

ISYE 6740, Fall 2021, Homework 3

100 points

Prof. Yao Xie

1. Conceptual questions. [20 points]

1. (10 points) Based on the outline given in the lecture, show that the maximum likelihood estimate (MLE) for Gaussian mean and variance parameters are given by

$$\hat{\mu} = \frac{1}{m} \sum_{i=1}^m x^i, \quad \hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m (x^i - \hat{\mu})^2,$$

respectively. Please make sure to complete details of the derivations.

2. (5 points) Please give one advantage of KDE over histogram, and give one advantage of histogram over KDE.
3. (5 points) For EM algorithm for GMM, please show how to use Bayes rule to drive τ_k^i in closed-form expression.

2. Density estimation: Psychological experiments. [40 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 third party histogram and KDE packages; 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

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

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) (5 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 $h > 0$.
- (b) (5 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.
- (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 what you have observed: is the distribution unimodal or bi-modal? Are there any outliers?

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

- (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**: $p(\text{amygdala} | \text{orientation} = c)$, $c = 2, \dots, 5$, using KDE. Set an appropriate kernel bandwidth $h > 0$. Do the same for the volume of the **acc**: plot $p(\text{acc} | \text{orientation} = c)$, $c = 2, \dots, 5$ 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.

- (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**: $p(\text{amygdala}, \text{acc} | \text{orientation} = c)$, $c = 2, \dots, 5$. 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.

3. Implementing EM for MNIST dataset. [40 points]

Implement the EM algorithm for fitting a Gaussian mixture model for the MNIST handwritten 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 1990 images in total, 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 mapping 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).

- (a) (10 points) Write down detailed expression of the E-step and M-step in the EM algorithm (hint: when computing τ_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).

(b) (15 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.

- (c) (5 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).
- (d) (10 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?