mathbf{x}, y) = P(y)P(\mathbf{x} | y) \quad (1)$$

$$P(y = j) = \begin{cases} \phi & \text{if } j = 1 \\ 1 - \phi & \text{if } j = 0 \end{cases} \quad (2)$$

$$P(\mathbf{x} | y = j) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j), j = 0 \text{ or } 1 \quad (3)$$

where  $\phi$  is a Bernoulli probability (i.e.,  $0 \leq \phi \leq 1$ ), and  $\boldsymbol{\mu}_j$  and  $\boldsymbol{\Sigma}_j$  are class-specific mean and covariance, respectively. For notational convenience, you can use  $\phi_1 = \phi$  and  $\phi_0 = 1 - \phi$ .

Further,  $\mathcal{N}(x; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$  is the multivariate Gaussian distribution which is defined as:

$$p(\mathbf{x} | y = j; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) = \mathcal{N}(x; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) = \frac{1}{(2\pi)^{\frac{M}{2}} |\boldsymbol{\Sigma}_j|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1}(\mathbf{x} - \boldsymbol{\mu}_j)\right) \quad (4)$$

Since we have unlabeled data, our goal is to maximize the following hybrid objective function:

$$\mathcal{J} = \sum_{i=1}^l \log p(\mathbf{x}^{(i)}, y^{(i)}) + \lambda \sum_{i=l+1}^{l+u} \log p(\mathbf{x}^{(i)}) \quad (5)$$

where  $\lambda$  is a hyperparameter that controls the weight of labeled and unlabeled data.

As we don't explicitly model the distribution  $p(\mathbf{x})$ , we use the law of total probability and rewrite the object function as:

$$\mathcal{J} = \sum_{i=1}^l \log p(\mathbf{x}^{(i)}, y^{(i)}) + \lambda \sum_{i=l+1}^{l+u} \log \sum_{j \in \{0, 1\}} p(\mathbf{x}^{(i)}, y^{(i)} = j) \quad (6)$$

This way the unlabeled training examples is using the same models as the labeled samples. Now we will be using EM algorithm to optimize this objective function.

*When deriving the solution, please show all the necessary steps in derivation and try to explain them to make the derivation as clearly understandable as possible.*

(Hint) You can use the fact:

$$p(\mathbf{x}^{(i)}, y^{(i)}) = \prod_{j \in \{0, 1\}} \left[ \frac{\phi_j}{(2\pi)^{\frac{M}{2}} |\boldsymbol{\Sigma}_j|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1}(\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)\right) \right]^{\mathbb{I}[y^{(i)} = j]} \quad (7)$$

---

<sup>1</sup>this type of learning formalism is called semi-supervised learning, which is a broad research field in machine learning.

## 2.1 Lower bound [3 points]

(a) Derive the variational lower bound  $\mathcal{L}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \phi)$  of the objective  $\mathcal{J}$ . Specifically, show that any arbitrary probability distribution  $q_i(y^{(i)} = j)$ , the lower bound of the objective function can be written as:

$$\mathcal{L}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \phi) = \sum_{i=1}^l \log p(\mathbf{x}^{(i)}, y^{(i)}) + \lambda \sum_{i=l+1}^{l+u} \sum_{j \in \{0,1\}} Q_{ij} \log \frac{p(\mathbf{x}^{(i)}, y^{(i)} = j)}{Q_{ij}} \quad (8)$$

where  $Q_{ij} \triangleq q_i(y^{(i)} = j)$  is a simplified shorthand notation. [Hint: you may want to use Jensen's Inequality.]

Solution

$$\mathcal{J} = \sum_{i=1}^l \log p(\mathbf{x}^{(i)}, y^{(i)}) + \lambda \sum_{i=l+1}^{l+u} \log \sum_{j \in \{0,1\}} p(\mathbf{x}^{(i)}, y^{(i)} = j) \quad (9)$$

$$= \sum_{i=1}^l \log p(\mathbf{x}^{(i)}, y^{(i)}) + \lambda \sum_{i=l+1}^{l+u} \log \sum_{j \in \{0,1\}} \boxed{Q_{ij}} \frac{p(\mathbf{x}^{(i)}, y^{(i)} = j)}{\boxed{Q_{ij}}} \quad (10)$$

$$\geq \sum_{i=1}^l \log p(\mathbf{x}^{(i)}, y^{(i)}) + \lambda \sum_{i=l+1}^{l+u} \sum_{j \in \{0,1\}} Q_{ij} \log \frac{p(\mathbf{x}^{(i)}, y^{(i)} = j)}{Q_{ij}} \quad (\because \text{Jensen's Inequality}) \quad (11)$$

$$= \mathcal{L}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \phi) \quad (12)$$

## 2.2 E-step [2 points]

(b) Write down the E-step; specifically, define the distribution  $Q_{ij} = q_i(y^{(i)} = j)$ .

Solution.

The lower bound is tight when  $Q_{ij}$  matches the posterior  $P(y^{(i)} = j | \mathbf{x}^{(i)})$ . Using the Bayes rule,

$$Q_{ij} = P(y^{(i)} = j | \mathbf{x}^{(i)}) = \frac{P(\mathbf{x}^{(i)} | y^{(i)} = j)P(y^{(i)} = j)}{P(\mathbf{x}^{(i)})} \quad (13)$$

$$= \frac{P(\mathbf{x}^{(i)} | y^{(i)} = j)P(y^{(i)} = j)}{\sum_{j=0,1} P(\mathbf{x}^{(i)} | y^{(i)} = j)P(y^{(i)} = j)} \quad (14)$$

$$= \frac{\phi_j \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\phi_1 \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + \phi_0 \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)} \quad (15)$$

We can clearly see that  $Q_{i0} + Q_{i1} = 1$  and  $Q_{ij} \geq 0$ .

## 2.3 M-step for $\boldsymbol{\mu}_k$ [6 points]

(c) Derive the M-step update rule for  $\boldsymbol{\mu}_k$  where  $k = 0$  or  $1$ , while holding  $Q_i$ 's (which you obtained in (a)) fixed. Also, explain in words (English) what *intuitively*  $\boldsymbol{\mu}_k$  looks like in terms of  $\mathbf{x}^{(i)}$ 's (each of labeled and unlabeled) and pseudo-counts.

## Solution

In the M-step (for  $\boldsymbol{\mu}_k$ ), we want to find  $\boldsymbol{\mu}_k$  that maximizes the variational lower bound  $\mathcal{L}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \phi)$ . Therefore, let's take a derivative of the lower bound with respect to  $\boldsymbol{\mu}_k$ :

$$\nabla_{\boldsymbol{\mu}_k} \mathcal{L}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \phi) \quad (16)$$

$$= \nabla_{\boldsymbol{\mu}_k} \left( \sum_{i=1}^l \log p(\mathbf{x}^{(i)}, y^{(i)}) + \lambda \sum_{i=l+1}^{l+u} \sum_{j \in \{0,1\}} Q_{ij} \log \frac{p(\mathbf{x}^{(i)}, y^{(i)} = j)}{Q_{ij}} \right) \quad (17)$$

$$= \nabla_{\boldsymbol{\mu}_k} \left( \sum_{i=1}^l \log p(\mathbf{x}^{(i)}, y^{(i)}) + \lambda \sum_{i=l+1}^{l+u} \sum_{j \in \{0,1\}} Q_{ij} \log p(\mathbf{x}^{(i)}, y^{(i)} = j) \right) \quad (18)$$

$$(\because \text{for a fixed } Q_{ij}, \text{ it can be treated as a constant}) \quad (19)$$

$$= \nabla_{\boldsymbol{\mu}_k} \left( \sum_{i=1}^l \log \underbrace{p(\mathbf{x}^{(i)}, y^{(i)})}_{\text{Apply hint here}} \right) + \nabla_{\boldsymbol{\mu}_k} \left( \lambda \sum_{i=l+1}^{l+u} \sum_{j \in \{0,1\}} Q_{ij} \log p(\mathbf{x}^{(i)}, y^{(i)} = j) \right) \quad (20)$$

$$= \nabla_{\boldsymbol{\mu}_k} \left( \sum_{i=1}^l \sum_{j \in \{0,1\}} \mathbb{I}[y^{(i)} = j] \log \left[ \frac{\phi_j}{(2\pi)^{\frac{M}{2}} |\boldsymbol{\Sigma}_j|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j) \right) \right] \right) + \quad (21)$$

$$\lambda \nabla_{\boldsymbol{\mu}_k} \left( \sum_{i=l+1}^{l+u} \sum_{j \in \{0,1\}} Q_{ij} \log \left[ \frac{\phi_j}{(2\pi)^{\frac{M}{2}} |\boldsymbol{\Sigma}_j|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j) \right) \right] \right) \quad (22)$$

$$= \nabla_{\boldsymbol{\mu}_k} \left( \sum_{i=1}^l \sum_{j \in \{0,1\}} \mathbb{I}[y^{(i)} = j] \left( -\frac{1}{2} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j) \right) + \text{const} \right) + \quad (23)$$

$$\lambda \nabla_{\boldsymbol{\mu}_k} \left( \sum_{i=l+1}^{l+u} \sum_{j \in \{0,1\}} Q_{ij} \left( -\frac{1}{2} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j) \right) + \text{const} \right) \quad (24)$$

$$= \left( \sum_{i=1}^l \mathbb{I}[y^{(i)} = k] \left( -\frac{1}{2} \cdot 2 \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k) (-1) \right) \right) + \lambda \left( \sum_{i=l+1}^{l+u} Q_{ik} \left( -\frac{1}{2} \cdot 2 \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k) (-1) \right) \right) \quad (25)$$

$$= \sum_{i=1}^l \mathbb{I}[y^{(i)} = k] \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k) + \lambda \sum_{i=l+1}^{l+u} Q_{ik} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k) \quad (26)$$

$$(27)$$

Setting the gradient  $\nabla_{\boldsymbol{\mu}_k} \mathcal{L}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \phi_k) = 0$  gives:

$$\sum_{i=1}^l \mathbb{I}[y^{(i)} = k] \boldsymbol{\Sigma}_k^{-1} \mathbf{x}^{(i)} - \sum_{i=1}^l \mathbb{I}[y^{(i)} = k] \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k + \lambda \sum_{i=l+1}^{l+u} Q_{ik} \boldsymbol{\Sigma}_k^{-1} \mathbf{x}^{(i)} - \lambda \sum_{i=l+1}^{l+u} Q_{ik} \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k = 0 \quad (28)$$

$$\left( \sum_{i=1}^l \mathbb{I}[y^{(i)} = k] \mathbf{x}^{(i)} + \lambda \sum_{i=l+1}^{l+u} Q_{ik} \mathbf{x}^{(i)} \right) - \left( \sum_{i=1}^l \mathbb{I}[y^{(i)} = k] + \lambda \sum_{i=l+1}^{l+u} Q_{ik} \right) \boldsymbol{\mu}_k = 0 \quad (29)$$

Finally,

$$\boldsymbol{\mu}_k = \frac{\sum_{i=1}^l \mathbb{I}[y^{(i)} = k] \mathbf{x}^{(i)} + \lambda \sum_{i=l+1}^{l+u} Q_{ik} \mathbf{x}^{(i)}}{\sum_{i=1}^l \mathbb{I}[y^{(i)} = k] + \lambda \sum_{i=l+1}^{l+u} Q_{ik}} \quad (30)$$

Intuitively speaking, this is a weighted average of  $\mathbf{x}^{(i)}$ 's where you give the count of 1 for each labeled example that belongs to label  $k$ , and  $\lambda Q_{ik}$  for each unlabeled example.

## 2.4 M-step for $\phi$ [6 points]

(d) Derive the M-step update rule for  $\phi \in \mathbb{R}$ , while holding  $Q_i$ 's (which you obtained in (a)) fixed. Also, explain in words (English) what *intuitively*  $\phi$  looks like in terms of  $\mathbf{x}^{(i)}$ 's (each of labeled and unlabeled) and pseudo-counts.

Solution

Similarly, let's take a derivative of the lower bound with respect to  $\phi$  (we skip some steps that are repetitive as in 4(c)).

$$\frac{\partial}{\partial \phi} \mathcal{L}(\boldsymbol{\mu}, \Sigma, \phi) \quad (31)$$

$$= \frac{\partial}{\partial \phi} \left( \sum_{i=1}^l \log p(\mathbf{x}^{(i)}, y^{(i)}) + \lambda \sum_{i=l+1}^{l+u} \sum_{j \in \{0,1\}} Q_{ij} \log p(\mathbf{x}^{(i)}, y^{(i)} = j) \right) \quad (32)$$

$$= \frac{\partial}{\partial \phi} \left( \sum_{i=1}^l \sum_{j \in \{0,1\}} \mathbb{I}[y^{(i)} = j] \log \left[ \frac{\phi_j}{(2\pi)^{\frac{M}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^\top \Sigma_j^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j) \right) \right] \right) + \quad (33)$$

$$\lambda \frac{\partial}{\partial \phi} \left( \sum_{i=l+1}^{l+u} \sum_{j \in \{0,1\}} Q_{ij} \log \left[ \frac{\phi_j}{(2\pi)^{\frac{M}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^\top \Sigma_j^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j) \right) \right] \right) \quad (34)$$

$$= \frac{\partial}{\partial \phi} \left( \sum_{i=1}^l \sum_{j \in \{0,1\}} \mathbb{I}[y^{(i)} = j] \log \phi_j + \text{const} \right) + \lambda \frac{\partial}{\partial \phi} \left( \sum_{i=l+1}^{l+u} \sum_{j \in \{0,1\}} Q_{ij} \log \phi_j + \text{const} \right) \quad (35)$$

$$= \sum_{i=1}^l \left( \mathbb{I}[y^{(i)} = 0] \frac{\partial}{\partial \phi} \log \phi_0 + \mathbb{I}[y^{(i)} = 1] \frac{\partial}{\partial \phi} \log \phi_1 \right) + \lambda \sum_{i=l+1}^{l+u} \left( Q_{i0} \frac{\partial}{\partial \phi} \log \phi_0 + Q_{i1} \frac{\partial}{\partial \phi} \log \phi_1 \right) \quad (36)$$

(37)

Now note that:

$$\frac{\partial}{\partial \phi} \log \phi_0 = \frac{\partial}{\partial \phi} \log(1 - \phi) = -\frac{1}{1 - \phi} \quad (38)$$

$$\frac{\partial}{\partial \phi} \log \phi_1 = \frac{\partial}{\partial \phi} \log \phi = \frac{1}{\phi} \quad (39)$$

Plugging these back and solving  $\frac{\partial}{\partial \phi} \mathcal{L}(\boldsymbol{\mu}_k, \Sigma_k, \phi) = 0$ :

$$\sum_{i=1}^l \left( \mathbb{I}[y^{(i)} = 1] \frac{1}{\phi} - \mathbb{I}[y^{(i)} = 0] \frac{1}{1 - \phi} \right) + \lambda \sum_{i=l+1}^{l+u} \left( Q_{i1} \frac{1}{\phi} - Q_{i0} \frac{1}{1 - \phi} \right) = 0 \quad (40)$$

$$\iff \left( \sum_{i=1}^l \mathbb{I}[y^{(i)} = 1] + \lambda \sum_{i=l+1}^{l+u} Q_{i1} \right) \frac{1}{\phi} = \left( \sum_{i=1}^l \mathbb{I}[y^{(i)} = 0] + \lambda \sum_{i=l+1}^{l+u} Q_{i0} \right) \frac{1}{1-\phi} \quad (41)$$

$$\iff (1-\phi) \cdot \left( \sum_{i=1}^l \mathbb{I}[y^{(i)} = 1] + \lambda \sum_{i=l+1}^{l+u} Q_{i1} \right) = \phi \cdot \left( \sum_{i=1}^l \mathbb{I}[y^{(i)} = 0] + \lambda \sum_{i=l+1}^{l+u} Q_{i0} \right) \quad (42)$$

$$\iff \left( \sum_{i=1}^l \mathbb{I}[y^{(i)} = 1] + \lambda \sum_{i=l+1}^{l+u} Q_{i1} \right) = \quad (43)$$

$$\phi \left( \sum_{i=1}^l \mathbb{I}[y^{(i)} = 1] + \lambda \sum_{i=l+1}^{l+u} Q_{i1} + \sum_{i=1}^l \mathbb{I}[y^{(i)} = 0] + \lambda \sum_{i=l+1}^{l+u} Q_{i0} \right) \quad (44)$$

(45)

Finally, we get

$$\phi = \frac{\sum_{i=1}^l \mathbb{I}[y^{(i)} = 1] + \lambda \sum_{i=l+1}^{l+u} Q_{i1}}{\sum_{i=1}^l \mathbb{I}[y^{(i)} = 1] + \lambda \sum_{i=l+1}^{l+u} Q_{i1} + \sum_{i=1}^l \mathbb{I}[y^{(i)} = 0] + \lambda \sum_{i=l+1}^{l+u} Q_{i0}} \quad (46)$$

$$= \frac{\sum_{i=1}^l \mathbb{I}[y^{(i)} = 1] + \lambda \sum_{i=l+1}^{l+u} Q_{i1}}{l + \lambda u} \quad (47)$$

Intuitively speaking, this is an weighted frequency where you give the count of 1 for each labeled example that belongs to label 1, and  $\lambda Q_{i1}$  for each unlabeled example.

## 2.5 M-step for $\Sigma_k$ [3 points]

- (e) Finally, let's think about the M-step update rule for  $\Sigma_k$  where  $k = 0$  or  $1$ . Since we know the derivation is very similar to the case of GDA (and GMM M-step), we do not require you to repeat the similar step as you have already worked on other two M-step update rules. Write down the M-step update rule for  $\Sigma_k$ , without derivation, based on your guess and the analogy we have seen. Also, explain in words (English) what *intuitively*  $\Sigma_k$  looks like in terms of  $\mathbf{x}^{(i)}$ 's (each of labeled and unlabeled) and pseudo-counts.

Solution

$$\Sigma_k = \frac{\sum_{i=1}^l \mathbb{I}[y^{(i)} = k] (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)^\top + \lambda \sum_{i=l+1}^{l+u} Q_{ik} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)^\top}{\sum_{i=1}^l \mathbb{I}[y^{(i)} = k] + \lambda \sum_{i=l+1}^{l+u} Q_{ik}} \quad (48)$$

Intuitively speaking, this is an weighted covariance where you give the count of 1 for each labeled example that belongs to label  $k$ , and  $\lambda Q_{ik}$  for each unlabeled example.

### 3 [20 points] PCA and eigenfaces

- (a) (8 pts) In lecture, we derived PCA from the “maximizing variance” viewpoint. In this problem, we will take the “minimizing squared error” viewpoint. Let  $K \in \{1, \dots, D\}$  be arbitrary and let  $\mathbf{x}^{(n)} \in \mathbb{R}^D$ . Let  $\mathcal{U} = \{\mathbf{U} = [\mathbf{u}_1 \cdots \mathbf{u}_K] \in \mathbb{R}^{D \times K} \mid \{\mathbf{u}_i\}_{i=1}^K \text{ are orthonormal vectors}\}$ , where  $\mathbf{u}_i$  is the  $i$ -th column vector of  $\mathbf{U}$ .

Let's define the objective function for minimizing the distortion error:

$$\mathcal{J} = \frac{1}{N} \sum_{n=1}^N \left\| \mathbf{x}^{(n)} - \mathbf{U} \mathbf{U}^\top \mathbf{x}^{(n)} \right\|^2 = \frac{1}{N} \sum_{n=1}^N \left\| \mathbf{x}^{(n)} - \sum_{i=1}^K \mathbf{u}_i \mathbf{u}_i^\top \mathbf{x}^{(n)} \right\|^2 \quad (49)$$

Here,  $\mathbf{U} \mathbf{U}^\top \mathbf{x}^{(n)}$  is called a projection of  $\mathbf{x}^{(n)}$  into the subspace spanned by  $\mathbf{u}_i$ 's, and we can denote the projection  $\tilde{\mathbf{x}}^{(n)} = \mathbf{U} \mathbf{U}^\top \mathbf{x}^{(n)}$  as in the lecture.

Specifically, show that:

$$\mathcal{J} = \sum_{i=1}^D \lambda_i - \sum_{i=1}^K \mathbf{u}_i^\top \mathbf{S} \mathbf{u}_i \quad (50)$$

where  $\mathbf{S}$  is the data covariance matrix  $\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}^{(n)} - \bar{\mathbf{x}})(\mathbf{x}^{(n)} - \bar{\mathbf{x}})^\top$ ,  $\bar{\mathbf{x}}$  is the data mean vector, and  $\lambda_1 \geq \dots \geq \lambda_d$  are the (ordered) eigenvalues of  $\mathbf{S}$ . Since the first term is a constant, the above equation implies that minimizing the squared error after projection is equivalent to maximizing the variance, as we have shown in the lecture slides.

With further simplification, show that the minimum distortion error corresponds to the sum of the  $D-K$  smallest eigenvalues of  $\mathbf{S}$ , i.e.,

$$\min_{\mathbf{U} \in \mathcal{U}} \mathcal{J} = \sum_{k=K+1}^D \lambda_k. \quad (51)$$

and that the  $\mathbf{u}_i$ 's that minimize  $\mathcal{J}$  are indeed the  $K$  eigenvectors of  $\mathbf{S}$  corresponding to the (ordered) eigenvalues  $\{\lambda_i\}_{i=1}^K$ . After showing Eq.(50), it is okay to use the fact (without proof) that the optimal solution  $\mathbf{u}_i$ 's that maximize  $\sum_{i=1}^K \mathbf{u}_i^\top \mathbf{S} \mathbf{u}_i$  is to pick the top- $K$  eigenvectors of  $\mathbf{S}$  (i.e.,  $\mathbf{u}_1, \dots, \mathbf{u}_K$  corresponding to the largest  $K$  eigenvalues of  $\mathbf{S}$  in descending order), as we already have seen in the lecture.

[Hint 1: You may assume that the data is zero-centered, i.e.,  $\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} = \mathbf{0} \in \mathbb{R}^D$  without loss of generality. Be sure to mention it if you make such an assumption.]

[Hint 2: You can rewrite the objective as  $\mathcal{J} = \frac{1}{N} \|\mathbf{X} - \mathbf{U} \mathbf{U}^\top \mathbf{X}\|_F^2$ , where  $\mathbf{X} \in \mathbb{R}^{D \times N}$  is the matrix that stacks all the data points  $\{\mathbf{x}^{(n)}\}_{n=1}^N$  as column vectors, and the  $\|\cdot\|_F$  denotes the Frobenius norm:

$$\|A\|_F = \sqrt{\sum_{i,j} (A_{ij})^2}, \quad (52)$$

and use the fact  $\|A\|_F^2 = \text{tr}(A^\top A)$  and  $\text{tr}(\mathbf{S}) = \sum_i \lambda_i$ . Also note that there are many possible approaches for proving the claim, so you do not have to use this fact if you take a different approach.]

Solution.

We will prove the equation Equation (51) and show that the optimal  $\mathbf{u}$  vectors are indeed the eigenvectors of  $\mathbf{S}$ .

Without loss of generality we can assume  $\bar{\mathbf{x}} = 0$ , i.e., the data is zero-centered. We can simplify

the reconstruction error as follows:

$$\mathcal{J} = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}^{(i)} - \mathbf{U} \mathbf{U}^\top \mathbf{x}^{(i)} \right\|^2 \quad (53)$$

$$= \frac{1}{N} \left\| \mathbf{X} - \mathbf{U} \mathbf{U}^\top \mathbf{X} \right\|_F^2 \quad (54)$$

$$= \text{tr} \left( \frac{1}{N} (\mathbf{X} - \mathbf{U} \mathbf{U}^\top \mathbf{X})^\top (\mathbf{X} - \mathbf{U} \mathbf{U}^\top \mathbf{X}) \right) \quad (55)$$

$$= \text{tr} \left( \frac{1}{N} \mathbf{X}^\top \mathbf{X} \right) - \text{tr} \left( \frac{2}{N} \mathbf{X} \mathbf{U} \mathbf{U}^\top \mathbf{X} \right) + \text{tr} \left( \frac{1}{N} \mathbf{X}^\top \mathbf{U} \underbrace{\mathbf{U}^\top \mathbf{U}}_{=\mathbf{I}_K \in \mathbb{R}^{K \times K}} \mathbf{U}^\top \mathbf{X} \right) \quad (56)$$

$$= \text{tr} \left( \frac{1}{N} \mathbf{X}^\top \mathbf{X} \right) - \text{tr} \left( \frac{2}{N} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X} \right) + \text{tr} \left( \frac{1}{N} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X} \right) \quad (57)$$

$$= \text{tr} \left( \frac{1}{N} \mathbf{X}^\top \mathbf{X} \right) - \text{tr} \left( \frac{1}{N} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X} \right). \quad (58)$$

We used the fact and that  $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_K$  because  $\mathbf{U}$  is an orthogonal matrix (Note that  $\mathbf{U} \mathbf{U}^\top \neq \mathbf{I}_D \in \mathbb{R}^{D \times D}$  since  $K < D$ ).

Now, recall  $\frac{1}{N} \mathbf{X} \mathbf{X}^\top = \mathbf{S}$ . Using the properties of the trace operator (e.g.,  $\text{tr}(ABCD) = \text{tr}(CDAB)$  and  $\text{tr}(AB) = \text{tr}(BA)$ ) we obtain:

$$\mathcal{J} = \text{tr} \left( \frac{1}{N} \mathbf{X}^\top \mathbf{X} \right) - \text{tr} \left( \frac{1}{N} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X} \right) = \text{tr} \left( \frac{1}{N} \mathbf{X} \mathbf{X}^\top \right) - \text{tr} \left( \mathbf{U}^\top \left( \frac{1}{N} \mathbf{X} \mathbf{X}^\top \right) \mathbf{U} \right) \quad (59)$$

$$= \text{tr}(\mathbf{S}) - \text{tr}(\mathbf{U}^\top \mathbf{S} \mathbf{U}) \quad (60)$$

$$= \sum_{i=1}^D \lambda_i - \sum_{i=1}^K \mathbf{u}_i^\top \mathbf{S} \mathbf{u}_i, \quad (61)$$

where  $\lambda_i$ 's are the (ordered) eigenvalues of  $\mathbf{S}$ .

Let's first show that  $\mathbf{u}_i$ 's that minimize the Eq. (61) are the eigenvectors of  $\mathbf{S}$ :

$$\begin{aligned} \arg \min_{\mathbf{U}: \|\mathbf{u}_j\|_2=1, \forall j} \mathcal{J} &= \arg \min_{\mathbf{U}: \|\mathbf{u}_j\|_2=1, \forall j} \left[ \sum_{i=1}^D \lambda_i - \sum_{i=1}^K \mathbf{u}_i^\top \mathbf{S} \mathbf{u}_i \right] \\ &= \arg \max_{\mathbf{U}: \|\mathbf{u}_j\|_2=1, \forall j} \left[ \sum_{i=1}^K \mathbf{u}_i^\top \mathbf{S} \mathbf{u}_i \right], \end{aligned}$$

because  $\sum_{i=1}^D \lambda_i$  it is constant w.r.t.  $\mathbf{u}_i$ 's. This is exactly the variance maximizing objective in the lecture note. We already showed that the  $\mathbf{u}_i$ 's that maximize the variance are the eigenvectors of  $\mathbf{S}$ ,  $\mathbf{u}_i$ 's. If we plug the eigenvectors  $\mathbf{u}_i$ 's into the  $\mathbf{u}_i$ 's, we get the following:

$$\min_{\mathbf{U} \in \mathcal{U}} \mathcal{J} = \min_{\|\mathbf{u}_j\|_2=1, \forall j} \left[ \sum_{i=1}^D \lambda_i - \sum_{i=1}^K \mathbf{u}_i^\top \mathbf{S} \mathbf{u}_i \right] = \sum_{i=1}^D \lambda_i - \sum_{i=1}^K \lambda_i = \sum_{i=K+1}^D \lambda_i. \quad (62)$$

Q.E.D.

Now, you will apply PCA to face images. The principal components (eigenvectors) of the face images are called eigenfaces.

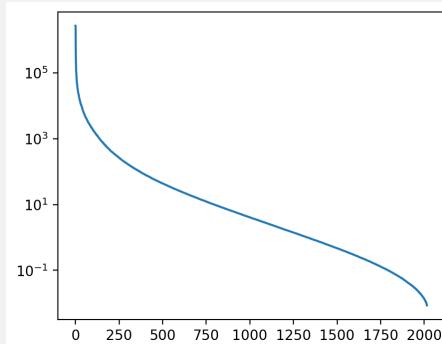
- (b) (4 pts) (**Autograder**) Work on the provided code `pca.ipynb` and `pca.py` to implement PCA. Your code will be graded by the correctness on the sample face dataset and some other randomly-generated dataset.
- (c) (3 pts) By regarding each image as a vector in a high dimensional space, perform PCA on the face images (sort the eigenvalues in descending order). In the write-up, report the eigenvalues corresponding to the first 10 principal components, and plot all the eigenvalues (in sorted order) where x-axis is the index of corresponding principal components and y-axis is the eigenvalue. Use log scale for the y-axis.

Solution.

The largest 10 eigenvalues are:

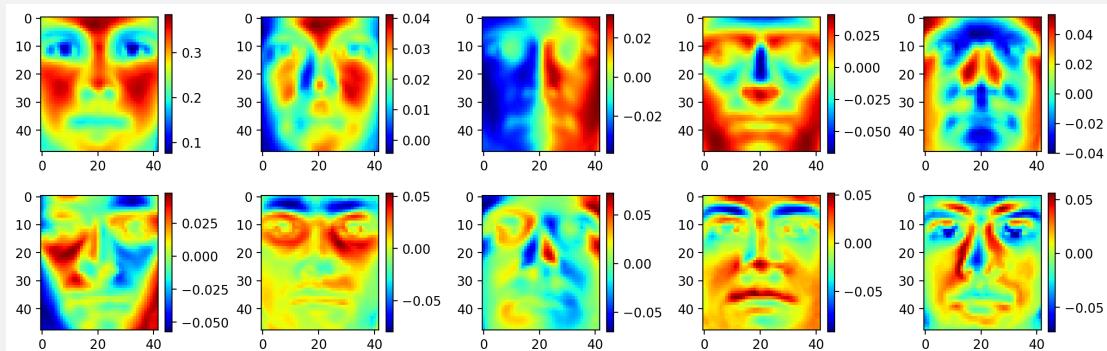
2718207.41, 2610783.99, 365363.88, 211062.37, 110557.41,  
104672.41, 78480.41, 67004.29, 52571.01, 49176.70

Plot:



- (d) (3 pts) Plot and attach to your write-up: a  $2 \times 5$  array of subplots showing the first 10 principal components/eigenvectors ("eigenfaces") (sorted according to the descending eigenvalues) as images, treating the mean of images as the first principal component. Comment on what facial or lighting variations some of the different principal components are capturing (Note: you don't need to comment for all the images. Just pick a few that capture some salient aspects of image).

Solution.



Examples of some aspects that the eigenfaces are possibly capturing:

- Third and sixth eigenfaces capture the difference between left and right side of face.
- Fourth eigenface captures shape of jaw
- Fifth eigenface captures shape of nose
- Seventh eigenface captures shape of eyes
- etc.

(e) (2 pts) Eigenfaces are a set of bases for all the images in the dataset (every image can be represented as a linear combination of these eigenfaces). Suppose we have  $L$  eigenfaces in total. Then we can use the  $L$  coefficients (of the bases, i.e. the eigenfaces) to represent an image. Moreover, we can use the first  $K (< L)$  eigenfaces to reconstruct the face image approximately (correspondingly use  $K$  coefficients to represent the image). In this case, we reduce the dimension of the representation of images from  $L$  to  $K$ . To determine the proper  $K$  to use, we will check the percentage of variance that has been preserved (recall that the basic idea of PCA is preserving variance). Specifically, we define total variance

$$v(K) = \sum_{i=1}^K \lambda_i$$

where  $1 \leq K < L$  and  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_L$  are eigenvalues. Then the percentage of total variance is

$$\frac{v(K)}{v(L)}$$

How many principal components are needed to represent 95% of the total variance? How about 99%? What is the percentage of reduction in dimension in each case?

Solution.

- 95% variance: 43 components. 97.87% reduction
- 99% variance: 167 components. 91.72% reduction

## 4 [10 points] Independent Component Analysis

In this problem, you will implement maximum-likelihood Independent Component Analysis (ICA) for blind audio separation. As we learned in the lecture, the maximum-likelihood ICA minimizes the following loss:

$$\ell(W) = \sum_{i=1}^N \left( \sum_{j=1}^m \log g' \left( w_j^\top x^{(i)} \right) + \log |W| \right), \quad (63)$$

where  $N$  is the number of time steps,  $m$  is the number of independent sources,  $W$  is the transformation matrix representing a concatenation of  $w_j$ 's, and  $g(s) = 1 / (1 + e^{-s})$  is the sigmoid function. This link has some nice demos of blind audio separation: [https://cnl.salk.edu/~tewon/Blind/blind\\_audio.html](https://cnl.salk.edu/~tewon/Blind/blind_audio.html).

We provided the starter code `ica.py` and the data `ica_data.dat`, which contains mixed sound signals from multiple microphones. Run the provided notebook `ica.ipynb` to load the data and run your ICA implementation from `ica.py`.

(a) (6 points) Autograder Implement ICA by filling in the `ica.py` file.

(b) (4 points) Run your ICA implementation in the `ica.ipynb` notebook. To make sure your code is correct, you should listen to the resulting unmixed sources. (Some overlap in the sources may be present, but the different sources should be pretty clearly separated.)

Report the  $W$  matrix you found and submit the notebook `ica.ipynb` (along with `ica.py`) to the autograder. Make sure the audio tracks are audible in the notebook before submitting. You do **not** need to submit your unmixed sound files (`ica_unmixed_track_X.wav`).

Solution.

The correct  $W$  is:

$$\begin{pmatrix} 72.15081922 & 28.62441682 & 25.91040458 & -17.2322227 & -21.191357 \\ 13.45886116 & 31.94398247 & -4.03003982 & -24.0095722 & 11.89906179 \\ 18.89688784 & -7.80435173 & 28.71469558 & 18.14356811 & -21.17474522 \\ -6.0119837 & -4.15743607 & -1.01692289 & 13.87321073 & -5.26252289 \\ -8.74061186 & 22.55821897 & 9.61289023 & 14.73637074 & 45.28841827 \end{pmatrix}$$

The five sources would be:

- Source 0: A man speaking
- Source 1: A woman speaking
- Source 2: A man speaking
- Source 3: Some music (Beethoven's Symphony No.5 in C minor)
- Source 4: A man speaking

## 5 [25 points] Conditional Variational Autoencoders

In this problem, you will implement a conditional variational autoencoder (CVAE) from [?] and train it on the MNIST dataset.

- (a) [5 points] Derive the variational lower bound of a conditional variational autoencoder. Show that:

$$\begin{aligned} \log p_\theta(\mathbf{x}|\mathbf{y}) &\geq \mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{y}) \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})} [\log p_\theta(\mathbf{x}|\mathbf{z}, \mathbf{y})] - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \| p_\theta(\mathbf{z}|\mathbf{y})) , \end{aligned} \quad (64)$$

where  $\mathbf{x}$  is a binary vector of dimension  $d$ ,  $\mathbf{y}$  is a one-hot vector of dimension  $c$  defining a class,  $\mathbf{z}$  is a vector of dimension  $m$  sampled from the posterior distribution  $q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})$ . The posterior distribution is modeled by a neural network of parameters  $\phi$ . The generative distribution  $p_\theta(\mathbf{x}|\mathbf{y})$  is modeled by another neural network of parameters  $\theta$ . Similar to the VAE that we learned in the class, we assume the conditional independence on the components of  $\mathbf{z}$ : i.e.,  $q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) = \prod_{j=1}^m q_\phi(z_j|\mathbf{x}, \mathbf{y})$ , and  $p_\theta(\mathbf{z}|\mathbf{y}) = \prod_{j=1}^m p_\theta(z_j|\mathbf{y})$ .

Solution.

From the second line in the following, we will denote the  $\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}, \mathbf{y})}$  as  $\mathbb{E}_\mathbf{z}$ ,  $q_\phi$  as  $q$ ,  $p_\theta$  as  $p$  for notational simplicity.

$$\log p_\theta(\mathbf{x}|\mathbf{y}) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}, \mathbf{y})} [\log p_\theta(\mathbf{x}|\mathbf{y})] \quad (65)$$

$$= \mathbb{E}_\mathbf{z} \left[ \log \frac{p(\mathbf{x}|\mathbf{y}, \mathbf{z})p(\mathbf{z}|\mathbf{y})}{p(\mathbf{z}|\mathbf{x}, \mathbf{y})} \right] \quad (66)$$

$$= \mathbb{E}_\mathbf{z} \left[ \log \frac{p(\mathbf{x}|\mathbf{y}, \mathbf{z})p(\mathbf{z}|\mathbf{y})}{p(\mathbf{z}|\mathbf{x}, \mathbf{y})} \frac{q(\mathbf{z}|\mathbf{x}, \mathbf{y})}{q(\mathbf{z}|\mathbf{x}, \mathbf{y})} \right] \quad (67)$$

$$= \mathbb{E}_\mathbf{z} [\log p(\mathbf{x}|\mathbf{y}, \mathbf{z})] - \mathbb{E}_\mathbf{z} \left[ \log \frac{q(\mathbf{z}|\mathbf{x}, \mathbf{y})}{p(\mathbf{z}|\mathbf{y})} \right] + \mathbb{E}_\mathbf{z} \left[ \log \frac{q(\mathbf{z}|\mathbf{x}, \mathbf{y})}{p(\mathbf{z}|\mathbf{x}, \mathbf{y})} \right] \quad (68)$$

$$= \mathbb{E}_\mathbf{z} [\log p(\mathbf{x}|\mathbf{y}, \mathbf{z})] - D_{KL}(q(\mathbf{z}|\mathbf{x}, \mathbf{y}) \| p(\mathbf{z}|\mathbf{y})) + \underbrace{D_{KL}(q(\mathbf{z}|\mathbf{x}, \mathbf{y}) \| p(\mathbf{z}|\mathbf{x}, \mathbf{y}))}_{\geq 0} \quad (69)$$

$$\geq \mathbb{E}_\mathbf{z} [\log p(\mathbf{x}|\mathbf{y}, \mathbf{z})] - D_{KL}(q(\mathbf{z}|\mathbf{x}, \mathbf{y}) \| p(\mathbf{z}|\mathbf{y})) \quad (70)$$

$$= \mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{y}) \quad (71)$$

- (b) [8 points] Derive the analytical solution to the KL-divergence between two Gaussian distributions  $D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \| p_\theta(\mathbf{z}|\mathbf{y}))$ . Let us assume that  $p_\theta(\mathbf{z}|\mathbf{y}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and show that:

$$D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \| p_\theta(\mathbf{z}|\mathbf{y})) = -\frac{1}{2} \sum_{j=1}^m (1 + \log(\sigma_j^2) - \mu_j^2 - \sigma_j^2), \quad (72)$$

where  $\mu_j$  and  $\sigma_j$  are the outputs of the neural network that estimates the parameters of the posterior distribution  $q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})$ .

You can assume without proof that

$$D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \| p_\theta(\mathbf{z}|\mathbf{y})) = \sum_{j=1}^m D_{KL}(q_\phi(z_j|\mathbf{x}, \mathbf{y}) \| p_\theta(z_j|\mathbf{y})). \quad (73)$$

This is a consequence of conditional independence of the components of  $\mathbf{z}$ .

Solution.

**Answer:** First, remember that all  $z_j$  are assumed to be independent, i.e.  $p_\theta(\mathbf{z}|\mathbf{y}) = \prod_j p_\theta(z_j|\mathbf{y})$ . Therefore, we have

$$D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \| p_\theta(\mathbf{z}|\mathbf{y})) \quad (74)$$

$$= \int_{\mathbf{z}} q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})}{p_\theta(\mathbf{z}|\mathbf{y})} d\mathbf{z} \quad (75)$$

$$= \int_{\mathbf{z}} q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \log q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) d\mathbf{z} - \int_{\mathbf{z}} q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \log p_\theta(\mathbf{z}|\mathbf{y}) d\mathbf{z} \quad (76)$$

$$= \int_{\mathbf{z}} q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \log \prod_j q_\phi(z_j|\mathbf{x}, \mathbf{y}) d\mathbf{z} - \int_{\mathbf{z}} q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \log \prod_j p_\theta(z_j|\mathbf{y}) d\mathbf{z} \quad (77)$$

$$= \int_{\mathbf{z}} q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \sum_j \log q_\phi(z_j|\mathbf{x}, \mathbf{y}) d\mathbf{z} - \int_{\mathbf{z}} q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \sum_j \log p_\theta(z_j|\mathbf{y}) d\mathbf{z} \quad (78)$$

$$= \sum_j \underbrace{\int_{\mathbf{z}} q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \log q_\phi(z_j|\mathbf{x}, \mathbf{y}) d\mathbf{z}}_{(\diamond)} - \sum_j \underbrace{\int_{\mathbf{z}} q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \log p_\theta(z_j|\mathbf{y}) d\mathbf{z}}_{(\heartsuit)} \quad (79)$$

Now, note that the term  $(\diamond)$  for each  $j$  in the sum, we have

$$\int_{\mathbf{z}} q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \log q_\phi(z_j|\mathbf{x}, \mathbf{y}) d\mathbf{z} = \int_{\mathbf{z}} \prod_{j'} q_\phi(z_{j'}|\mathbf{x}, \mathbf{y}) \log q_\phi(z_j|\mathbf{x}, \mathbf{y}) d\mathbf{z} \quad (80)$$

$$= \int_{z_j} q_\phi(z_j|\mathbf{x}, \mathbf{y}) \log q_\phi(z_j|\mathbf{x}, \mathbf{y}) dz_j \prod_{j' \neq j} \int_{z_{j'}} q_\phi(z_{j'}|\mathbf{x}, \mathbf{y}) dz_{j'} \quad (81)$$

$$= \int_{z_j} q_\phi(z_j|\mathbf{x}, \mathbf{y}) \log q_\phi(z_j|\mathbf{x}, \mathbf{y}) dz_j \quad (82)$$

Similarly, we can rewrite the term  $(\heartsuit)$  as:

$$\int_{\mathbf{z}} q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \log p_\theta(z_j|\mathbf{y}) d\mathbf{z} = \int_{z_j} q_\phi(z_j|\mathbf{x}, \mathbf{y}) \log p_\theta(z_j|\mathbf{y}) dz_j \quad (83)$$

Therefore, we have

$$D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \| p_\theta(\mathbf{z}|\mathbf{y})) \quad (84)$$

$$= \underbrace{\sum_j \int_{z_j} q_\phi(z_j|\mathbf{x}, \mathbf{y}) \log q_\phi(z_j|\mathbf{x}, \mathbf{y}) dz_j}_{(1)} - \underbrace{\sum_j \int_{z_j} q_\phi(z_j|\mathbf{x}, \mathbf{y}) \log p_\theta(z_j|\mathbf{y}) dz_j}_{(2)} \quad (85)$$

Now, let us plug in Gaussian distributions. Let's start with the left side:

$$(1) = \sum_j \int_{z_j} q_\phi(z_j|\mathbf{x}, \mathbf{y}) \log q_\phi(z_j|\mathbf{x}, \mathbf{y}) dz_j \quad (86)$$

$$= \sum_j \int_{z_j} q_\phi(z_j|\mathbf{x}, \mathbf{y}) \log \left( \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp \left( -\frac{(z_j - \mu_j)^2}{2\sigma_j^2} \right) \right) dz_j \quad (87)$$

$$= \sum_j \int_{z_j} q_\phi(z_j | \mathbf{x}, \mathbf{y}) \left( \log \frac{1}{\sqrt{2\pi\sigma_j^2}} - \frac{(z_j - \mu_j)^2}{2\sigma_j^2} \right) dz_j \quad (88)$$

$$= \sum_j -\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma_j^2 - \frac{1}{2\sigma_j^2} \int_{z_j} q_\phi(z_j | \mathbf{x}, \mathbf{y}) [(z_j - \mu_j)^2] dz_j \quad (89)$$

$$= \sum_j -\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma_j^2 - \frac{1}{2\sigma_j^2} \sigma_j^2 \quad (90)$$

$$= \sum_j -\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma_j^2 - \frac{1}{2} \quad (91)$$

Now, let us solve the right side. Remember that we assume that  $p_\theta(\mathbf{z}|\mathbf{y}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Therefore:

$$\textcircled{2} = \sum_j \int_{z_j} q_\phi(z_j | \mathbf{x}, \mathbf{y}) \log p_\theta(z_j | \mathbf{y}) dz_j \quad (92)$$

$$= \sum_j \int_{z_j} q_\phi(z_j | \mathbf{x}, \mathbf{y}) \log \left( \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{z_j^2}{2} \right) \right) dz_j \quad (93)$$

$$= \sum_j \int_{z_j} q_\phi(z_j | \mathbf{x}, \mathbf{y}) \left( \log \frac{1}{\sqrt{2\pi}} - \frac{z_j^2}{2} \right) dz_j \quad (94)$$

$$= \sum_j -\frac{1}{2} \log 2\pi - \frac{1}{2} \int_{z_j} q_\phi(z_j | \mathbf{x}, \mathbf{y}) z_j^2 dz_j \quad (95)$$

$$= \sum_j -\frac{1}{2} \log 2\pi - \frac{1}{2} (\sigma_j^2 + \mu_j^2) \quad (96)$$

Therefore, combining the solutions of the left and right side we have:

$$D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) \| p_\theta(\mathbf{z}|\mathbf{y})) = \left( \sum_j -\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma_j^2 - \frac{1}{2} \right) - \left( \sum_j -\frac{1}{2} \log 2\pi - \frac{1}{2} (\sigma_j^2 + \mu_j^2) \right) \quad (97)$$

$$= \sum_j -\frac{1}{2} \log \sigma_j^2 - \frac{1}{2} + \frac{1}{2} (\sigma_j^2 + \mu_j^2) \quad (98)$$

$$= -\frac{1}{2} \sum_j (1 + \log \sigma_j^2 - \sigma_j^2 - \mu_j^2) \quad (99)$$

- (c) [12 points] Fill in code for CVAE network as a `nn.Module` class called `CVAE` in the starter code `cvae.py` and the notebook `cvae.ipynb`:

- Implement the `recognition_model` function  $q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})$ .
- Implement the `generative_model` function  $p_\theta(\mathbf{x}|\mathbf{z}, \mathbf{y})$ .
- Implement the `forward` function by inferring the Gaussian parameters using the recognition model, sampling a latent variable using the reparametrization trick and generating the data using the generative model.
- Implement the variational lowerbound `loss_function`  $\mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{y})$ .
- Train the CVAE and visualize the generated image for each class (i.e., 10 images per class).

- Repeat the image generation 10 times with different random noise. In the write-up, attach and submit  $10 \times 10$  array of images showing all the generated images, where the images in the same row are generated from the same random noise, and images in the same column are generated from the same class label.
- The hyperparameters and training setups provided in the code should work well for learning a CVAE on the MNIST dataset, but please feel free to make any changes as needed and you think appropriate to make CVAE work. Please discuss (if any) there are some notable changes you have made.

If trained successfully, you should be able to sample images  $\mathbf{x}$  that look like MNIST digits reflecting the given label  $\mathbf{y}$ , and the noise vector  $\mathbf{z}$ .

Solution.

Example Output:

