

**Due: Wednesday, February 24 at 11:59 pm**

This homework consists of coding assignments and math problems.

**Begin early; you can submit models to Kaggle only twice a day!**

**DELIVERABLES:**

1. Submit your predictions for the test sets to Kaggle as early as possible. Include your Kaggle scores in your write-up. The Kaggle competition for this assignment can be found at
  - MNIST: <https://www.kaggle.com/c/spring21-cs189-hw3-mnist>
  - SPAM: <https://www.kaggle.com/c/spring21-cs189-hw3-spam>
2. Write-up: Submit your solution in **PDF** format to “Homework 3 Write-Up” in Gradescope.
  - On the first page of your write-up, please list students with whom you collaborated
  - Start each question on a new page. If there are graphs, include those graphs on the same pages as the question write-up. **DO NOT** put them in an appendix. We need each solution to be self-contained on pages of its own.
  - **Only PDF uploads to Gradescope will be accepted.** You are encouraged use  $\text{\LaTeX}$  or Word to typeset your solution. You may also scan a neatly handwritten solution to produce the PDF.
  - **Replicate all your code in an appendix.** Begin code for each coding question in a fresh page. Do not put code from multiple questions in the same page. When you upload this PDF on Gradescope, *make sure* that you assign the relevant pages of your code from appendix to correct questions.
  - While collaboration is encouraged, *everything* in your solution must be your (and only your) creation. Copying the answers or code of another student is strictly forbidden. Furthermore, all external material (i.e., *anything* outside lectures and assigned readings, including figures and pictures) should be cited properly. We wish to remind you that consequences of academic misconduct are *particularly severe*!
3. Code: Submit your code as a .zip file to “Homework 3 Code”.
  - **Set a seed for all pseudo-random numbers generated in your code.** This ensures your results are replicated when readers run your code.
  - Include a README with your name, student ID, the values of random seed (above) you used, and any instructions for compilation.
  - Do NOT provide any data files. Supply instructions on how to add data to your code.
  - Code requiring exorbitant memory or execution time might not be considered.

- Code submitted here must match that in the PDF Write-up. The Kaggle score will not be accepted if the code provided a) does not compile or b) compiles but does not produce the file submitted to Kaggle.
4. The assignment covers concepts on Gaussian distributions and classifiers. Some of the material may not have been covered in lecture; you are responsible for finding resources to understand it.

## 1 Honor Code

Declare and sign the following statement (Mac Preview and FoxIt PDF Reader, among others, have tools to let you sign a PDF file):

*“I certify that all solutions are entirely my own and that I have not looked at anyone else’s solution. I have given credit to all external sources I consulted.”*

Signature: \_\_\_\_\_

## 2 Gaussian Classification

Let  $f(x | C_i) \sim \mathcal{N}(\mu_i, \sigma^2)$  for a two-class, one-dimensional classification problem with classes  $C_1$  and  $C_2$ ,  $P(C_1) = P(C_2) = 1/2$ , and  $\mu_2 > \mu_1$ .

1. Find the Bayes optimal decision boundary and the corresponding Bayes decision rule by finding the point(s) at which the posterior probabilities are equal.
2. Suppose the decision boundary for your classifier is  $x = b$ . The Bayes error is the probability of misclassification, namely,

$$P_e = P((C_1 \text{ misclassified as } C_2) \cup (C_2 \text{ misclassified as } C_1)).$$

Show that the Bayes error associated with this decision rule, in terms of  $b$ , is

$$P_e(b) = \frac{1}{2\sqrt{2\pi}\sigma} \left( \int_{-\infty}^b \exp\left(-\frac{(x-\mu_2)^2}{2\sigma^2}\right) dx + \int_b^{\infty} \exp\left(-\frac{(x-\mu_1)^2}{2\sigma^2}\right) dx \right).$$

3. Using the expression above for the Bayes error, calculate the optimal decision boundary  $b^*$  that minimizes  $P_e(b)$ . How does this value compare to that found in part 1? *Hint:  $P_e(b)$  is convex for  $\mu_1 < b < \mu_2$ .*

## 3 Isocontours of Normal Distributions

Let  $f(\mu, \Sigma)$  be the probability density function of a normally distributed random variable in  $\mathbb{R}^2$ . Write code to plot the isocontours of the following functions, each on its own separate figure. Make sure it is clear which figure belongs to which part. You’re free to use any plotting libraries or stats utilities available in your programming language; for instance, in Python you can use Matplotlib and SciPy.

1.  $f(\mu, \Sigma)$ , where  $\mu = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$  and  $\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$ .
2.  $f(\mu, \Sigma)$ , where  $\mu = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$  and  $\Sigma = \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix}$ .
3.  $f(\mu_1, \Sigma_1) - f(\mu_2, \Sigma_2)$ , where  $\mu_1 = \begin{bmatrix} 0 \\ 2 \end{bmatrix}$ ,  $\mu_2 = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$  and  $\Sigma_1 = \Sigma_2 = \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix}$ .
4.  $f(\mu_1, \Sigma_1) - f(\mu_2, \Sigma_2)$ , where  $\mu_1 = \begin{bmatrix} 0 \\ 2 \end{bmatrix}$ ,  $\mu_2 = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$ ,  $\Sigma_1 = \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix}$  and  $\Sigma_2 = \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix}$ .
5.  $f(\mu_1, \Sigma_1) - f(\mu_2, \Sigma_2)$ , where  $\mu_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ ,  $\mu_2 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$ ,  $\Sigma_1 = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$  and  $\Sigma_2 = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$ .

## 4 Eigenvectors of the Gaussian Covariance Matrix

Consider two one-dimensional random variables  $X_1 \sim \mathcal{N}(3, 9)$  and  $X_2 \sim \frac{1}{2}X_1 + \mathcal{N}(4, 4)$ , where  $\mathcal{N}(\mu, \sigma^2)$  is a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ . Write a program that draws  $n = 100$  random two-dimensional sample points from  $(X_1, X_2)$  such that the  $i$ th value sampled from  $X_2$  is calculated based on the  $i$ th value sampled from  $X_1$ . In your code, make sure to choose and set a fixed random number seed for whatever random number generator you use, so your simulation is reproducible, and document your choice of random number seed and random number generator in your write-up. For each of the following parts, include the corresponding output of your program.

- (a) Compute the mean (in  $\mathbb{R}^2$ ) of the sample.
- (b) Compute the  $2 \times 2$  covariance matrix of the sample.
- (c) Compute the eigenvectors and eigenvalues of this covariance matrix.
- (d) On a two-dimensional grid with a horizontal axis for  $X_1$  with range  $[-15, 15]$  and a vertical axis for  $X_2$  with range  $[-15, 15]$ , plot
  - (i) all  $n = 100$  data points, and
  - (ii) arrows representing both covariance eigenvectors. The eigenvector arrows should originate at the mean and have magnitudes equal to their corresponding eigenvalues.
- (e) Let  $U = [v_1 \ v_2]$  be a  $2 \times 2$  matrix whose columns are the eigenvectors of the covariance matrix, where  $v_1$  is the eigenvector with the larger eigenvalue. We use  $U^T$  as a rotation matrix to rotate each sample point from the  $(X_1, X_2)$  coordinate system to a coordinate system aligned with the eigenvectors. (As  $U^T = U^{-1}$ , the matrix  $U$  reverses this rotation, moving back from the eigenvector coordinate system to the original coordinate system). *Center* your sample points by subtracting the mean  $\mu$  from each point; then rotate each point by  $U^T$ , giving  $x_{\text{rotated}} = U^T(x - \mu)$ . Plot these rotated points on a new two dimensional-grid, again with both axes having range  $[-15, 15]$ .

In your plots, **clearly label the axes and include a title**. Moreover, **make sure the horizontal and vertical axis have the same scale!** The aspect ratio should be one.

## 5 Classification and Risk

Suppose we have a classification problem with classes labeled  $1, \dots, c$  and an additional “doubt” category labeled  $c + 1$ . Let  $r : \mathbb{R}^d \rightarrow \{1, \dots, c + 1\}$  be a decision rule. Define the loss function

$$L(r(x) = i, y = j) = \begin{cases} 0 & \text{if } i = j \quad i, j \in \{1, \dots, c\}, \\ \lambda_r & \text{if } i = c + 1, \\ \lambda_s & \text{otherwise,} \end{cases}$$

where  $\lambda_r \geq 0$  is the loss incurred for choosing doubt and  $\lambda_s \geq 0$  is the loss incurred for making a misclassification. Hence the risk of classifying a new data point  $x$  as class  $i \in \{1, 2, \dots, c + 1\}$  is

$$R(r(x) = i|x) = \sum_{j=1}^c L(r(x) = i, y = j) P(Y = j|x).$$

1. Show that the following policy obtains the minimum risk when  $\lambda_r \leq \lambda_s$ .
  - (a) Choose class  $i$  if  $P(Y = i|x) \geq P(Y = j|x)$  for all  $j$  and  $P(Y = i|x) \geq 1 - \lambda_r/\lambda_s$ ;
  - (b) Choose doubt otherwise.
2. What happens if  $\lambda_r = 0$ ? What happens if  $\lambda_r > \lambda_s$ ? Explain why this is consistent with what one would expect intuitively.

## 6 Maximum Likelihood Estimation and Bias

Let  $X_1, \dots, X_n \in \mathbb{R}$  be  $n$  sample points drawn independently from normal distributions such that  $X_i \sim \mathcal{N}(\mu, \sigma_i^2)$ , where  $\sigma_i = \sigma/\sqrt{i}$  for some parameter  $\sigma$ . (Every sample point comes from a distribution with a different variance.)

- (a) Derive the maximum likelihood estimates, denoted  $\hat{\mu}$  and  $\hat{\sigma}$ , for the mean  $\mu$  and the parameter  $\sigma$ . You may write an expression for  $\hat{\sigma}^2$  rather than  $\hat{\sigma}$  if you wish—it’s probably simpler that way. Show all your work.
- (b) Given the true value of a statistic  $\theta$  and an estimator  $\hat{\theta}$  of that statistic, we define the *bias* of the estimator to be the the expected difference from the true value. That is,

$$\text{bias}(\hat{\theta}) = \mathbb{E}[\hat{\theta}] - \theta.$$

We say that an estimator is *unbiased* if its bias is 0.

Either prove or disprove the following statement: *The MLE sample estimator  $\hat{\mu}$  is unbiased. Hint: Neither the true  $\mu$  nor true  $\sigma^2$  are known when estimating sample statistics, thus we need to plug in appropriate estimators.*

1. Either prove or disprove the following statement: *The MLE sample estimator  $\hat{\sigma}^2$  is unbiased. Hint: Neither the true  $\mu$  nor true  $\sigma^2$  are known when estimating sample statistics, thus we need to plug in appropriate estimators.*

## 7 Covariance Matrices and Decompositions

As described in lecture, the covariance matrix  $\text{Var}(R) \in \mathbb{R}^{d \times d}$  for a random variable  $R \in \mathbb{R}^d$  with mean  $\mu$  is

$$\text{Var}(R) = \text{Cov}(R, R) = \mathbb{E}[(R - \mu)(R - \mu)^\top] = \begin{bmatrix} \text{Var}(R_1) & \text{Cov}(R_1, R_2) & \dots & \text{Cov}(R_1, R_d) \\ \text{Cov}(R_2, R_1) & \text{Var}(R_2) & & \text{Cov}(R_2, R_d) \\ \vdots & & \ddots & \vdots \\ \text{Cov}(R_d, R_1) & \text{Cov}(R_d, R_2) & \dots & \text{Var}(R_d) \end{bmatrix},$$

where  $\text{Cov}(R_i, R_j) = \mathbb{E}[(R_i - \mu_i)(R_j - \mu_j)]$  and  $\text{Var}(R_i) = \text{Cov}(R_i, R_i)$ .

If the random variable  $R$  is sampled from the multivariate normal distribution  $\mathcal{N}(\mu, \Sigma)$  with the PDF

$$f(x) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-((x-\mu)^\top \Sigma^{-1} (x-\mu))/2},$$

then  $\text{Var}(R) = \Sigma$ .

Given  $n$  points  $X_1, X_2, \dots, X_n$  sampled from  $\mathcal{N}(\mu, \Sigma)$ , we can estimate  $\Sigma$  with the maximum likelihood estimator

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)(X_i - \mu)^\top,$$

which is also known as the covariance matrix of the sample.

- The estimate  $\hat{\Sigma}$  makes sense as an approximation of  $\Sigma$  only if  $\hat{\Sigma}$  is invertible. Under what circumstances is  $\hat{\Sigma}$  not invertible? Make sure your answer is complete; i.e., it includes all cases in which the covariance matrix of the sample is singular. Express your answer in terms of the geometric arrangement of the sample points  $X_i$ .
- Suggest a way to fix a singular covariance matrix estimator  $\hat{\Sigma}$  by replacing it with a similar but invertible matrix. Your suggestion may be a kludge, but it should not change the covariance matrix too much. Note that infinitesimal numbers do not exist; if your solution uses a very small number, explain how to calculate a number that is sufficiently small for your purposes.
- Consider the normal distribution  $\mathcal{N}(0, \Sigma)$  with mean  $\mu = 0$ . Consider all vectors of length 1; i.e., any vector  $x$  for which  $\|x\| = 1$ . Which vector(s)  $x$  of length 1 maximizes the PDF  $f(x)$ ? Which vector(s)  $x$  of length 1 minimizes  $f(x)$ ? Your answers should depend on the properties of  $\Sigma$ . Explain your answer.

## 8 Gaussian Classifiers for Digits and Spam

In this problem, you will build classifiers based on Gaussian discriminant analysis. Unlike Homework 1, you are NOT allowed to use any libraries for out-of-the-box classification (e.g. `sklearn`). You may use anything in `numpy` and `scipy`.

The training and test data can be found with this homework. Don't use the training/test data from Homework 1, as they have changed for this homework. You can verify that your data files and python environment are setup properly by running `sanity.py`. Submit your predicted class labels for the test data on the Kaggle competition website and be sure to include your Kaggle display name and scores in your writeup. Also be sure to include an appendix of your code at the end of your writeup.

1. Taking pixel values as features (no new features yet, please), fit a Gaussian distribution to each digit class using maximum likelihood estimation. This involves computing a mean and a covariance matrix for each digit class, as discussed in lecture.

*Hint:* You may, and probably should, contrast-normalize the images before using their pixel values. One way to normalize is to divide the pixel values of an image by the  $l_2$ -norm of its pixel values.

2. (Written answer + graph) Visualize the covariance matrix for a particular class (digit). How do the diagonal terms compare with the off-diagonal terms? What do you conclude from this?
3. Classify the digits in the test set on the basis of posterior probabilities with two different approaches. Feel free to either use the starter code provided in `starter.py` or write your own implementation from scratch

- (a) (Graphs) Linear discriminant analysis (LDA). Model the class conditional probabilities as Gaussians  $\mathcal{N}(\mu_C, \Sigma)$  with different means  $\mu_C$  (for class C) and the same covariance matrix  $\Sigma$ , which you compute by averaging the 10 covariance matrices from the 10 classes.

To implement LDA, you will sometimes need to compute a matrix-vector product of the form  $\Sigma^{-1}x$  for some vector  $x$ . You should **not** compute the inverse of  $\Sigma$  (nor the determinant of  $\Sigma$ ) as it is not guaranteed to be invertable. Instead, you should find a way to solve the implied linear system without computing the inverse. *Hint: How do we solve OLS when the data matrix is singular?*

Hold out 10,000 randomly chosen training points for a validation set (You may reuse your homework 1 solution or an out-of-the-box library for dataset splitting *only*). Classify each image in the validation set into one of the 10 classes. Compute the error rate  $(1 - \frac{\text{\# points correctly classified}}{\text{\# total points}})$  on the validation set and plot it over the following numbers of randomly chosen training points: 100, 200, 500, 1,000, 2,000, 5,000, 10,000, 30,000, 50,000. (Expect some variation in your error rate when few training points are used.)

- (b) (Graphs) Quadratic discriminant analysis (QDA). Model the class conditional probabilities as Gaussians  $\mathcal{N}(\mu_C, \Sigma_C)$ , where  $\Sigma_C$  is the estimated covariance matrix for class C.

(If any of these covariance matrices turn out singular, implement the trick you described in Q7(b). You are welcome to use  $k$ -fold cross validation to choose the right constant(s) for that trick.) Repeat the same tests and error rate calculations you did for LDA.

- (c) (Written answer) Which of LDA and QDA performed better? Why?
- (d) (Written answer + graph) Include a plot of validation error versus the number of training points for each digit. Plot all the 10 curves on the same graph as shown in Figure 1. Which digit is easiest to classify? Include written answers where indicated.

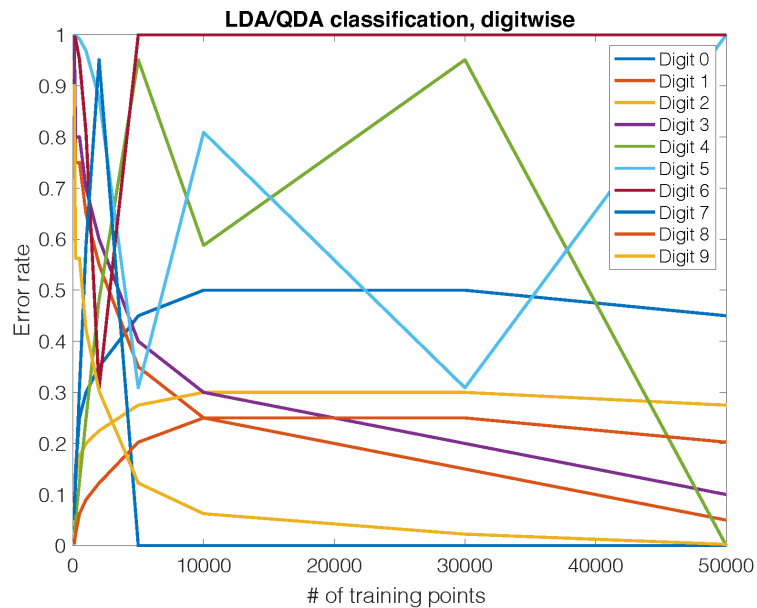


Figure 1: Sample graph with 10 plots

4. Using the `mnist_data.mat`, train your best classifier for the `training_data` and classify the images in the `test_data`. Submit your labels to the online Kaggle competition. Record your optimum prediction rate in your submission. You are welcome to compute extra features for the Kaggle competition, as long as they do not use an exterior learned model for their computation (no transfer learning!). If you do so, please describe your implementation in your assignment. Please use extra features **only** for the Kaggle portion of the assignment.
5. Next, apply LDA or QDA (your choice) to spam. Submit your test results to the online Kaggle competition. Record your optimum prediction rate in your submission. If you use additional features (or omit features), please describe them.

*Optional:* If you use the defaults, expect relatively low classification rates. We suggest using a Bag-Of-Words model. You are encouraged to explore alternative hand-crafted features, and are welcome to use any third-party library to implement them, as long as they do not use a separate model for their computation (no word-2-vec!).