

ISYE 6740, Fall 2025, Homework 3

100 points + 10 bonus points

Student Name Here

Provided Data:

Questions marked with GS must be submitted to Gradescope. You must still include all your results and explanations in this PDF, and include your code in your canvas submission to receive credit. Failure to pass all gradescope tests will result in a 50% penalty to the points for that question.

This assignment does not have any gradescope requirements

- Q2: Density Estimation: Psychological experiments [n90pol.csv]
- Q3: Implementing EM for MNIST Dataset [data.dat, data.mat, label.mat, label.dat]

1. Conceptual questions. [20 points]

1. (5 points) Please compare the pros and cons of KDE over histogram, and give at least one advantage and disadvantage to each.
2. (5 points) Explain why the maximum likelihood estimation (MLE) cannot be applied directly to estimate the parameters of a GMM. Additionally, what is the standard approach used to fit a GMM effectively?
3. (5 points) For the EM algorithm for GMM, please show how to use the Bayes rule to derive τ_k^i in a closed-form expression.
4. (5 points) Based on the outline given in the lecture, show that the maximum likelihood estimate (MLE) for Gaussian random variable using n -dimensional observations x^1, \dots, x^m , that are *i.i.d.* (independent and identically distributed) following the distribution $\mathcal{N}(\mu, \Sigma)$, and the mean and variance parameters are given by

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

respectively. Please show the work for your derivations in full detail. Discuss what is the rank of $\hat{\Sigma}$ when $m < n$? In the extreme case, when $m = 1$ (i.e., only one sample), what is the rank of $\hat{\Sigma}$? This question shows that, when the data is high-dimensional, the sample size needs to be sufficiently large relative to the dimension (otherwise we need special treatment to handle low-rank covariance matrix.)

2. Density estimation: Psychological experiments. [25 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

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

1. (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 $h > 0$.
2. (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.

3. (5 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 $h > 0$. Keep in mind that your choice of bandwidth can heavily affect your distributions, and you should be experimenting to ensure you are capturing the true distribution.

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?

Are the two variables (**amygdala**, **acc**) likely to be independent or not? Please support your argument with reasonable investigations.

4. (5 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.

5. (5 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 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 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).

(In this question, we use the same set of data from the provided data files for training and testing)

1. (10 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.

2. (20 points) Report the fitted GMM model when EM has terminated in your algorithms as follows:
 - The numerical weights for each component
 - The mean of each component by mapping it back to the original space and reformatting the vectors into 28-by-28 matrices. These should be displayed as images, ideally corresponding to a kind of “average” of the images.
 - Two 4-by-4 covariance matrices by visualizing their intensities (i.e. a gray-scaled image or heatmap.)
3. (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 (1 - Accuracy) 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?

4. Correlated features (motivation for PCA) (15 points)

We will consider a simple question related to our house price prediction with correlated features. You are predicting house prices with two features:

- x_1 = number of **bedrooms**
- x_2 = number of **bathrooms**

Suppose the (standardized) features are perfectly correlated: $x_1 = 1.5x_2$ for every sample. You collect $m = 10$ samples and fit a linear model **with an intercept**:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon.$$

Let the $n \times 3$ design matrix be $X = [\mathbf{1}^T; \mathbf{x}_1^T; \mathbf{x}_2^T] \in \mathbb{R}^{3 \times m}$. Define the **sample covariance matrix**

$$G = \frac{1}{m} X X^T.$$

Hint

Professor Xie covers a similar example in detail in her recorded office hours, and the slides for this are under Module Two: `pca_example.pdf`.

1. (5 points) Express G analytically in terms of

$$T = \sum_{i=1}^m x_{2,i}, \quad S = \sum_{i=1}^m x_{2,i}^2,$$

using the relation $x_{1,i} = 1.5x_{2,i}$. Show that $\text{rank}(G) = 2$.

2. (5 points) Find the the largest eigenvector for G , i.e., the weights to combine features, using a mathematical proof.
3. (5 points) (Numerical check) Use the bathroom counts $[1, 2, 2, 3, 1, 2, 3, 2, 1, 2]$ for the 10 homes. Compute G and verify its rank.

5. De-bias review system using EM. [Bonus, 10 points]

In this question, we will develop an algorithm to remove individual reviewer's bias from their score. Consider the following problem. There are P papers submitted to a machine learning conference. Each of R reviewers reads each paper, and gives it a score indicating how good he/she thought that paper was. We let $x^{(pr)}$ denote the score that reviewer r gave to paper p . A high score means the reviewer liked the paper, and represents a recommendation from that reviewer that it be accepted for the conference. A low score means the reviewer did not like the paper.

We imagine that each paper has some “intrinsic” true value that we denote by μ_p , where a large value means it's a good paper. Each reviewer is trying to estimate, based on reading the paper, what μ_p is; the score reported $x^{(pr)}$ is then reviewer r 's guess of μ_p .

However, some reviewers are just generally inclined to think all papers are good and tend to give all papers high scores; other reviewers may be particularly nasty and tend to give low scores to everything. (Similarly, different reviewers may have different amounts of variance in the way they review papers, making some reviewers more consistent/reliable than others.) We let ν_r denote the “bias” of reviewer r . A reviewer with bias ν_r is one whose scores generally tend to be ν_r higher than they should be.

All sorts of different random factors influence the reviewing process, and hence we will use a model that incorporates several sources of noise. Specifically, we assume that reviewers's scores are generated by a random process given as follows:

$$\begin{aligned} y^{(p)} &\sim \mathcal{N}(\mu_p, \sigma_p^2) \\ z^{(r)} &\sim \mathcal{N}(\nu_r, \tau_r^2) \\ x^{(pr)} | y^{(p)}, z^{(r)} &\sim \mathcal{N}(y^{(p)} + z^{(r)}, \sigma^2). \end{aligned}$$

The variables $y^{(p)}$ and $z^{(r)}$ are independent; the variables (x, y, z) for different paper-reviewer pairs are also jointly independent. Also, we only ever observe the $x^{(pr)}$ s; thus, the $y^{(p)}$ s and $z^{(r)}$ s are all latent random variables.

We would like to estimate the parameters $\mu_p, \sigma_p^2, \nu_r, \tau_r^2$. If we obtain good estimates of the papers “intrinsic values” μ_p , these can then be used to make acceptance/rejection decisions for the conference.

We will estimate the parameters by maximizing the marginal likelihood of the data $\{x^{(pr)}; p = 1, \dots, P, r = 1, \dots, R\}$. This problem has latent variables $y^{(p)}$ s and $z^{(r)}$ s, and the maximum likelihood problem cannot be solved in closed form. So, we will use EM.

Your task is to derive the EM update equations. For simplicity, you need to treat only $\{\mu_p, \sigma_p^2; p = 1 \dots, P\}$ and $\{\nu_r, \tau_r^2; r = 1 \dots R\}$ as parameters, i.e. treat σ^2 (the conditional variance of $x^{(pr)}$ given $y^{(p)}$ and $z^{(r)}$) as a fixed, known constant.

1. Derive the E-step (5 points)

- 1.1. The joint distribution $p(y^{(p)}, z^{(r)}, x^{(pr)})$ has the form of a multivariate Gaussian density. Find its associated mean vector and covariance matrix in terms of the

parameters $\mu_p, \sigma_p^2, \nu_r, \tau_r^2$ and σ^2 . [Hint: Recognize that $x^{(pr)}$ can be written as $x^{(pr)} = y^{(p)} + z^{(r)} + \epsilon^{(pr)}$, where $\epsilon^{(pr)} \sim \mathcal{N}(0, \sigma^2)$ is independent Gaussian noise.]

- 1.2. Derive an expression for $Q_{pr}(\theta'|\theta) = \mathbb{E}[\log p(y^{(p)}, z^{(r)}, x^{(pr)})|x^{(pr)}, \theta]$ using the conditional distribution $p(y^{(p)}, z^{(r)}|x^{(pr)})$ (E-step) (Hint, you may use the rules for conditioning on subsets of jointly Gaussian random variables.)
2. (5 points) Derive the M-step to update the parameters μ_p, σ_p^2, ν_r , and τ_r^2 . [Hint: It may help to express an approximation to the likelihood in terms of an expectation with respect to $(y^{(p)}, z^{(r)})$ drawn from a distribution with density $Q_{pr}(y^{(p)}, z^{(r)})$.]

Remark: John Platt (whose SMO algorithm you've seen) implemented a method quite similar to this one to estimate the papers' true scores. (There, the problem was a bit more complicated because not all reviewers reviewed every paper, but the essential ideas are the same.) Because the model tried to estimate and correct for reviewers' biases, its estimates of the paper's value were significantly more useful for making accept/reject decisions than the reviewers' raw scores for a paper.