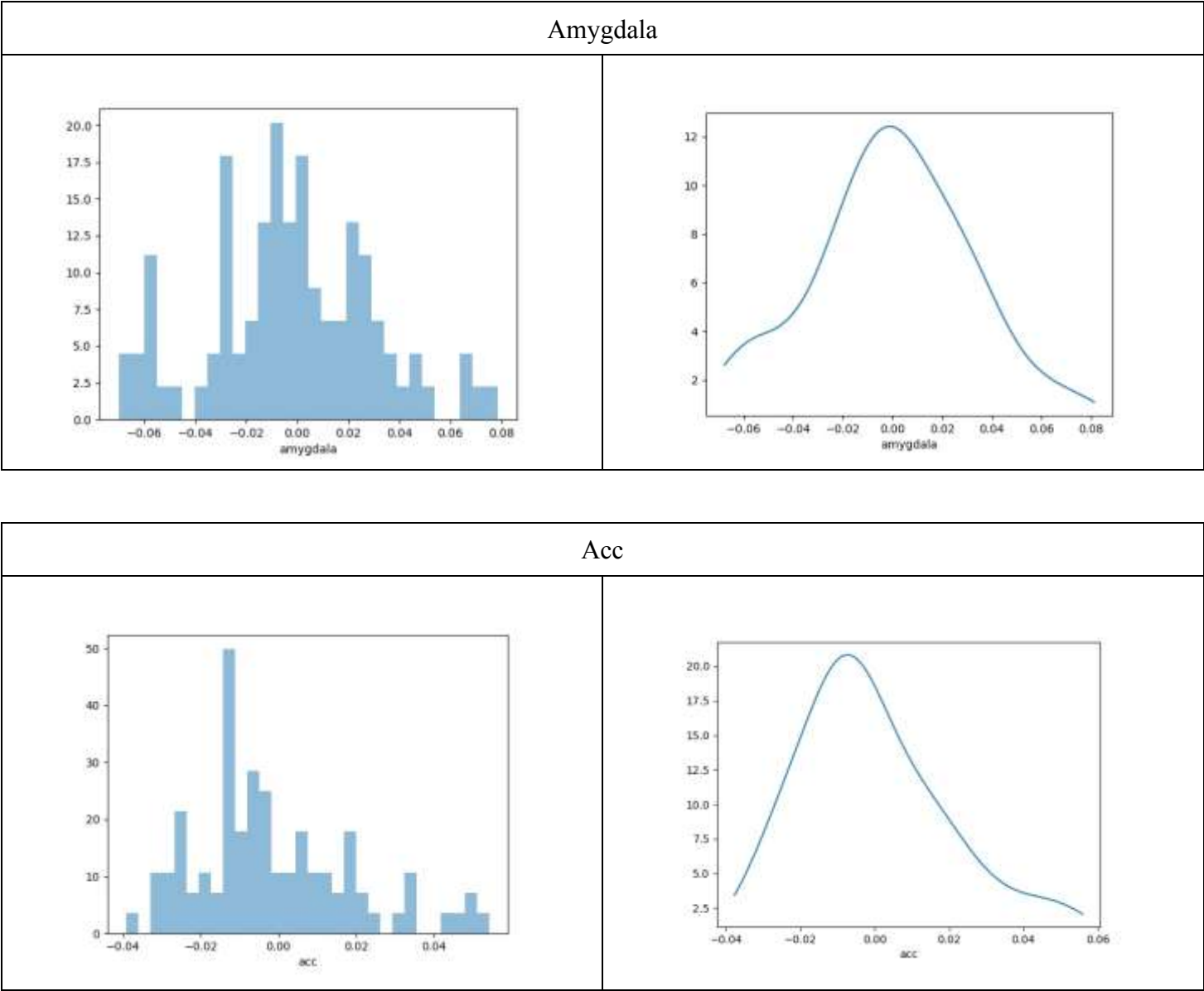


1.Density estimation

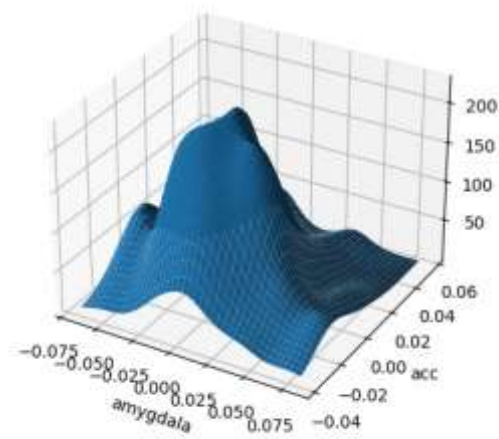
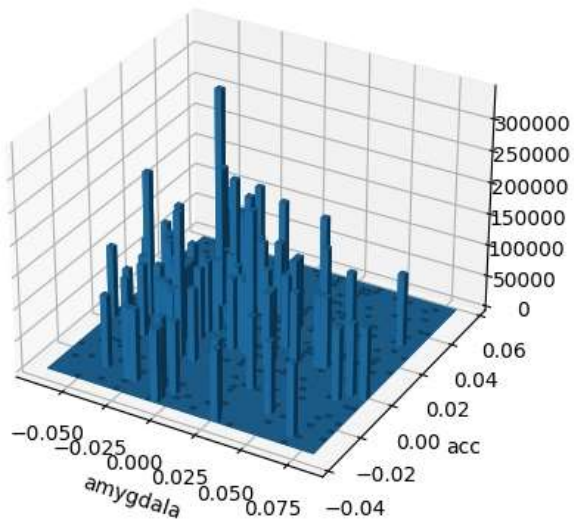
(a)

I use the rules of thumb to decide the bandwidth for KDE. The graph is shown as below:



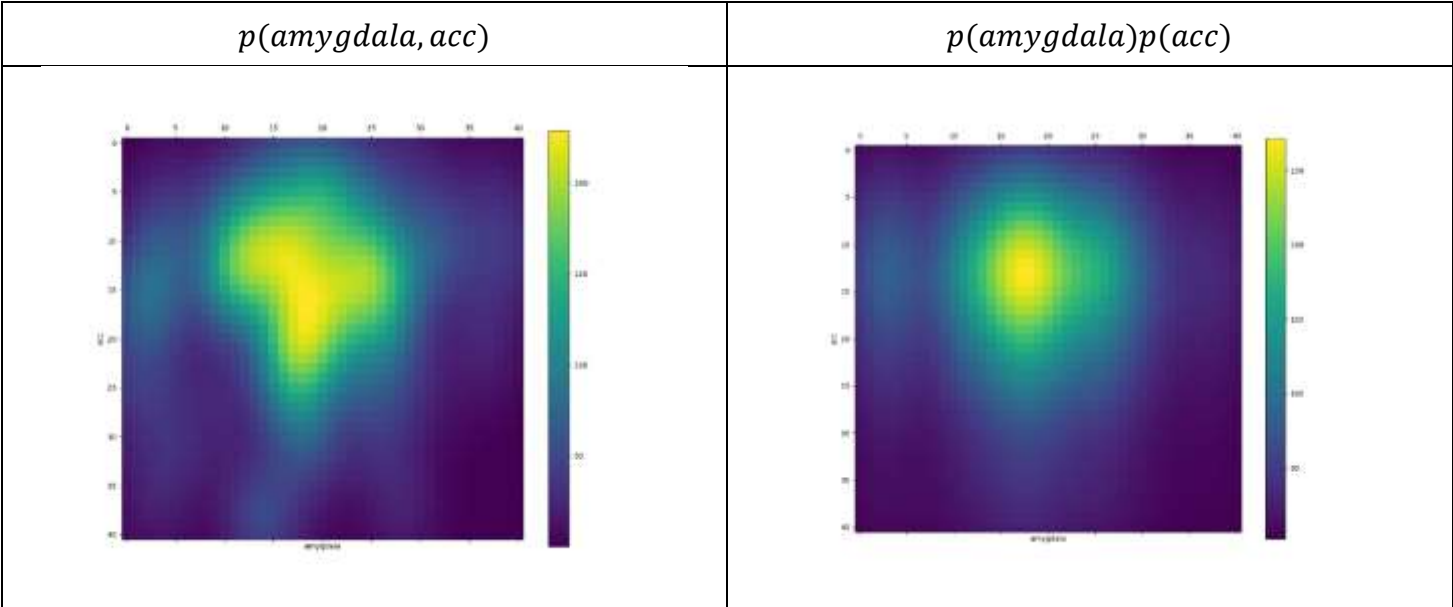
(b)

I set bandwidth=0.01 for KDE. The graph is shown as below:

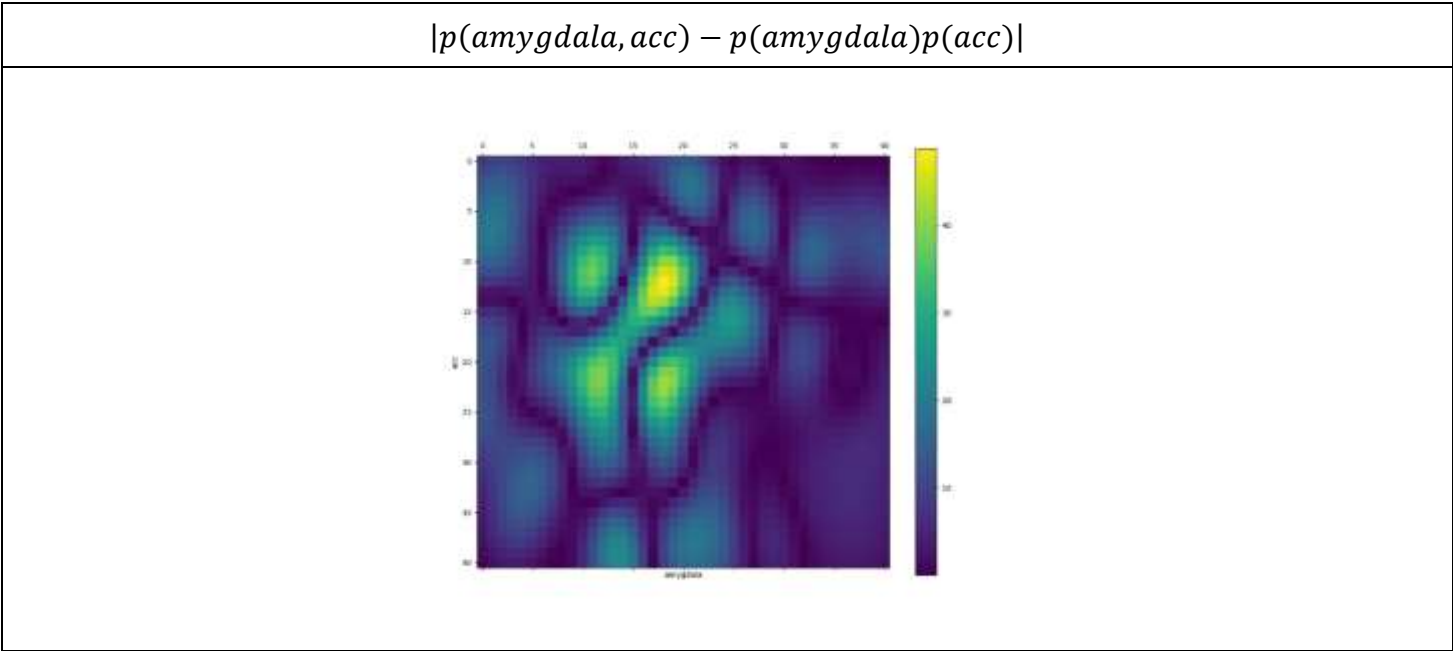


(c)

I set bandwidth=0.01 both for joint distribution and marginal distribution in KDE. The graph is shown as below:

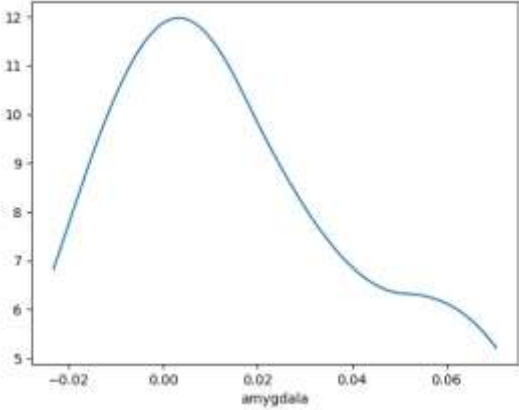
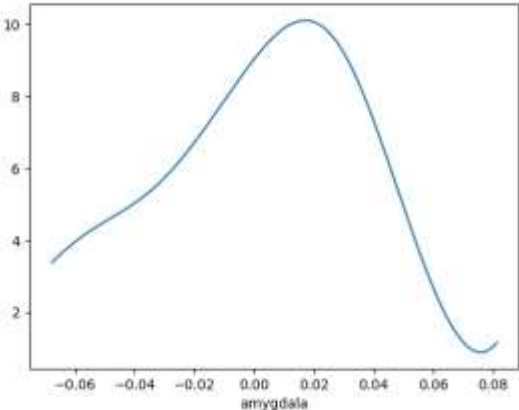
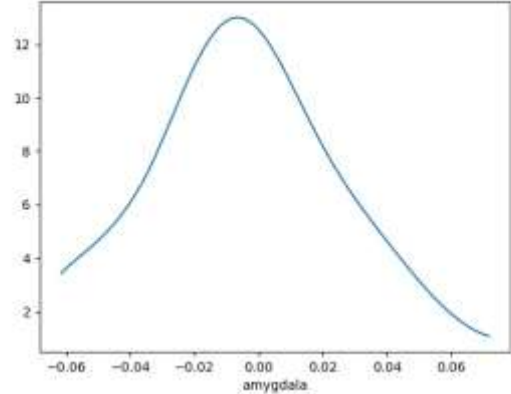


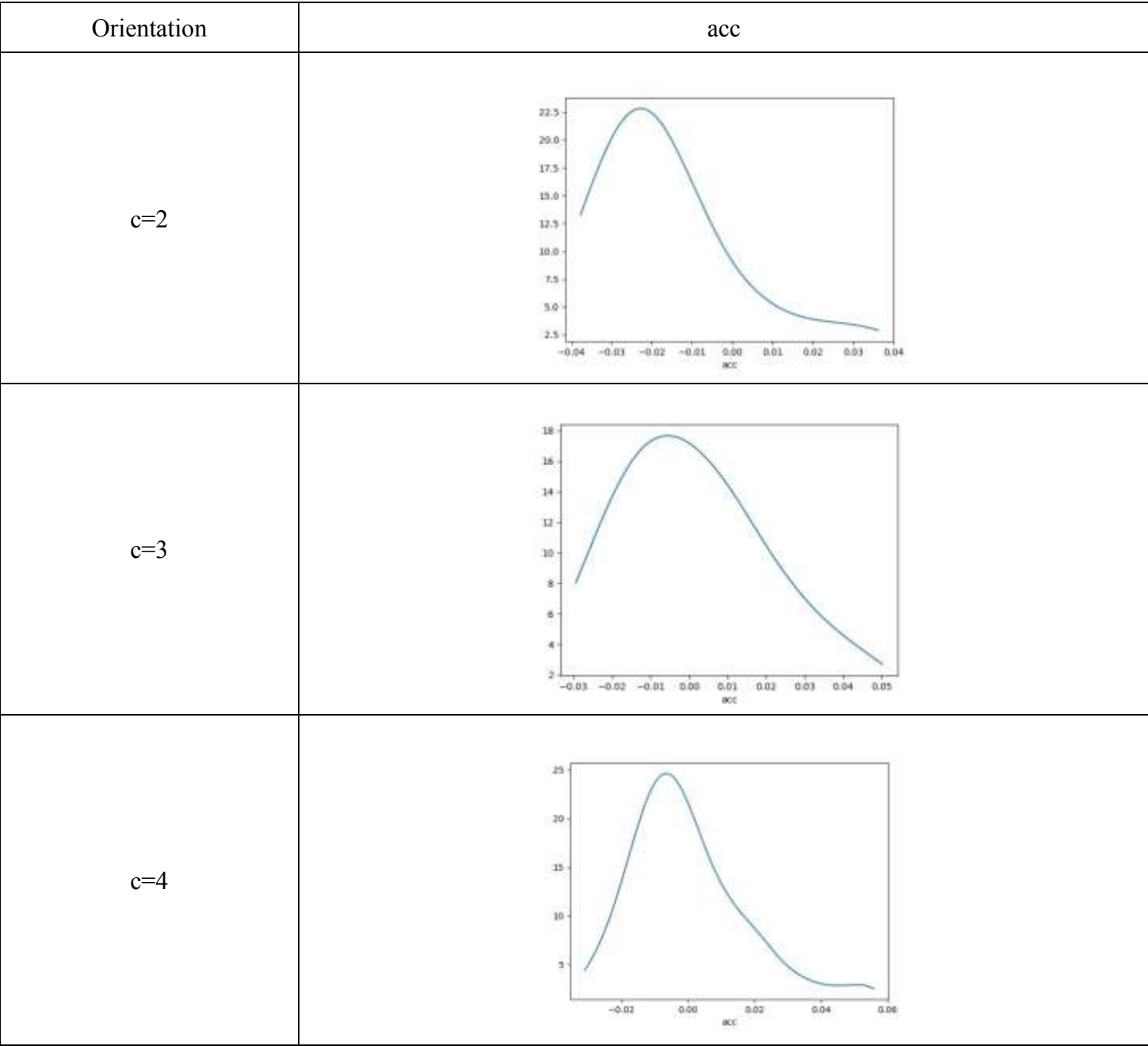
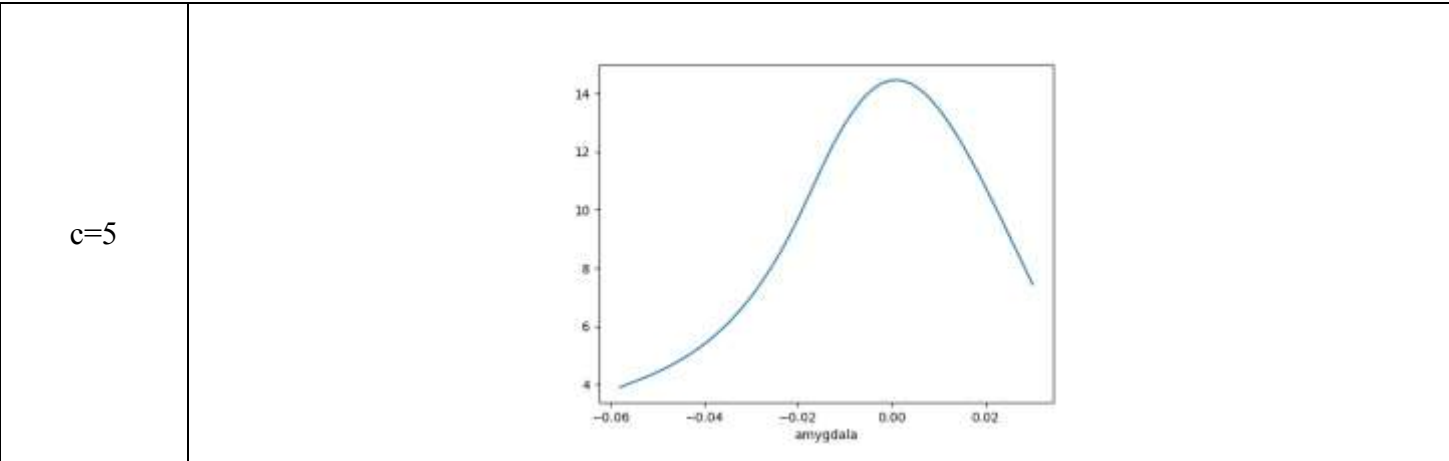
The difference for these two graphs is:

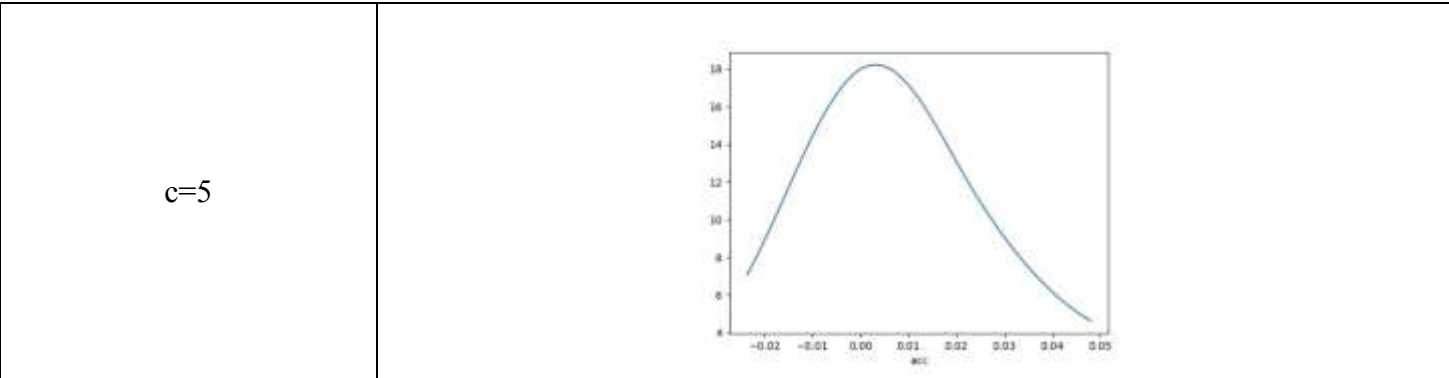


As you can see, the difference between $p(amygdala, acc)$ and $p(amygdala)p(acc)$ roughly doesn't equal to 0, which means that the two parts of brains are related with each other in this area..

(d)

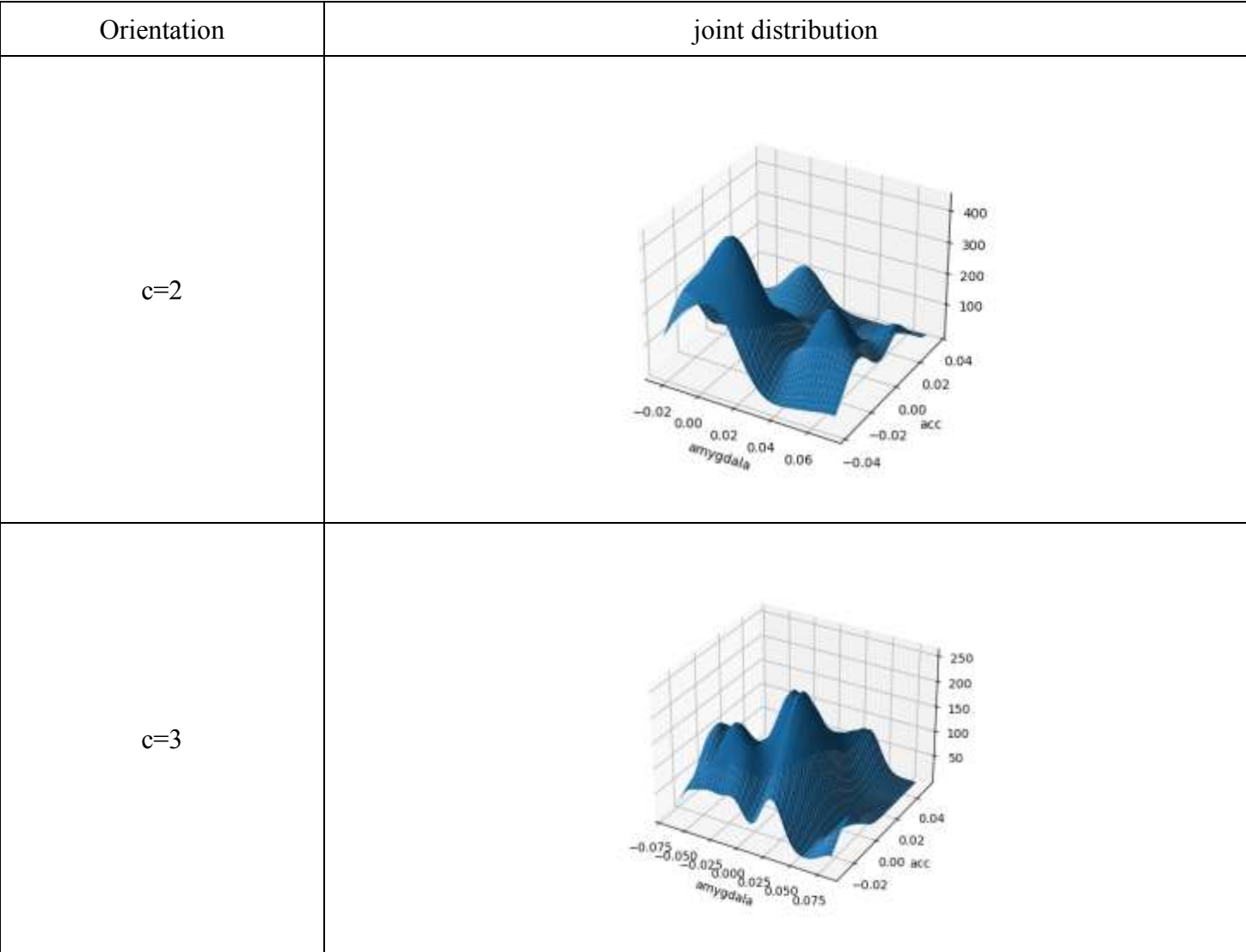
Orientation	amygdala
c=2	 <p>A line graph showing a distribution curve. The x-axis is labeled 'amygdala' and ranges from -0.02 to 0.06 with major ticks every 0.02. The y-axis ranges from 5 to 12 with major ticks every 1 unit. The curve starts at approximately (-0.02, 6.8), rises to a peak of 12 at x ≈ 0.005, and then descends to approximately (0.07, 5.2).</p>
c=3	 <p>A line graph showing a distribution curve. The x-axis is labeled 'amygdala' and ranges from -0.06 to 0.08 with major ticks every 0.02. The y-axis ranges from 2 to 10 with major ticks every 2 units. The curve starts at approximately (-0.06, 3.5), rises to a peak of 10 at x ≈ 0.015, and then descends to approximately (0.075, 1.2).</p>
c=4	 <p>A line graph showing a distribution curve. The x-axis is labeled 'amygdala' and ranges from -0.06 to 0.06 with major ticks every 0.02. The y-axis ranges from 2 to 12 with major ticks every 2 units. The curve starts at approximately (-0.06, 3.5), rises to a peak of 12 at x ≈ -0.005, and then descends to approximately (0.07, 1.2).</p>

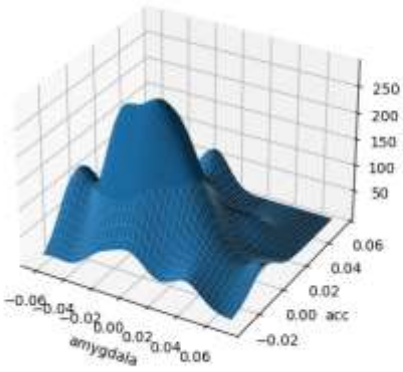
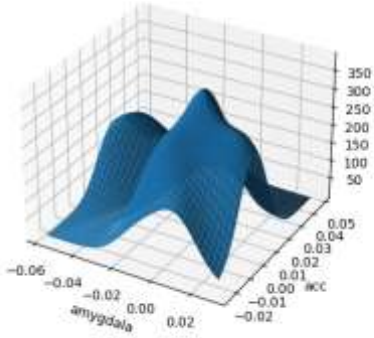




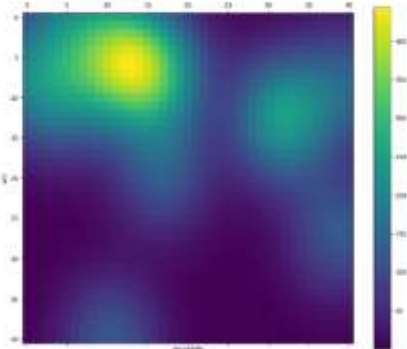
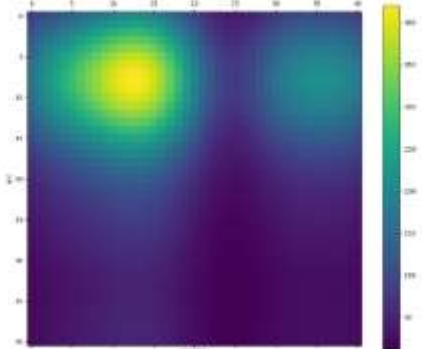
(e)

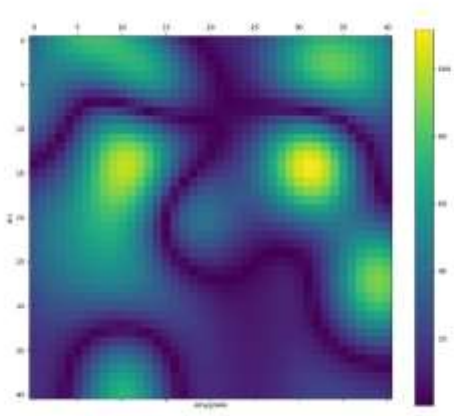
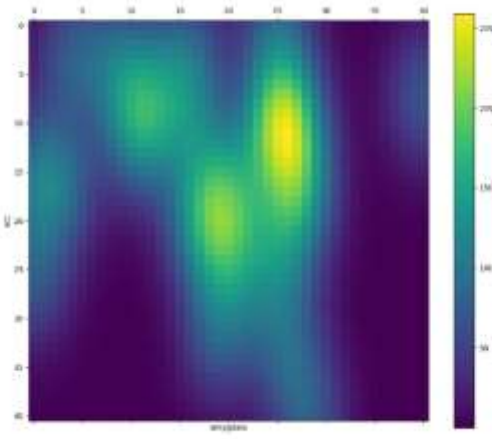
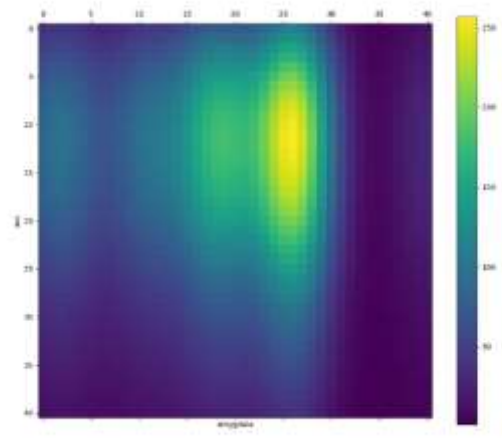
The joint distribution is shown as below:

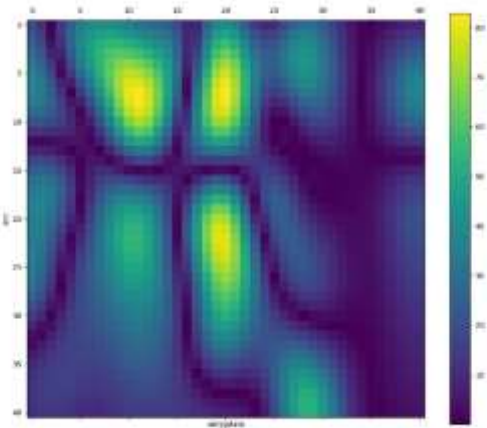
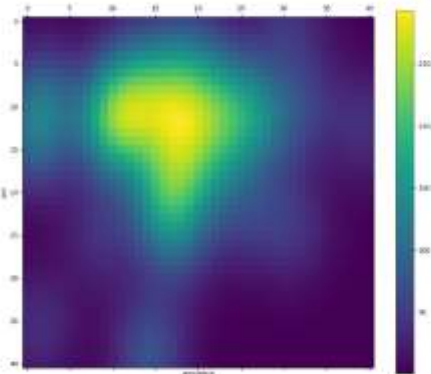
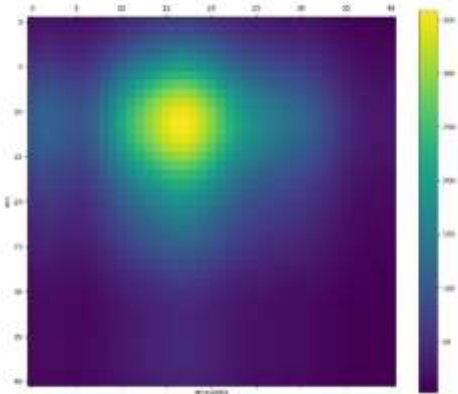
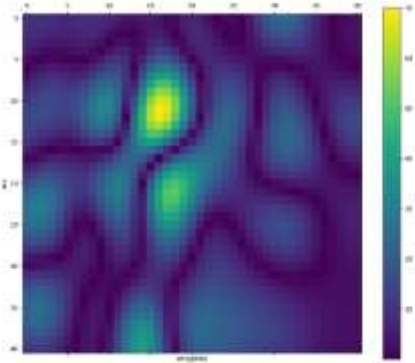


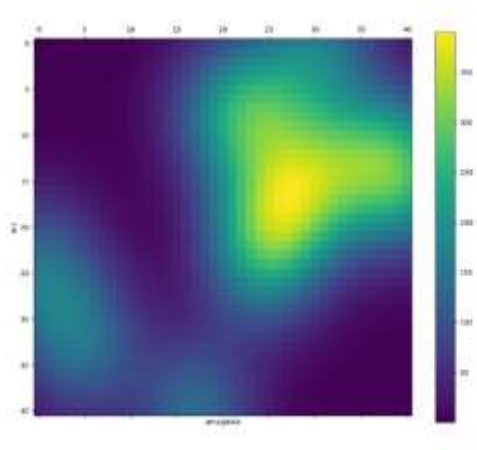
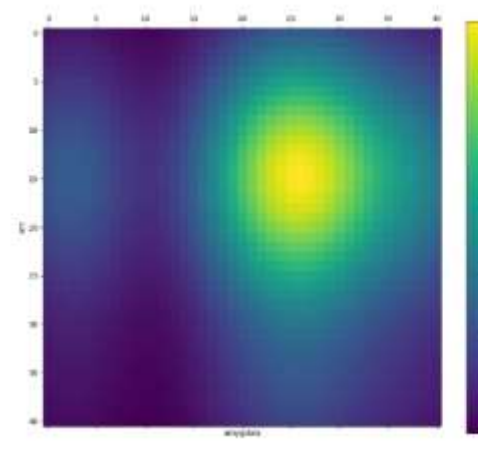
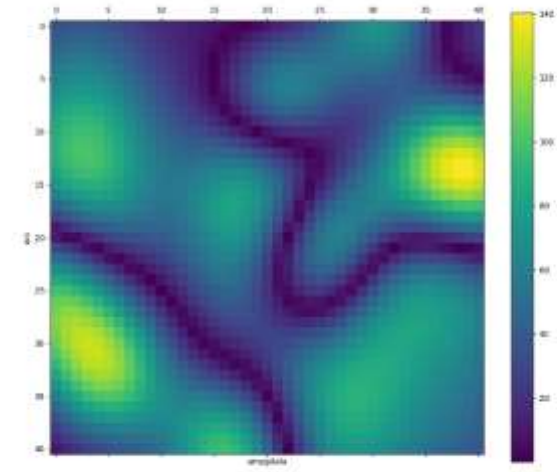
c=4	
c=5	

(f)

Orientation	$p(amygdala, acc)$	$p(amygdala)p(acc)$
c=2		

	$ p(amygdala, acc) - p(amygdala)p(acc) $	
		
c=3	$p(amygdala, acc)$	$p(amygdala)p(acc)$
		
	$ p(amygdala, acc) - p(amygdala)p(acc) $	

		
c=4	$p(amygdala, acc)$	$p(amygdala)p(acc)$
		
	$ p(amygdala, acc) - p(amygdala)p(acc) $	
		

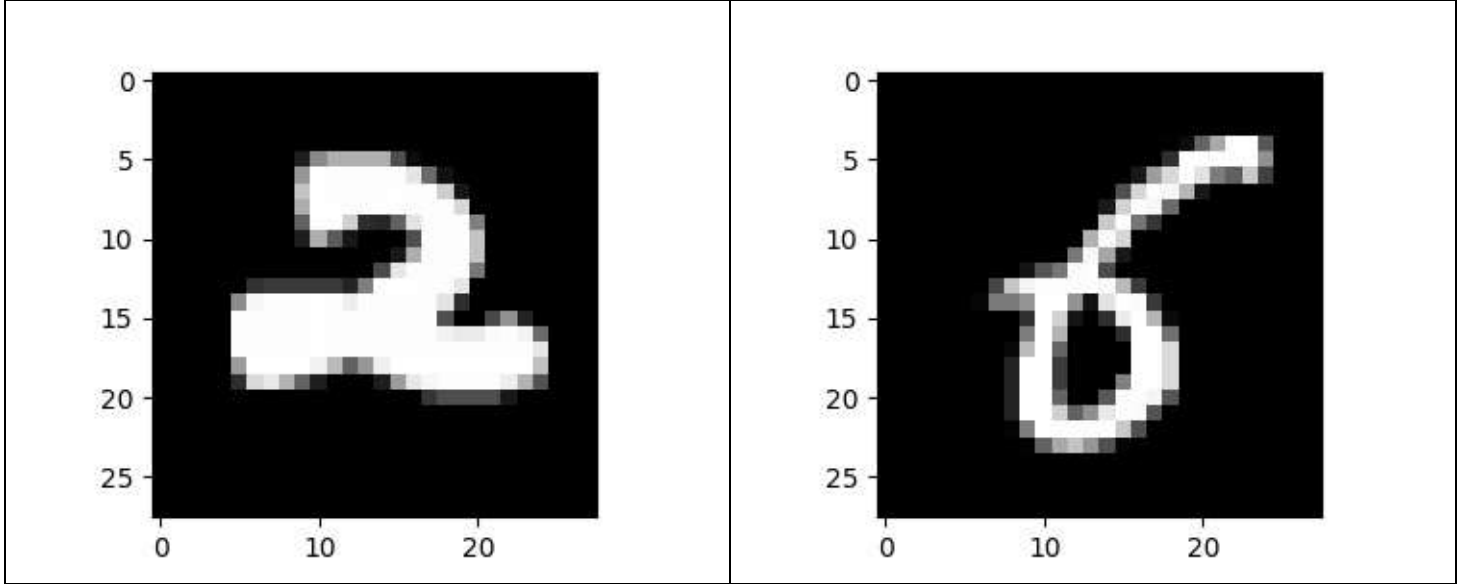
c=5	$p(amygdala, acc)$	
		
	$ p(amygdala, acc) - p(amygdala)p(acc) $	
		

As you can see, for each Orientation, the difference $|p(amygdala, acc) - p(amygdala)p(acc)|$ does not equal to 0, which means the two parts of brain are related conditionally on the political orientation.

2 Implementing EM

(a).

The raw images for “2” and “6” are like the following:



(b).

For the E-step:

Since z^i is i.i.d and x^i is i.i.d, so we could get posterior distribution of $q(z^1, z^2, \dots, z^m)$ like this:

$$q(z^1, z^2, \dots, z^m) = \prod_{i=1}^m p(z^i | x^i, \theta^t)$$

To be detailed, I define:

$$\begin{aligned} \tau_k^i = p(z^i = k | x^i, \theta^t) &= \frac{p(x^i | z^i = k, \theta^t) p(z^i = k)}{\sum_{j=1..K} p(z^i = j, x^i)} \\ &= \frac{\pi_k N(x^i | \mu_k, \Sigma_k)}{\sum_{j=1..K} \pi_j N(x^i | \mu_j, \Sigma_j)} \\ &= \frac{\pi_k \frac{1}{(2\pi)^{n/2} |\Sigma_k|^{1/2}} \exp(-\frac{1}{2} (x^i - \mu_k)^T \Sigma_k^{-1} (x^i - \mu_k))}{\sum_{j=1..K} \pi_j \frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp(-\frac{1}{2} (x^i - \mu_j)^T \Sigma_j^{-1} (x^i - \mu_j))} \end{aligned}$$

$$= \frac{\pi_k \frac{1}{|\Sigma_k|^{1/2}} \exp \left(-\frac{1}{2} (x^i - \mu_k)^T \Sigma_k^{-1} (x^i - \mu_k) \right)}{\sum_{j=1..K} \pi_j \frac{1}{|\Sigma_j|^{1/2}} \exp \left(-\frac{1}{2} (x^i - \mu_j)^T \Sigma_j^{-1} (x^i - \mu_j) \right)}$$

If I take the expectation over $q(z^1, z^2, \dots, z^m)$ with respect to likelihood function $f(\theta)$, the expectation should be the lower bound of the maximum value of $f(\theta)$, since expectation is something like take the average. The lower bound of $f(\theta)$ is shown like the following:

$$\begin{aligned} f(\theta) &= E_{q(z^1, z^2, \dots, z^m)} \left[\log \prod_{i=1}^m p(z^i, x^i | \theta^t) \right] \\ &= E_{q(z^1, z^2, \dots, z^m)} \sum_{i=1}^m \log [p(z^i, x^i | \theta^t)] \\ &= \sum_{i=1}^m E_{p(z^i | x^i, \theta^t)} \log [p(z^i, x^i | \theta^t)] \\ &= \sum_{i=1}^m \sum_{k=1}^K \tau_k^i \left[\log [p(z^i, x^i | \theta^t)] \right] \\ &= \sum_{i=1}^m \sum_{k=1}^K \tau_k^i \left[\log \pi_k - \frac{1}{2} (x^i - \mu_k)^T \Sigma_k^{-1} (x^i - \mu_k) - \frac{1}{2} \log |\Sigma_k| - \frac{n}{2} \log (2\pi) \right] \end{aligned}$$

The M step is like this:

Now we want to maximize our lower bound $f(\theta)$, and now our unknown variables are π_k, μ_k, Σ_k . We want to find the expression of the unknown variables which can lead to the maximum of $f(\theta)$. We can notice that there is a constraint for one variable, that is $\sum \pi_k = 1$. So we could use Lagrange Multiplier Method:

$$L = \sum_{i=1}^m \sum_{k=1}^K \tau_k^i \left[\log \pi_k - \frac{1}{2} (x^i - \mu_k)^T \Sigma_k^{-1} (x^i - \mu_k) - \frac{1}{2} \log |\Sigma_k| - \frac{n}{2} \log (2\pi) \right] + \lambda \left(1 - \sum_{k=1}^K \pi_k \right)$$

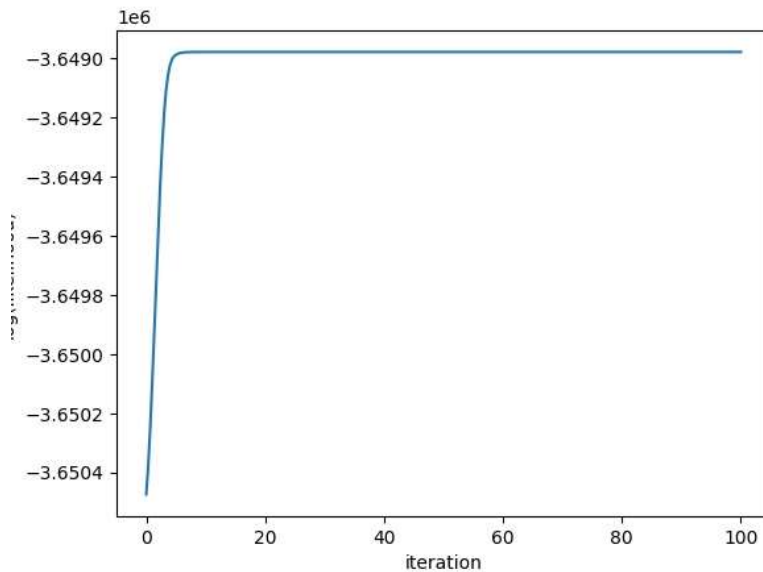
Now I take the derivative with respect to each variable:

$$\begin{aligned} \frac{\partial L}{\partial \pi_k} &= \sum_{i=1}^m \frac{\tau_k^i}{\pi_k} - \lambda = 0 \\ \Rightarrow \pi_k &= \frac{1}{\lambda} \sum_{i=1}^m \tau_k^i \\ \frac{\partial L}{\partial \mu_k} &= \sum_{i=1}^m \tau_k^i [\Sigma_k^{-1} (x^i - \mu_k)] = 0 \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^m \tau_k^i [(x^i - \mu_k)] = 0 \\
\Rightarrow \mu_k &= \frac{\sum_{i=1}^m \tau_k^i x^i}{\sum_{i=1}^m \tau_k^i} \\
\frac{\partial L}{\partial \Sigma_k} &= \sum_{i=1}^m \tau_k^i \left[-\frac{1}{2} (\Sigma_k^{-1})^{-1} (x^i - \mu_k)(x^i - \mu_k)^T - \frac{1}{2} \Sigma_k^{-1} \right] = 0 \\
\Rightarrow \Sigma_k &= \frac{\sum_{i=1}^m \tau_k^i \left[-\frac{1}{2} (x^i - \mu_k)(x^i - \mu_k)^T \right]}{\sum_{i=1}^m \tau_k^i}
\end{aligned}$$

(c)

As what can be shown in the picture, the x Axis is the EM algorithm iteration times, and the y axis is the value of $\log(\text{likelihood})$. The graph shows that with the iteration times increase, the value of $\log(\text{likelihood})$ become steady and meet the local maximum value, which indicate that the EM algorithm converges.

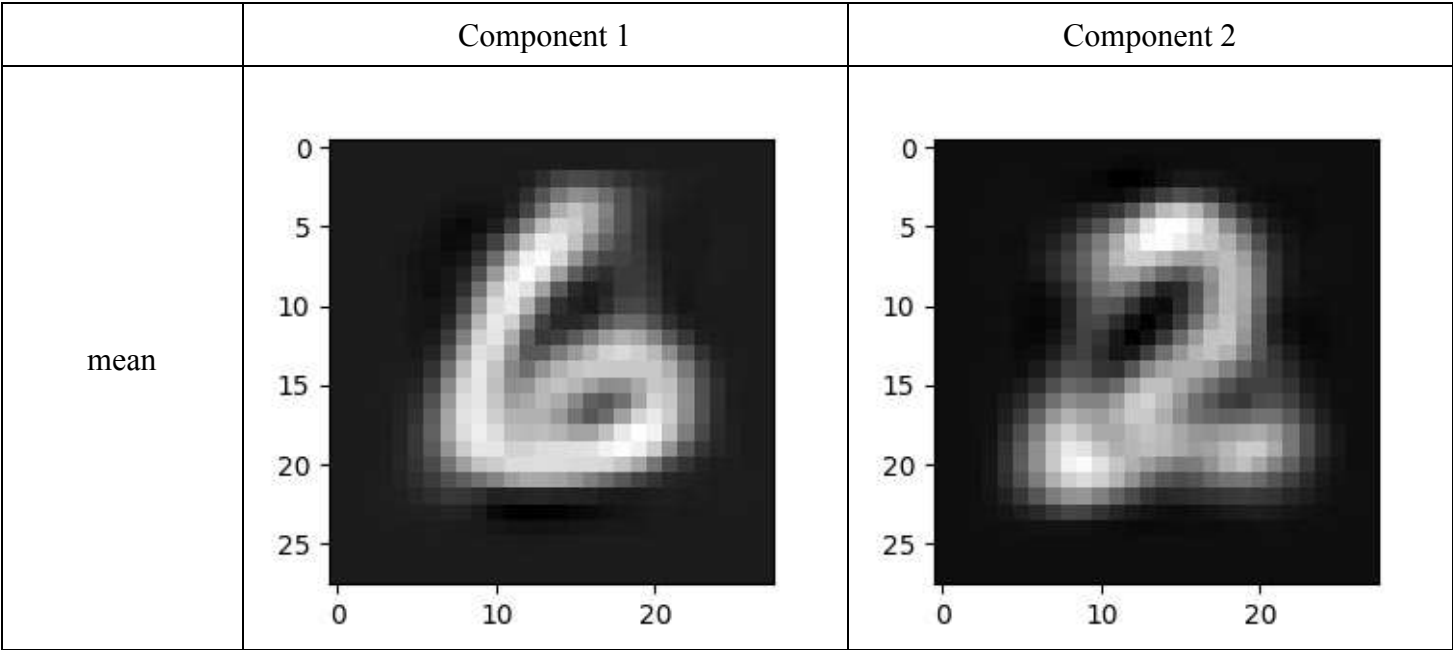


(d)

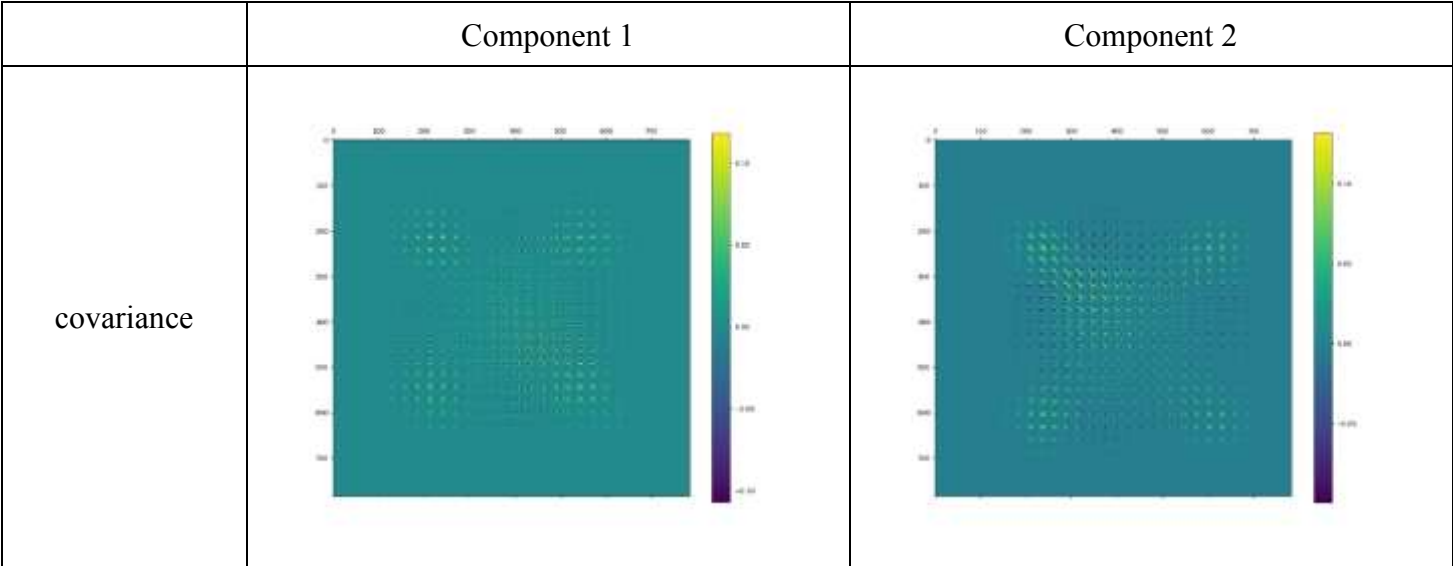
There are two components for the EM algorithm. The weights are shown as below:

	Component 1	Component 2
weight	0.5069	0.4931

The means for each component:



The covariance for each component:



(e)

For the original data, I find that the 0~1031 sample is label “2”, which should belong to component 2. The 1032~1989 sample is label “6”, which should belong to component 1. The mismatch rates for Kmeans and EM are shown as below:

	Kmeans	EM
Mismatch Rate	6.3%	3.4%

So, we could conclude that the EM shows better behavior in clustering.

