

ISYE 6740, Summer 2021, Homework 3

100 points

Prof. Yao Xie

1. Density estimation: Psychological experiments (50 points).

In Kanai, R., Feilden, T., Firth, C. and Rees, G., 2011. *Political orientations are correlated with brain structure in young adults. Current biology, 21(8), pp.677-680.*, data are collected to study whether or not the two brain regions are likely to be independent of each other and considering different types of political view **For this question; you can use the proper package for histogram and KDE; no need to write your own.** The data set n90pol.csv contains information on 90 university students who participated in a psychological experiment designed to look for relationships between the size of different regions of the brain and political views. The variables `amygdala` and `acc` indicate the volume of two particular brain regions known to be involved in emotions and decision-making, the amygdala and the anterior cingulate cortex; more exactly, these are residuals from the predicted volume, after adjusting for height, sex, and similar body-type variables. The variable `orientation` gives the students' locations on a five-point scale from 1 (very conservative) to 5 (very liberal). Note that in the dataset, we only have observations for orientation from 2 to 5.

Recall in this case, the kernel density estimator (KDE) for a density is given by

$$p(x) = \frac{1}{m} \sum_{i=1}^m \frac{1}{h} K\left(\frac{x^i - x}{h}\right),$$

where x^i are two-dimensional vectors, $h > 0$ is the kernel bandwidth, based on the criterion we discussed in lecture. For one-dimensional KDE, use a one-dimensional Gaussian kernel

$$K(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}.$$

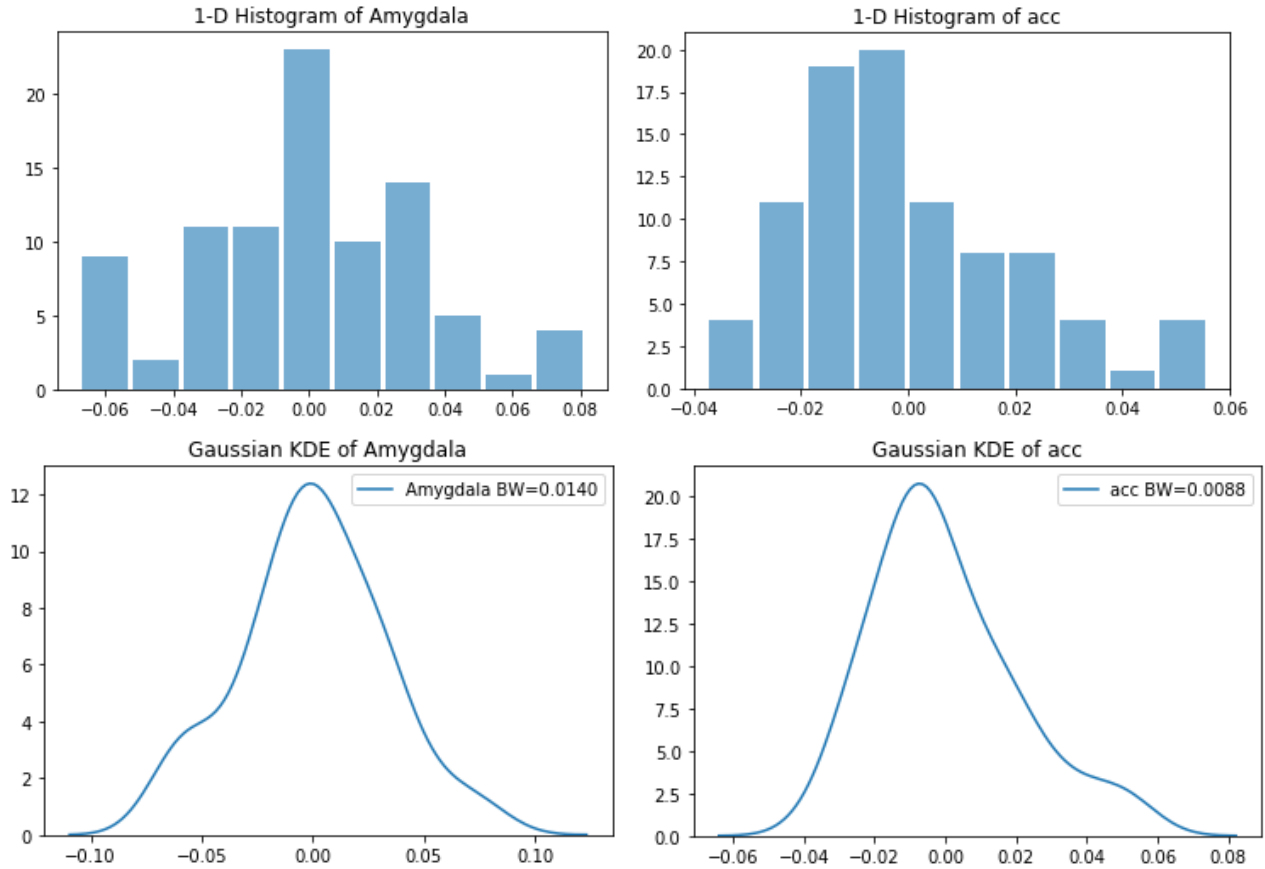
For two-dimensional KDE, use a two-dimensional Gaussian kernel: for

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \in \mathbb{R}^2,$$

where x_1 and x_2 are the two dimensions respectively

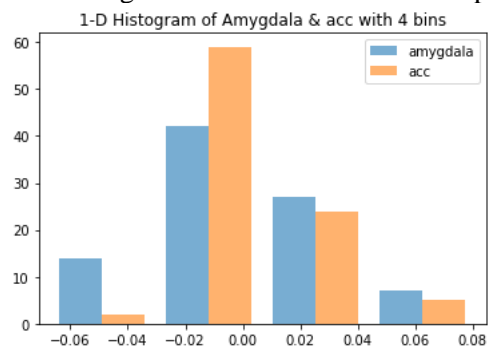
$$K(x) = \frac{1}{2\pi} e^{-\frac{(x_1)^2 + (x_2)^2}{2}}.$$

- (a) (10 points) Form the 1-dimensional histogram and KDE to estimate the distributions of amygdala and acc, respectively. For this question, you can ignore the variable orientation. Decide on a suitable number of bins so you can see the shape of the distribution clearly. Set an appropriate kernel bandwidth $h > 0$.

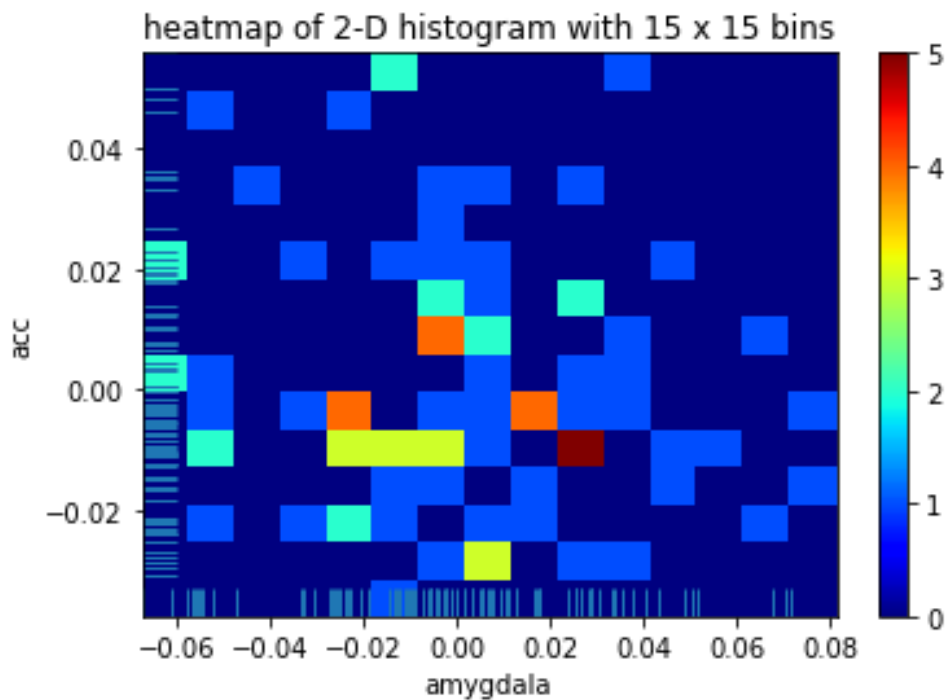


- nbins of histogram = 10
- The kernel bandwidths of Amygdala and acc are based on Silverman's rule of thumb $h = 1.6 * \sigma * m^{(0.2)}$

1-D Histogram of the variables in same plot for comparison:



- (b) (10 points) Form 2-dimensional histogram for the pairs of variables (`amygdala`, `acc`).
Decide on a suitable number of bins so you can see the shape of the distribution clearly.



We can see the concentration of the data points as shown in the heat map and followed by the rug plot on the axes.

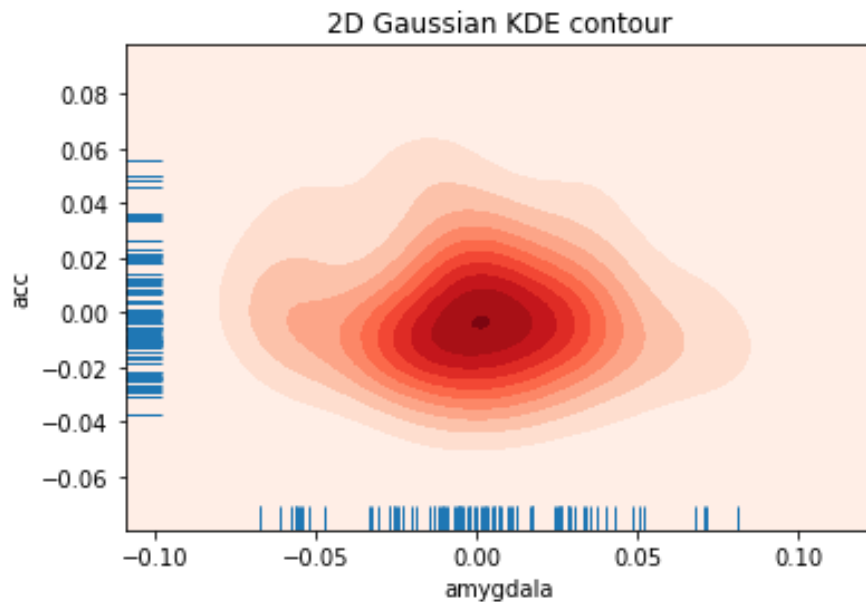
The data is crowded around Amygdala (-0.02,0) and acc (0,-0.01)

- (c) (10 points) Use kernel-density-estimation (KDE) to estimate the 2-dimensional density function of (**amygdala**, **acc**) (this means for this question, you can ignore the variable orientation). Set an appropriate kernel bandwidth $h > 0$.

Please show the two-dimensional KDE (e.g., two-dimensional heat-map, two-dimensional contour plot, etc.)

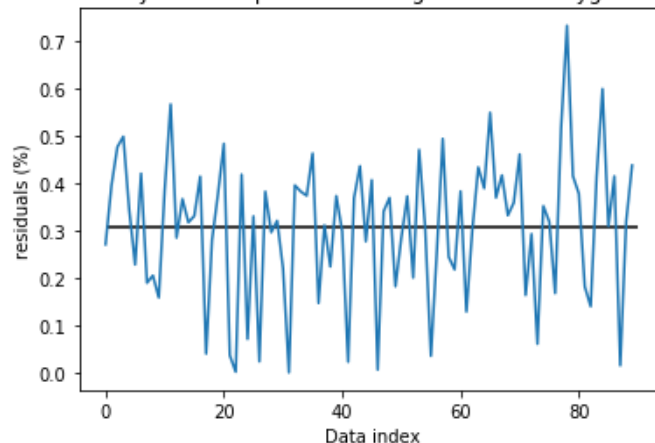
Please explain based on the results, can you infer that the two variables (**amygdala**, **acc**) are likely to be independent or not?

- Kernel bandwidth $h = 0.01397$ is used.



To check for independence, I would like to check the difference (residuals) between the joint probabilities of the variables and the product of their marginal probabilities. Below is the %residuals for each of the data points and on an avg there is a **30%** variability. This implies that the variables are not totally independent. There is some level of dependance between them.

Residuals of Joint dist - product of marginal dist of amygdala and acc



- (d) (10 points) We will consider the variable **orientation** and consider conditional distributions. Please plot the estimated conditional distribution of **amygdala** conditioning on political orientation: $p(\text{amygdala} \mid \text{orientation} = c)$, $c = 2, \dots, 5$, using KDE. Set an appropriate kernel bandwidth $h > 0$. Do the same for the volume of the **acc**: plot $p(\text{acc} \mid \text{orientation} = c)$, $c = 2, \dots, 5$ using KDE. (Note that the conditional distribution can be understood as fitting a distribution for the data with the same **orientation**. Thus you should plot 8 one-dimensional distribution functions in total for this question.)

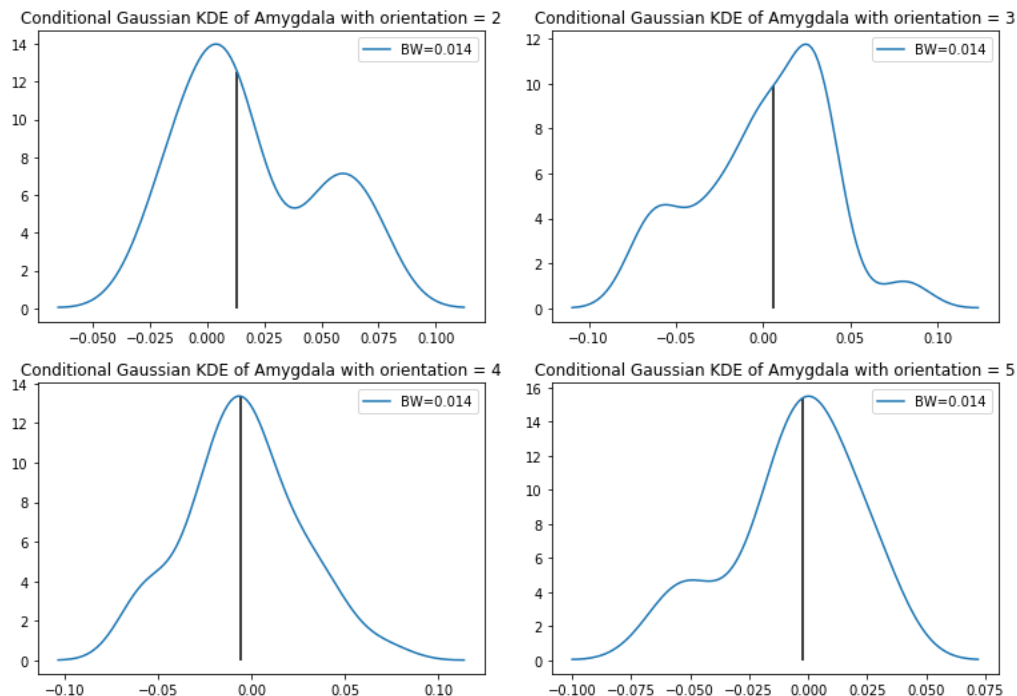
Now please explain based on the results, can you infer that the conditional distribution of **amygdala** and **acc**, respectively, are different from $c = 2, \dots, 5$? This is a type of scientific question one could infer from the data: Whether or not there is a difference between brain structure and political view.

Now please also fill out the *conditional sample mean* for the two variables:

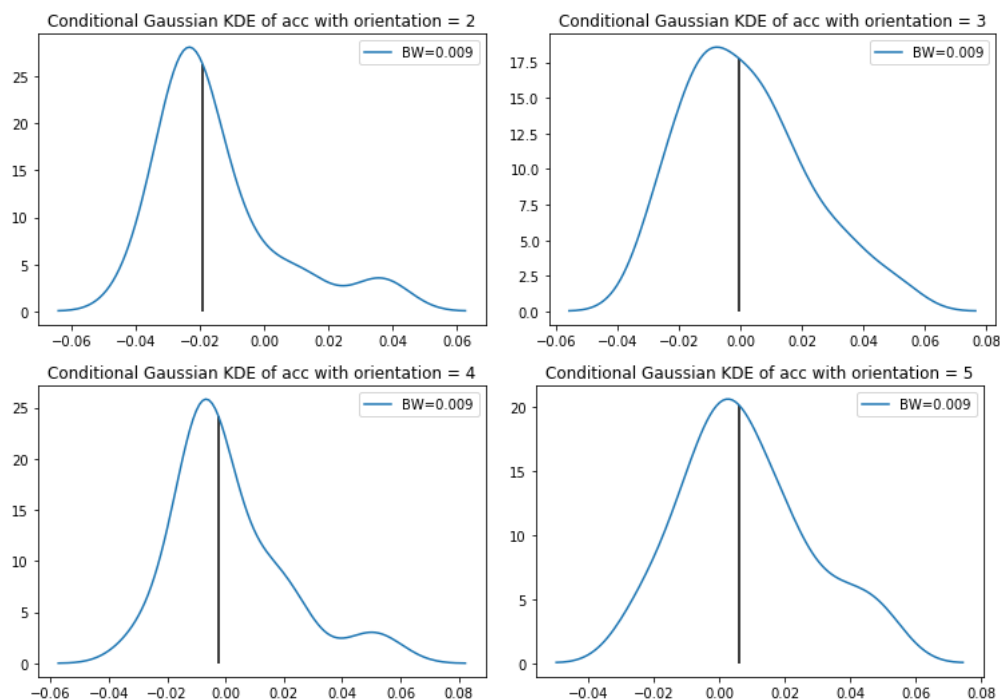
	$c = 2$	$c = 3$	$c = 4$	$c = 5$
amygdala				
acc				

Remark: As you can see this exercise, you can extract so much more information from density estimation than simple summary statistics (e.g., the sample mean) in terms of explorable data analysis.

- Amygdala responses are varying quite strongly for people with different political orientations.
- Specifically the extreme liberal (5) and extreme conservative(2) orientations have strong left inclined and right inclined Gaussian KDEs respectively



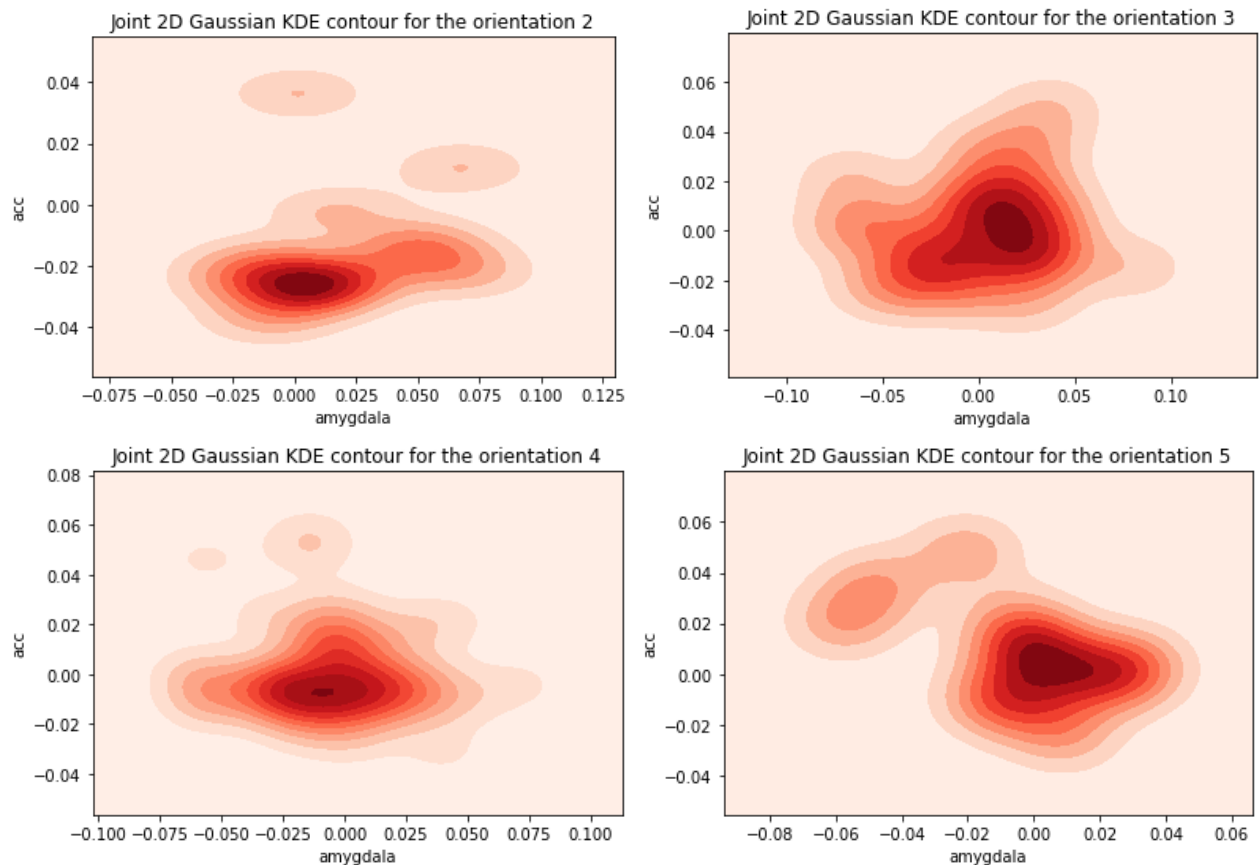
- Acc responses are not varying a lot for people with different political orientation. We do not see as strong orientations as in the amygdala area.



Mean	$c = 2$	$c = 3$	$c = 4$	$c = 5$
amygdala	0.0190	0.0005875	-0.004719	-0.00569
acc	-0.01477	0.00167	0.0013	0.00814

- (e) (10 points) Again we will consider the variable **orientation**. We will estimate the conditional *joint* distribution of the volume of the **amygdala** and **acc**, conditioning on a function of political orientation: $p(\text{amygdala}, \text{acc} \mid \text{orientation} = c)$, $c = 2, \dots, 5$. You will use two-dimensional KDE to achieve the goal; et an appropriate kernel bandwidth $h > 0$. Please show the two-dimensional KDE (e.g., two-dimensional heat-map, two-dimensional contour plot, etc.).

Please explain based on the results, can you infer that the conditional distribution of two variables (**amygdala**, **acc**) are different from $c = 2, \dots, 5$? This is a type of scientific question one could infer from the data: Whether or not there is a difference between brain structure and political view.



The brain structure political orientations of scale 3 and 4 are close, but the orientation 2 and orientation 5 differ quite a lot in the shape of the conditional KDE distribution this highlights that the brain regions amygdala and acc are correlated with the political orientation.

2. Implementing EM for MNIST dataset (50 points).

Implement the EM algorithm for fitting a Gaussian mixture model for the MNIST hand-written digits dataset. For this question, we reduce the dataset to be only two cases, of digits “2” and “6” only. Thus, you will fit GMM with $C = 2$. Use the data file `data.mat` or `data.dat`. True label of the data are also provided in `label.mat` and `label.dat`.

The matrix `images` is of size 784-by-1990, i.e., there are totally 1990 images, and each column of the matrix corresponds to one image of size 28-by-28 pixels (the image is vectorized; the original image can be recovered by map the vector into a matrix).

First use PCA to reduce the dimensionality of the data before applying to EM. We will put all “6” and “2” digits together, to project the original data into 4-dimensional vectors.

Now implement EM algorithm for the projected data (with 4-dimensions).

- (b) (10 points) Write down detailed expression of the E-step and M-step in the EM algorithm (hint: when computing γ_k^i , you can drop the $(2\pi)^{n/2}$ factor from the numerator and denominator expression, since it will be canceled out; this can help avoid some numerical issues in computation).

Expectation Step:-

$\gamma_k^i \rightarrow$ the vector that has the probabilities of i^{th} data in k^{th} component. \therefore From Bayes rule

$$\gamma_k^i = p(z_k^i = 1 | D, \mu, \Sigma) = \frac{\pi_k N(x^i | \mu_k, \Sigma_k)}{\sum_{k'=1}^K \pi_{k'} N(x^i | \mu_{k'}, \Sigma_{k'})}$$

We expand this unimodal gaussian density function

$$\Rightarrow N(x^i | \mu_k, \Sigma_k) = \frac{1}{|\Sigma_k|^{1/2} (2\pi)^{n/2}} e^{\left(-\frac{1}{2} (x^i - \mu_k)^T \Sigma_k^{-1} (x^i - \mu_k)\right)}$$

This can be substituted in the equation. $(2\pi)^{n/2}$ gets cancelled in the process.

$$\therefore \gamma_k^i = \frac{\pi_k \left[\frac{1}{|\Sigma_k|^{1/2}} e^{\left(-\frac{1}{2} (x^i - \mu_k)^T \Sigma_k^{-1} (x^i - \mu_k)\right)} \right]}{\sum_{k'=1}^K \pi_{k'} \left[\frac{1}{|\Sigma_{k'}|^{1/2}} e^{\left(-\frac{1}{2} (x^i - \mu_{k'})^T \Sigma_{k'}^{-1} (x^i - \mu_{k'})\right)} \right]}$$

for $k=1, \dots, K, i=1, \dots, m$

Maximization step:-

Update (π_k, μ_k, Σ_k) with γ_k^i . $\pi_k = \frac{\sum_i \gamma_k^i}{m}$

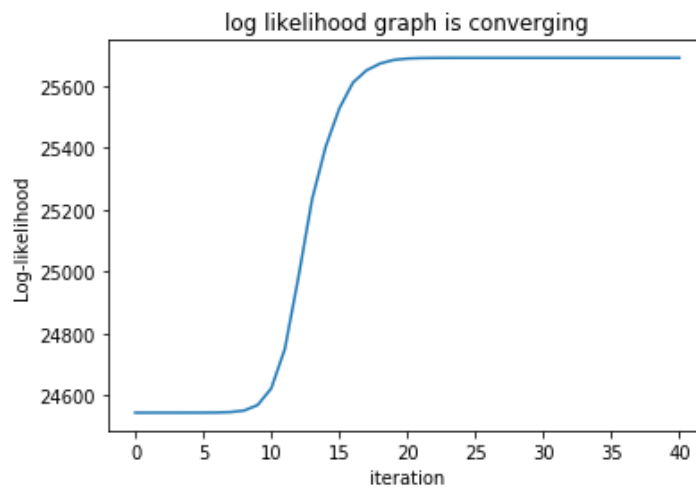
$$\mu_k = \frac{\sum_i \gamma_k^i x^i}{\sum_i \gamma_k^i} \quad \Sigma_k = \frac{\sum_i \gamma_k^i (x^i - \mu_k)(x^i - \mu_k)^T}{\sum_i \gamma_k^i}$$

$(k=1, \dots, K, i=1, \dots, m)$

(c) (16 points) Implement EM algorithm yourself. Use the following initialization

- initialization for mean: random Gaussian vector with zero mean
- initialization for covariance: generate two Gaussian random matrix of size n -by- n : S_1 and S_2 , and initialize the covariance matrix for the two components are $\Sigma_1 = S_1 S_1^T + I_n$, and $\Sigma_2 = S_2 S_2^T + I_n$, where I_n is an identity matrix of size n -by- n .

Plot the log-likelihood function versus the number of iterations to show your algorithm is converging.

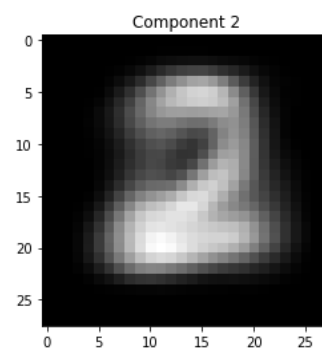
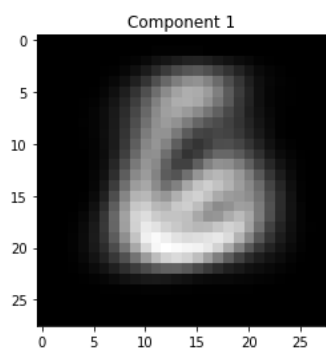


Given a threshold of 10^{-9} the Log-likelihood graph converges pretty soon, after about 20 iterations of E-M steps

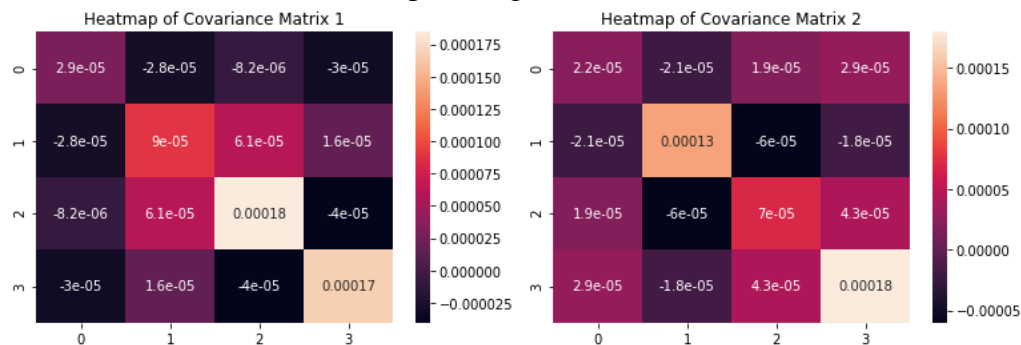
- (d) (12 points) Report, the fitted GMM model when EM has terminated in your algorithms as follows. Report the weights for each component, and the mean of each component, by mapping them back to the original space and reformat the vector to make them into 28-by-28 matrices and show images. Ideally, you should be able to see these means corresponds to some kind of “average” images. You can report the two 4-by-4 covariance matrices by visualizing their intensities (e.g., using a gray scaled image or heat map).

Mixture Model of two gaussian models (GMM):

- Weights of each component is [0.51309328 0.48690672]
- GMM Model with Mean:
 - $\mu_1 = [0.00739311 \ 0.00312245 \ -0.0006553 \ 0.00018435]$
 - $\mu_2 = [-0.00779072 \ -0.00329038 \ 0.00069055 \ -0.00019426]$
- The images corresponding to means of each of the gaussian surfaces are below



- The Covariance matrices of the respective gaussian surfaces are below:



- (e) (12 points) Use the τ_k^i to infer the labels of the images, and compare with the true labels. Report the mis-classification rate for digits “2” and “6” respectively. Perform K -means clustering with $K = 2$ (you may call a package or use the code from your previous homework). Find out the mis-classification rate for digits “2” and “6” respectively, and compare with GMM. Which one achieves the better performance?

• **GMM:**

- Misclassification of 2's as 6's = $(67/(67+965)*100) = \sim 6.49\%$
- Misclassification of 6's as 2's = $(9/(9+949)*100) = \sim 0.94\%$

GMM	Predicted_2	Predicted_6
True_2	965	67
True_6	9	949

• **Kmeans:**

- Misclassification of 2's as 6's = $(65/(65+967)*100) = \sim 6.30\%$
- Misclassification of 6's as 2's = $(82/(82+876)*100) = \sim 8.56\%$

K-means	Predicted_2	Predicted_6
True_2	967	65
True_6	82	876

While the classification of the digit “2” is almost similar by both methods, GMM Achieves much better performance than K-means in classifying the digits “6” as shown by the misclassification rates above.