

Feel free to work with other students, but make sure you write up the homework and code on your own (no copying homework *or* code; no pair programming). Feel free to ask students or instructors for help debugging code or whatever else, though.

The starter files for problem 2 can be found under the Resource tab on course website. Please print out all the graphs generated by your own code and submit them together with the written part, and make sure you upload the code to your Github repository.

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

$$\mu_k = \frac{\sum_i r_{ik} \mathbf{x}_i}{r_k}$$
$$\Sigma_k = \frac{1}{r_k} \sum_i r_{ik} (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^\top = \frac{1}{r_k} \sum_i r_{ik} \mathbf{x}_i \mathbf{x}_i^\top - r_k \mu_k \mu_k^\top.$$

We have the complete data log likelihood $l(\mu_k, \Sigma_k) = \sum_k \sum_i r_{ik} \log \mathbb{P}(\mathbf{x}_i | \theta_k)$
 $= \frac{-1}{2} \sum_i r_{ik} (\log |\Sigma_k| + (\mathbf{x}_i - \mu_k)^\top \Sigma_k^{-1} (\mathbf{x}_i - \mu_k))$ to a constant.

Differentiate with respect to μ_k . The $\log |\Sigma_k|$ term becomes a constant that gets eliminated. The Σ_k^{-1} term becomes a constant that remains due to similar operations as the fact that for example, $\frac{\partial}{\partial x} x^\top A x = 2Ax$.

$$\frac{\partial l}{\partial \mu_k} = \sum_i r_{ik} \Sigma_k^{-1} (\mathbf{x}_i - \mu_k) = \Sigma_k^{-1} \sum_i r_{ik} (\mathbf{x}_i - \mu_k) = 0 \text{ since } \Sigma_k^{-1} \text{ is linear.}$$

Simplify in order to achieve optimization.

First, Σ_k^{-1} can be eliminated as a constant.

Then, $\sum_i r_{ik} \mathbf{x}_i = \sum_i r_{ik} \mu_k$.

At optimality, this becomes $\sum_i r_{ik} \mathbf{x}_i = \mu_k \sum_i r_{ik}$.

Isolate for μ_k , so $\mu_k = \frac{\sum_i r_{ik} \mathbf{x}_i}{r_k}$ since $\sum_i r_{ik} = r_k$, as desired.

Next, differentiate with respect to Σ_k , so

$$\frac{\partial l}{\partial \Sigma_k} = \frac{-1}{2} \sum_i r_{ik} (\Sigma_k^{-1} - \Sigma_k^{-1} (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^\top \Sigma_k^{-1}) = 0.$$

We did this by using the fact that $\frac{\partial^{-1}}{\partial x} = -Y^{-1} \frac{\partial Y}{\partial x} Y^{-1}$. In this case, Y^{-1} is Σ_k^{-1} , and x is Σ_k .

This gives the optimality condition that $\sum_i r_{ik} I = (\sum_i r_{ik} (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^\top) \Sigma_k^{-1}$.

Multiply by Σ_k on the right and divide by $r_k = \sum_i r_{ik}$ to get

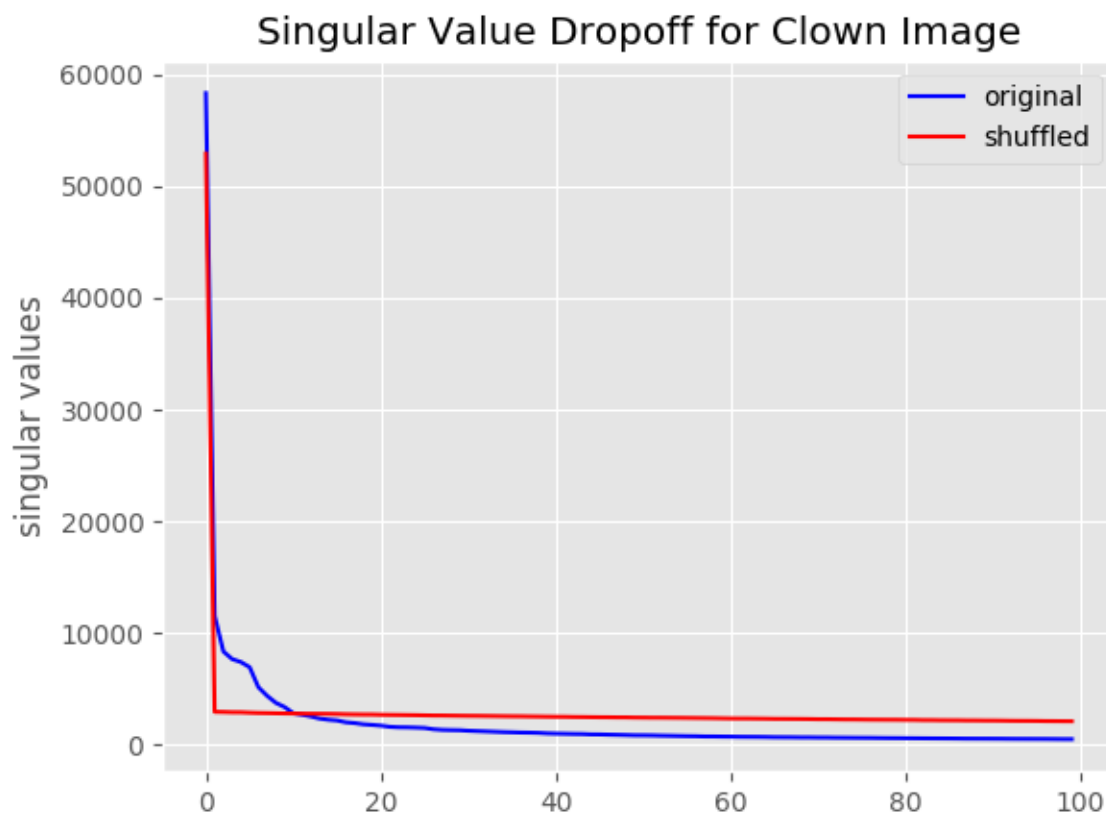
$$\frac{\Sigma_{i r_{ik} I}}{\Sigma_{i r_{ik}}} = (\Sigma_i r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T) \Sigma_k^{-1} \Sigma_k, \text{ which becomes}$$

$$\Sigma_k = (\Sigma_i r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T), \text{ as desired.}$$

■

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 $k \in \{2, 10, 20\}$.



The reconstruction plot of the original image and rank 2, 10, and 20 truncated stochastic gradient descent approximations of the original image are below:

Original Image



Rank 2 Approximation



Rank 10 Approximation



Rank 20 Approximation

