

**1 (Murphy 11.2 - EM for Mixtures of Gaussians)** Show that the M step for ML estimation of a mixture of Gaussians is given by

$$\begin{aligned}\boldsymbol{\mu}_k &= \frac{\sum_i r_{ik} \mathbf{x}_i}{r_k} \\ \boldsymbol{\Sigma}_k &= \frac{1}{r_k} \sum_i r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^\top = \frac{1}{r_k} \sum_i r_{ik} \mathbf{x}_i \mathbf{x}_i^\top - r_k \boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top.\end{aligned}$$

First, we begin with the result of the M step presented in Murphy Equation 11.30 (Section 11.4.2.3):

$$\begin{aligned}l(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) &= -\frac{1}{2} \sum_i r_{ik} [\log |\boldsymbol{\Sigma}_k| + (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)] \\ &= -\frac{1}{2} \sum_i r_{ik} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} \sum_i r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k).\end{aligned}$$

Next, we will find the optimal values for  $\boldsymbol{\mu}_k$  and  $\boldsymbol{\Sigma}_k$  by finding their critical points. Taking the derivative of  $l(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  with respect to  $\boldsymbol{\mu}_k$ , we have:

$$\begin{aligned}\frac{\partial l(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} &= \frac{\partial}{\partial \boldsymbol{\mu}_k} \left( -\frac{1}{2} \sum_i r_{ik} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} \sum_i r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right) \\ &= \frac{\partial}{\partial \boldsymbol{\mu}_k} \left( -\frac{1}{2} \sum_i r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right) \\ &= \frac{\partial}{\partial U} \left( -\frac{1}{2} \sum_i r_{ik} U^\top \boldsymbol{\Sigma}_k^{-1} U \right) \frac{dU}{d\boldsymbol{\mu}_k}, \text{ where } U = (\mathbf{x}_i - \boldsymbol{\mu}_k), \\ &= \left( -\frac{1}{2} \sum_i r_{ik} \boldsymbol{\Sigma}_k^{-1} U - \frac{1}{2} \sum_i r_{ik} \boldsymbol{\Sigma}_k^{-1} U \right) (-1), \text{ since } \frac{\partial U^\top x}{\partial U} = \frac{\partial x^\top U}{\partial U} = x, \\ &= \sum_i r_{ik} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k).\end{aligned}$$

Now we must find the critical point for  $\boldsymbol{\mu}_k$ . Setting the derivative of  $l(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  with respect

to  $\mu_k$  to 0 and solving for  $\mu_k$ , we have:

$$\begin{aligned}
\sum_i r_{ik} \Sigma_k^{-1} (\mathbf{x}_i - \mu_k) &= 0 \\
\sum_i r_{ik} \Sigma_k^{-1} \mathbf{x}_i &= \sum_i r_{ik} \Sigma_k^{-1} \mu_k \\
\sum_i r_{ik} \mathbf{x}_i &= \sum_i r_{ik} \mu_k, \text{ since the covariance matrix is a linear operator} \\
\sum_i r_{ik} \mathbf{x}_i &= \mu_k \sum_i r_{ik}, \text{ since } \mu_k \text{ does not depend on } i, \\
\sum_i r_{ik} \mathbf{x}_i &= \mu_k r_k, \text{ so} \\
\mu_k &= \frac{\sum_i r_{ik} \mathbf{x}_i}{r_k}.
\end{aligned}$$

We see that this yields the value of  $\mu_k$  presented in Murphy, Equation 11.31.

Next, we will take the derivative of  $l(\mu_k, \Sigma_k)$  with respect to  $\Sigma_k$ :

$$\begin{aligned}
\frac{\partial l(\mu_k, \Sigma_k)}{\partial \Sigma_k} &= \frac{\partial}{\partial \Sigma_k} \left( -\frac{1}{2} \sum_i r_{ik} \log |\Sigma_k| - \frac{1}{2} \sum_i r_{ik} (\mathbf{x}_i - \mu_k)^\top \Sigma_k^{-1} (\mathbf{x}_i - \mu_k) \right) \\
&= \frac{\partial}{\partial \Sigma_k} \left( -\frac{1}{2} \sum_i r_{ik} \log |\Sigma_k| \right) - \frac{\partial}{\partial \Sigma_k} \left( \frac{1}{2} \sum_i r_{ik} (\mathbf{x}_i - \mu_k)^\top \Sigma_k^{-1} (\mathbf{x}_i - \mu_k) \right) \\
&= -\frac{1}{2} \sum_i r_{ik} \Sigma_k^{-1} - \frac{\partial}{\partial \Sigma_k} \left( \frac{1}{2} \sum_i r_{ik} (\mathbf{x}_i - \mu_k)^\top \Sigma_k^{-1} (\mathbf{x}_i - \mu_k) \right) \\
&= -\frac{1}{2} \sum_i r_{ik} \Sigma_k^{-1} - \frac{1}{2} \sum_i r_{ik} (-\Sigma_k^{-\top}) (\mathbf{x}_i - \mu_k) (\mathbf{x}_i - \mu_k)^\top (\Sigma_k^{-\top}) \\
&\quad \text{since } \frac{\partial \mathbf{a}^\top \mathbf{X}^{-1} \mathbf{b}}{\partial \mathbf{X}} = -\mathbf{X}^{-\top} \mathbf{a} \mathbf{b}^\top \mathbf{X}^{-\top}.
\end{aligned}$$

Now we must find the critical point for  $\Sigma_k$ . Setting the derivative of  $l(\mu_k, \Sigma_k)$  with respect

to  $\Sigma_k$  to 0 and solving for  $\Sigma_k$ , we have:

$$\begin{aligned}
0 &= -\frac{1}{2} \sum_i r_{ik} \Sigma_k^{-1} - \frac{1}{2} \sum_i r_{ik} (-\Sigma_k^{-\top})(\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^\top (\Sigma_k^{-\top}) \\
-\frac{1}{2} \sum_i r_{ik} \Sigma_k^{-1} &= \frac{1}{2} \sum_i r_{ik} (-\Sigma_k^{-\top})(\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^\top (\Sigma_k^{-\top}) \\
r_k \Sigma_k^{-1} &= \sum_i r_{ik} (\Sigma_k^{-\top})(\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^\top (\Sigma_k^{-\top}) \\
r_k \Sigma_k^{-1} \Sigma_k &= \sum_i r_{ik} \Sigma_k^{-\top} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \Sigma_k^{-\top} \Sigma_k \\
r_k &= \sum_i r_{ik} \Sigma_k^{-\top} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \\
\Sigma_k r_k &= \sum_i r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^\top, \\
&\text{since } \Sigma_k \text{ is a linear operator, so} \\
\Sigma_k &= \frac{1}{r_k} \sum_i r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^\top.
\end{aligned}$$

We see that this yields the value of  $\Sigma_k$  presented in Murphy, Equation 11.32. ■

**2 (SVD Image Compression)** In this problem, we will use the image of a scary clown online to perform image compression. In the starter code, we have already load the image into a matrix/array for you. However, you might need internet connection to access the image and therefore successfully run the starter code. The code requires Python library Pillow in order to run.

Plot the progression of the 100 largest singular values for the original image and a randomly shuffled version of the same image (all on the same plot). In a single figure plot a grid of four images: the original image, and a rank  $k$  truncated SVD approximation of the original image for  $k \in \{2, 10, 20\}$ .

■