

**Due: Friday, February 23 at 11:59 pm**

- Homework 3 consists of coding assignments and math problems.
- We prefer that you typeset your answers using L<sup>A</sup>T<sub>E</sub>X or other word processing software. If you haven't yet learned L<sup>A</sup>T<sub>E</sub>X, one of the crown jewels of computer science, now is a good time! Neatly handwritten and scanned solutions will also be accepted.
- In all of the questions, **show your work**, not just the final answer.
- 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.
- **Start 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/competitions/cs189-hw3-mnist-spring-2024/>
  - SPAM: <https://www.kaggle.com/competitions/cs189-hw3-spam-spring-2024/>
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 L<sup>A</sup>T<sub>E</sub>X 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”. The code must be in a form that enables the readers to compile (if necessary) and run it to produce your Kaggle submissions.

- **Set a seed for all pseudo-random numbers generated in your code.** This ensures your results are replicated when readers run your code. For example, you can seed numpy with `np.random.seed(42)`.
- Include a README with your name, student ID, the values of random seed you used, and instructions for compiling (if necessary) and running your code. If the data files need to be anywhere other than the main directory for your code to run, let us know where.
- Do **not** submit 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/run or b) runs but does not produce the file submitted to Kaggle.

# 1 Honor Code

Declare and sign the following statement (Mac Preview, PDF Expert, 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  $f_{X|Y=C_i}(x) \sim \mathcal{N}(\mu_i, \sigma^2)$  for a two-class, one-dimensional ( $d = 1$ ) classification problem with classes  $C_1$  and  $C_2$ ,  $P(Y = C_1) = P(Y = C_2) = 1/2$ , 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. Use the 0-1 loss function.
2. Suppose the decision boundary for your classifier is  $x = b$ . The Bayes error is the probability of misclassification, namely,

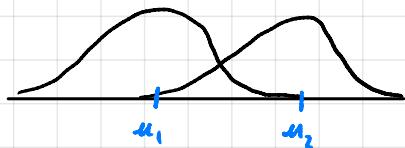
$$P_e = P((C_1 \text{ misclassified as } C_2) \cup (C_2 \text{ misclassified as } C_1)).$$

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  $\mu_1 < b < \mu_2$ .*

## Section 2



$$L(z,y) = \begin{cases} 1 & z \neq y \\ 0 & \text{else} \end{cases}$$

1. Bayes' Optimal decision boundary:

$$P(Y=1|X) = P(Y=-1|X)$$

From gaussian:

$$\frac{P(X|Y=1)}{P(X)} = \frac{P(X|Y=-1)}{P(X)}$$

$$P(X|Y=1) = P(X|Y=-1)$$

$$\left(\frac{x-\mu_1}{\sigma}\right)^2 = \left(\frac{x-\mu_2}{\sigma}\right)^2$$

$$x^2 - 2x\mu_1 - \mu_1^2 = x^2 - 2x\mu_2 - \mu_2^2$$

$$2x(-\mu_1 + \mu_2) = -\mu_2^2 + \mu_1^2$$

$$x = \frac{-\mu_2^2 + \mu_1^2}{2(-\mu_1 + \mu_2)}$$

$$= \frac{(\mu_2^2 - \mu_1^2)}{2(-\mu_1 + \mu_2)}$$

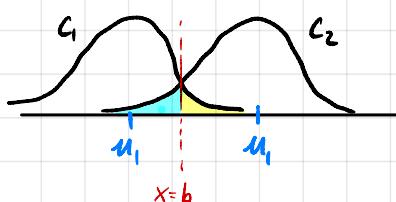
$$= \frac{\mu_2 + \mu_1}{2}$$

Bayes' Decision Rule:

$$r^*(x) = \begin{cases} 1 & \text{if } x > \frac{\mu_2 + \mu_1}{2} \\ -1 & \text{else} \end{cases}$$

given 1 is the class of  $N(\mu_2, \sigma^2)$

2.



$$P_e(C_0) = \int_{-\infty}^b N(\mu_2, \sigma^2) dx + \int_b^\infty N(\mu_1, \sigma^2) dx$$

The highlighted area and eq. shows the probability of making a mistake. Essentially we are calculating the cumulative prob that the classifier will make a mistake.

The blue highlight is the prob that a sample that's  $C_2$  will be classified as  $C_1$ . Similarly the yellow is the prob that a sample in  $C_1$  will be classified as  $C_2$ .

3. Using the FTC

$$0 = e^{\frac{(b-\mu_2)^2}{2\sigma^2}} \cdot 1 - e^{\frac{(b-\mu_2)^2}{2\sigma^2}} \cdot \frac{d}{dx} e^{\frac{(b-\mu_2)^2}{2\sigma^2}} + e^{\frac{(b-\mu_1)^2}{2\sigma^2}} \cdot \frac{d}{dx} e^{\frac{(b-\mu_1)^2}{2\sigma^2}} - e^{\frac{(b-\mu_1)^2}{2\sigma^2}} \cdot 1$$

$$0 = e^{\frac{(b-\mu_2)^2}{2\sigma^2}} - e^{\frac{(b-\mu_1)^2}{2\sigma^2}}$$

$$b = \frac{\mu_2 + \mu_1}{2}$$

The optimal value of b is the same.

### 3 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,  $\lambda_s \geq 0$  is the loss incurred for making a misclassification, and at least one of them is nonzero. Hence the risk of classifying a new data point  $x$  as class  $i \in \{1, 2, \dots, c + 1\}$  is

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

To be clear, the actual label  $Y$  can *never* be  $c + 1$ .

1. Show that the following predictor obtains the minimum risk when  $\lambda_r \leq \lambda_s$ :
  - Choose class  $i$  if  $P(Y = i|x) \geq P(Y = j|x)$  for all  $j$  and  $P(Y = i|x) \geq 1 - \lambda_r/\lambda_s$ ;
  - 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.

## Section 3

1. In English: the predictor chooses  $i$  if the prob of  $x$  being in class  $i$  is greater or equal to any other class  $j$ .

AND

the prob of  $i$  is greater than  $1 - \frac{\text{punishment for doubt}}{\text{punishment for wrong}}$ .

$$\boxed{< 1}$$

Take the complement:

$$-P(Y \neq i | x) \geq -\frac{\lambda_r}{\lambda_s}$$

$$\lambda_s P(Y \neq i | x) \leq \lambda_r$$

$$E[\text{loss} | \text{misclassified}] \leq E[\text{loss} | \text{classify as unsure}]$$

This means, if the expected loss for choosing  $Y=i$  is less than the loss for being unsure, then classify  $Y$  as  $i$ .

Otherwise, if the loss of misclassifying is greater, then categorize it as C+1 / unsure.

This method must obtain the minimum risk, as any other choice increases the expected loss for any  $x, y$ .

2. If  $\lambda_r = 0$

$$\lambda_s P(Y \neq i | x) \leq 0$$

since  $\lambda_s P(Y \neq i | x) > 0$ , this case will never be true. Intuitively this makes it that the cost of getting the prediction wrong is always higher than classifying as unsure. Since all classifications have some uncertainty, that means the classifier will always classify  $y$  as C+1

If  $\lambda_r > \lambda_s$

$$P(Y=i | x) \geq 1 - \frac{\lambda_r}{\lambda_s} \geq 0$$

This is guaranteed to be true, and so the classifier will never choose C+1 for  $y$ .

Intuitively this is when the cost of classifying as C+1 is always greater than getting it wrong.

## 4 Maximum Likelihood Estimation and Bias

Let  $X_1, \dots, X_n \in \mathbb{R}$  be  $n$  sample points drawn independently from univariate 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**.) Note the word “univariate”; we are working in dimension  $d = 1$ , and each “point” is just a real number.

1. Derive the maximum likelihood estimates, denoted  $\hat{\mu}$  and  $\hat{\sigma}$ , for the mean  $\mu$  and the parameter  $\sigma$ . (The formulae from class don’t apply here, because every point has a different variance.) 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.
2. 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.*

3. 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.*
4. Suppose the Variance Fairy drops by to give us the true value of  $\sigma^2$ , so that we only have to estimate  $\mu$ . Given the loss function  $L(\hat{\mu}, \mu) = (\hat{\mu} - \mu)^2$ , what is the risk of our MLE estimator  $\hat{\mu}$ ?

## Section 4

$X_1, \dots, X_n$  when  $X_i \sim N(\mu, (\frac{\sigma}{\sqrt{n}})^2)$   $d=1$

$$f_{X_i}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

$$\ln(f_{X_i}(x)) = -\frac{1}{2}\ln(2\pi) - \ln(\sigma) - \frac{1}{2}\ln(\frac{1}{n}) - \frac{(x_i-\mu)^2}{2\sigma^2}$$

↓

$$l(\mu, \sigma, X_1, \dots, X_n) = \sum_{i=1}^n \left( -\frac{(x_i-\mu)^2}{2\sigma^2} - \ln(\sigma) - \frac{1}{2}\ln(2\pi) - \frac{1}{2}\ln(\frac{1}{n}) \right)$$

$$\begin{aligned} \frac{\partial}{\partial \mu} l(\mu, \sigma, X_1, \dots, X_n) &= \sum_{i=1}^n \left( -\frac{(x_i-\mu)^2}{2\sigma^2} - \ln(\sigma) - \frac{1}{2}\ln(2\pi) - \frac{1}{2}\ln(\frac{1}{n}) \right) \\ &= \sum_{i=1}^n \left( -\frac{1}{2\sigma^2} \frac{d}{dx} (x_i-\mu)^2 \right) \\ &= -\frac{1}{2\sigma^2} \left( \sum_{i=1}^n 2(x_i-\mu) \right) \\ &= -\frac{1}{\sigma^2} \left( \sum_{i=1}^n (x_i-\mu) \right) \\ &= -\frac{1}{\sigma^2} \left( \sum_{i=1}^n (x_i) - \mu \cdot \frac{n(1+n)}{2} \right) \end{aligned}$$

Find critical point  $\hat{\mu} = \bar{x}$

$$\hat{\mu} = \frac{1}{n(1+n)} \sum_{i=1}^n x_i$$

$$\begin{aligned} \frac{\partial l}{\partial \sigma} &= \frac{\partial}{\partial \sigma} \left( \sum_{i=1}^n \left( -\frac{(x_i-\mu)^2}{2\sigma^2} - \ln(\sigma) - \frac{1}{2}\ln(2\pi) - \frac{1}{2}\ln(\frac{1}{n}) \right) \right) \\ &= \frac{\partial}{\partial \sigma} (-n\ln(\sigma)) - \sum_{i=1}^n \frac{\partial}{\partial \sigma} \left( -\frac{(x_i-\mu)^2}{2\sigma} \right) \\ &= -\frac{n}{\sigma} - \sum_{i=1}^n \left( -\frac{(x_i-\mu)^2}{2} \frac{\partial}{\partial \sigma} (\sigma^{-2}) \right) \\ &= -\frac{n}{\sigma} - \sum_{i=1}^n -\frac{(x_i-\mu)^2}{2} \cdot -2\sigma^{-3} \\ &= -\frac{n}{\sigma} - \frac{1}{\sigma^3} \sum_{i=1}^n (x_i-\mu)^2 \end{aligned}$$

Critical points

$$\frac{n}{\sigma} = \frac{1}{\sigma^3} \sum_{i=1}^n (x_i-\mu)^2$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i-\mu)^2$$

2.

$$\frac{2}{n(1+n)} \sum_{i=1}^n x_i$$

$$\begin{aligned} E[\hat{\mu}] &= \frac{2}{n(1+n)} \cdot \sum_{i=1}^n E[x_i] \\ &= \frac{2}{n(1+n)} \cdot n \cdot \frac{n(1+n)}{2} \\ &= \mu \end{aligned}$$

$$\text{Bias}(\hat{x}) = E[\hat{x}] - \mu$$

$$= 0$$

$\therefore \hat{\mu}$  is not biased.

$$\begin{aligned} 3. E[\hat{\sigma}^2] &= \frac{1}{n} \left( E \left[ \sum_i x_i^2 - 2\bar{x}\mu + \mu^2 \right] \right) \\ &= \frac{1}{n} \left( E \left[ \sum_i x_i^2 - 2\bar{x}\mu + \sum_i \mu^2 \right] \right) \\ &= \frac{1}{n} \left( E \left[ \sum_i x_i^2 - 2\mu \sum_i x_i + \frac{\mu^2 n(1+n)}{2} \right] \right) \\ &= \frac{1}{n} \left( E \left[ \sum_i x_i^2 \right] - 2\mu \sum_i x_i + \frac{\mu^2 n(1+n)}{2} \right) \\ &= \frac{1}{n} \left( \sum_{i=1}^n E[x_i^2] - 2\mu \cdot \frac{n\mu(1+n)}{2} + \mu^2 \frac{n(1+n)}{2} \right) \\ &= \frac{1}{n} \left( \frac{E[x_i^2]n(1+n)}{2} - \frac{\mu^2 n(1+n)}{2} \right) \\ &= \frac{(1+n)}{2} (E[x_i^2] - E[x_i]^2) \\ &= \frac{(1+n)}{2} \frac{\sigma^2}{i} \end{aligned}$$

$$\text{Bias}(\hat{\sigma}^2) = E[\hat{\sigma}^2] - \sigma^2$$

$$= \frac{(1+n)}{2} \sigma^2 - \sigma^2$$

$$\neq 0$$

$\therefore$  Variance is biased

$$4. R(\hat{\mu}) = E[L(\hat{\mu}, \mu)]$$

$$= E[\mu^2 - 2\bar{x}\mu + \mu^2]$$

$$= E[\hat{\mu}^2] - 2\mu\hat{\mu} + \mu^2$$

$$\begin{aligned} &= \mu^2 - 2\mu\hat{\mu} + \mu^2 \\ &= 2\mu^2 - 2\mu\hat{\mu} \end{aligned}$$

$$\begin{aligned} &= E \left[ \left( \frac{2}{n(1+n)} \sum_{i=1}^n x_i \right)^2 \right] \\ &= E \left[ \frac{4}{n^2(1+n)^2} \sum_{i=1}^n x_i \cdot \sum_{j=1}^n x_j \right] \\ &= \frac{4}{n^2(1+n)^2} \sum_{i=1}^n E[x_i] \cdot \sum_{j=1}^n E[x_j] \\ &= \frac{4}{n^2(1+n)^2} \cdot \frac{mn(1+n)}{2} \cdot \frac{mn(1+n)}{2} \\ &= \mu^2 \end{aligned}$$

## 5 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  $\mu \in \mathbb{R}^d$  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

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

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

1. 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? Express your answer in terms of the geometric arrangement of the sample points  $X_i$ . We want a geometric characterization, not an algebraic one. Make sure your answer is complete; i.e., it includes all cases in which the covariance matrix of the sample is singular.
2. 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.
3. 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.
4. Suppose we have  $X \sim \mathcal{N}(0, \Sigma)$ ,  $X \in \mathbb{R}^n$  and a unit vector  $y \in \mathbb{R}^n$ . We can compute the projection of the random vector  $X$  onto a unit direction vector  $y$  as  $p = y^\top X$ . First, compute the variance of  $p$ . Second, with this information, what does the largest eigenvalue  $\lambda_{\max}$  of the covariance matrix tell us about the variances of expressions of the form  $y^\top X$ ?

## Section 5

1. Collinear samples: When sample points lie on the same hyperplane.

If the dimensions are  $>$  than observations: The cov matrix can't be full rank.

If features are redundant: If points are the same so that the # of unique points are  $< d$ , then there will be dependent columns in cov matrix.

Repeated Eigenvalues: If samples cause the eigenvalues to be repeated.

2. I would first try to reduce dimensionality, to reduce redundant dimensions. This may change the covar. matrix a lot, but we aren't really losing anything since those dimensions were redundant.

The other way (suggested in q8), is to add a small constant across the diagonal. This is limited by how fast the program runs. A smaller value makes the inverse have really big values, hence slowing down calculations.

3. Max f(x): Pick a  $x$  that is the unit vector of the eigen-vector with the largest correspondingly eigen-value of  $\Sigma$ . This is direction of max variance.

Min f(x): Similar to max f(x), but choose the smallest eigenvalue/eigenvector pair. This is the path of least variance.

4.  $\text{Var}(\hat{y}) = \text{Var}(y^T X) = y^T \Sigma y$

$\lambda_{\max}$  tells us about a direction with high variance. This means that many points are concentrated at on  $y^T X$ . If  $\lambda_{\max}$  is small, that means the variance along  $y^T X$  is small.

## 6 Isocontours of Normal Distributions

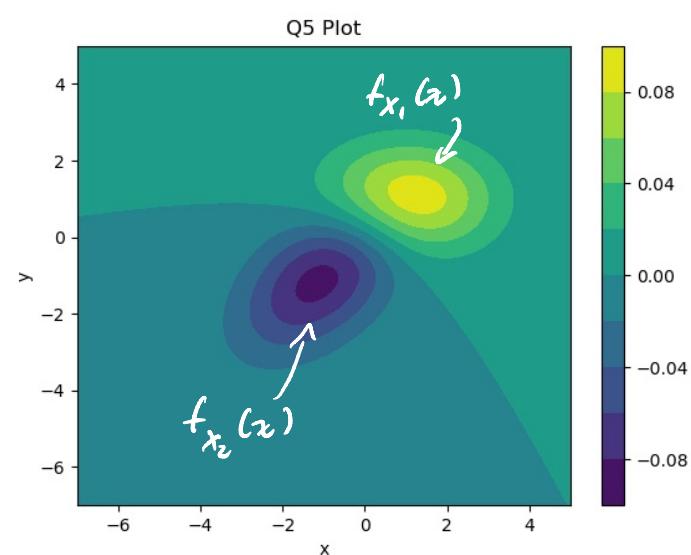
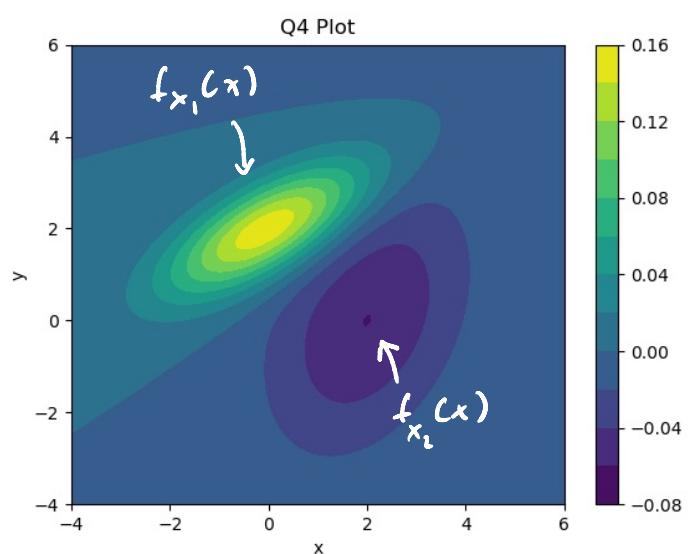
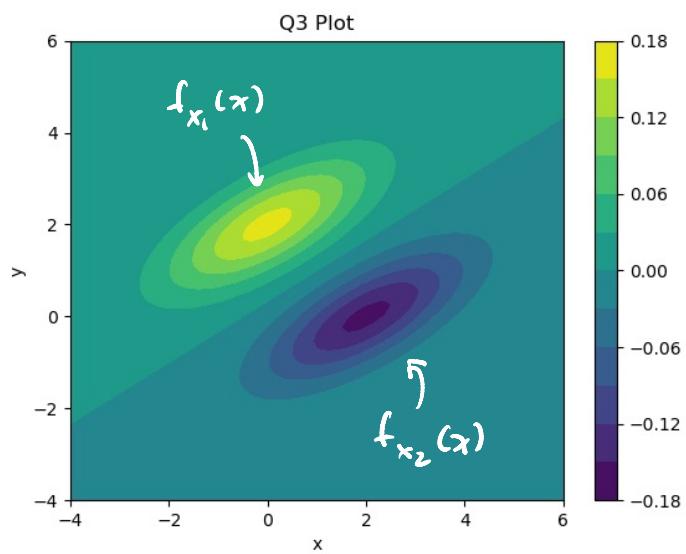
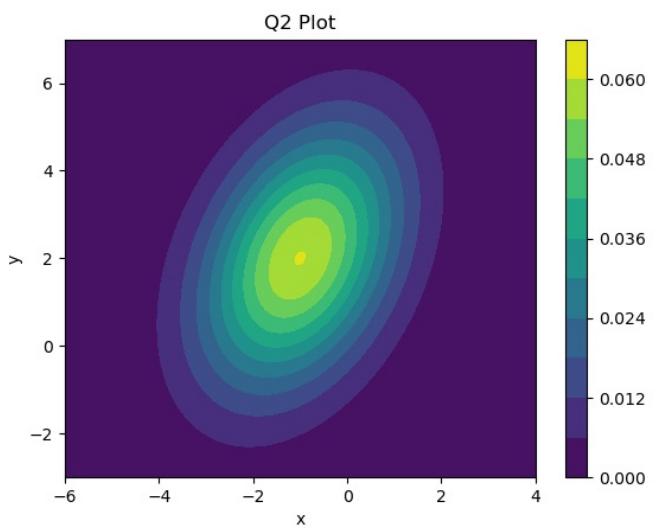
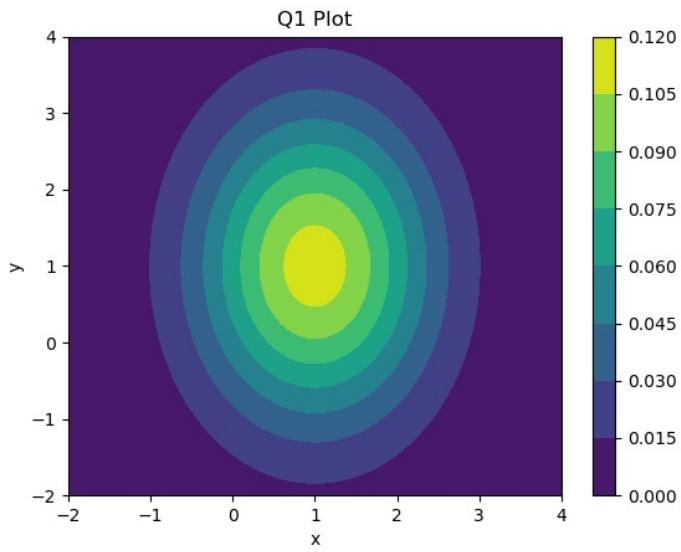
Let  $f(\mu, \Sigma)$  be the probability density function of a normally distributed random variable in  $\mathbb{R}^2$ . Write code to plot the isocontours of the following functions, each on its own separate figure. Make sure it is clear which figure belongs to which part. You're free to use any plotting libraries or stats utilities you like; for instance, in Python you can use Matplotlib and SciPy. Choose the boundaries of the domain you plot large enough to show the interesting characteristics of the isocontours (use your judgment). Make sure we can tell what isovalue each contour is associated with—you can do this with labels or a colorbar/legend.

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  $\Sigma_1 = \Sigma_2 = \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix}$ .

(Yes, this is a difference between two PDFs. No, it is not itself a valid PDF. Just plot its isocontours anyway.)

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  $\Sigma_2 = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$ .

# Section 6



## 7 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  $\sigma^2$ . (This means that you have to draw  $X_1$  first and use it to compute a random  $X_2$ .)

Write a program that draws  $n = 100$  random two-dimensional sample points from  $(X_1, X_2)$ . For each sample point, the value of  $X_2$  is a function of the value of  $X_1$  for that *same* sample point, but the sample points are independent of each other. 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.

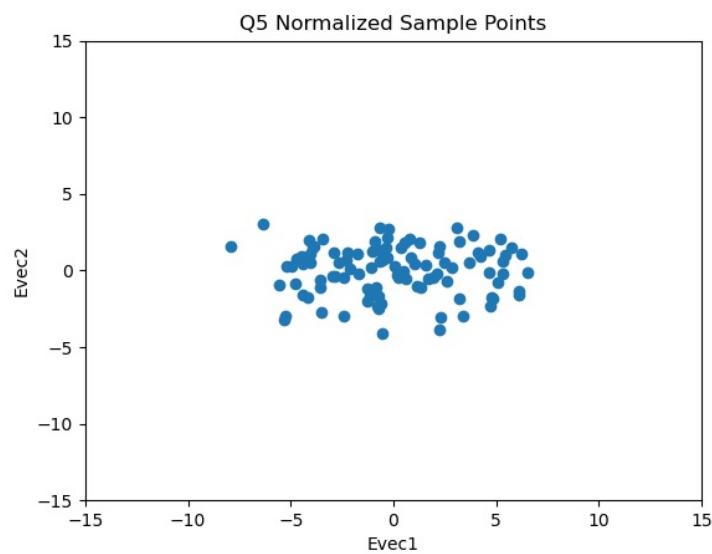
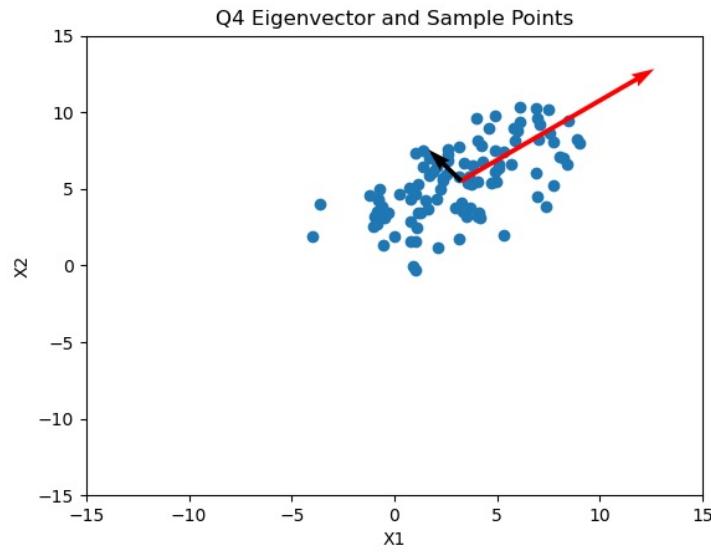
1. Compute the mean (in  $\mathbb{R}^2$ ) of the sample.
2. Compute the  $2 \times 2$  covariance matrix of the sample (based on the sample mean, not the true mean—which you would not know given real-world data).
3. Compute the eigenvectors and eigenvalues of this covariance matrix.
4. On a two-dimensional 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.

Hint: make *sure* your plotting software is set so the figure is square (i.e., the horizontal and vertical scales are the same). Not doing that may lead to hours of frustration!

5. Let  $U = [v_1 \ v_2]$  be a  $2 \times 2$  matrix whose columns are the **unit** eigenvectors of the covariance matrix, where  $v_1$  is the eigenvector with the larger eigenvalue. We use  $U^\top$  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  $U^\top = U^{-1}$ , 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  $U^\top$ , giving  $x_{\text{rotated}} = U^\top(x - \mu)$ . Plot these rotated points on a new two dimensional-grid, again with both axes having range  $[-15, 15]$ . (You are not required to plot the eigenvectors, which would be horizontal and vertical.)

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.

# Section 7



## 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. **Do NOT use the training/test data from Homework 1, as they have changed for this homework.** The starter code is similar to HW1's; we provide `check.py` and `save_csv.py` files for you to produce your Kaggle submission files. 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.

Reminder: please also select relevant code from the appendix on Gradescope for your answer to each question.

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 9 and Section 4.4 of *An Introduction to Statistical Learning*. Attach the relevant code as your answer to this part.

*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). Tell us which digit and include your visualization in your write-up. 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.

- (a) (Graphs) Linear discriminant analysis (LDA). Model the class conditional probabilities as Gaussians  $N(\mu_C, \Sigma)$  with different means  $\mu_C$  (for class C) and the same pooled within-class covariance matrix  $\Sigma$ , which you compute from a weighted average of the 10 covariance matrices from the 10 classes, as described in Lecture 9.

In your implementation, you might run into issues of determinants overflowing or underflowing, or normal PDF probabilities underflowing. These problems might be solved by learning about `numpy.linalg.slogdet` and/or `scipy.stats.multivariate_normal.logpdf`.

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 invertible. Instead, you should find a way to solve the implied linear system without computing the inverse.

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 unpredictability 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 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? Write down your answer and suggest why you think it's the easiest digit.

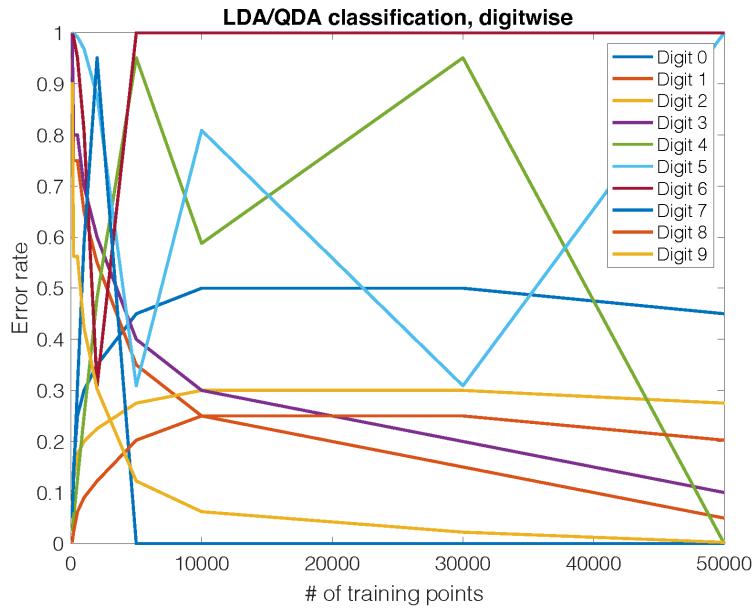


Figure 1: Sample graph with 10 plots

- 4. (Written answer) With `mnist-data-hw3.npz`, 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 write-up and include your Kaggle username. Don't forget to use the "submissions" tab or link on Kaggle to select your best submission!

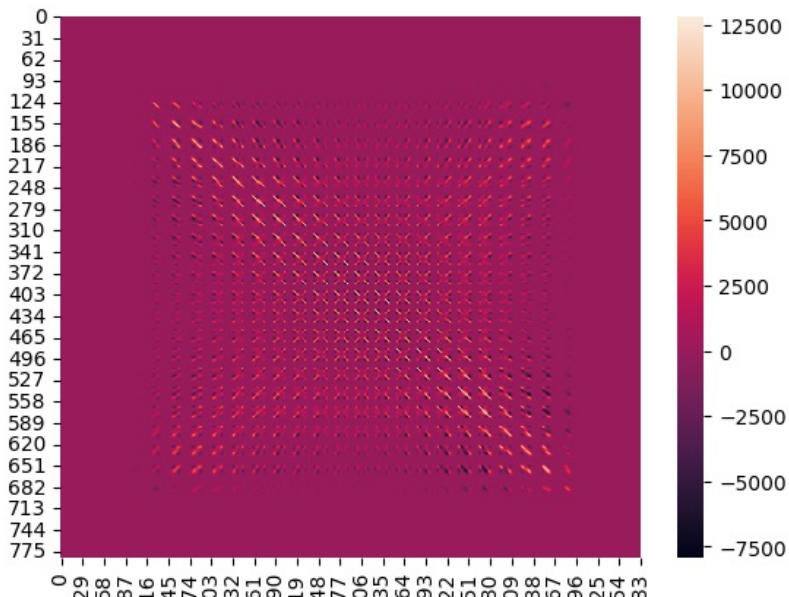
You are welcome to compute extra features for the Kaggle competition, as long as they do not use an exterior learned model for their computation (no transfer learning!). If you do so, please describe your implementation in your assignment. Please use extra features **only** for the Kaggle portion of the assignment.

5. (Written answer) Next, apply LDA or QDA (your choice) to spam (`spam-data-hw3.npz`). 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. We include a `featurize.py` file (similar to HW1's) that you may modify to create new features.

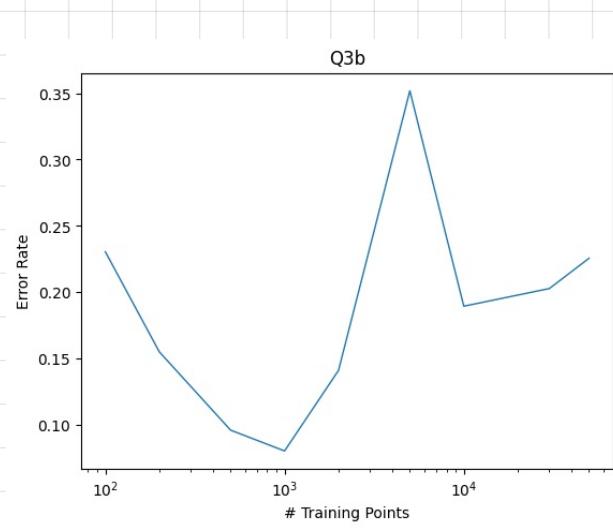
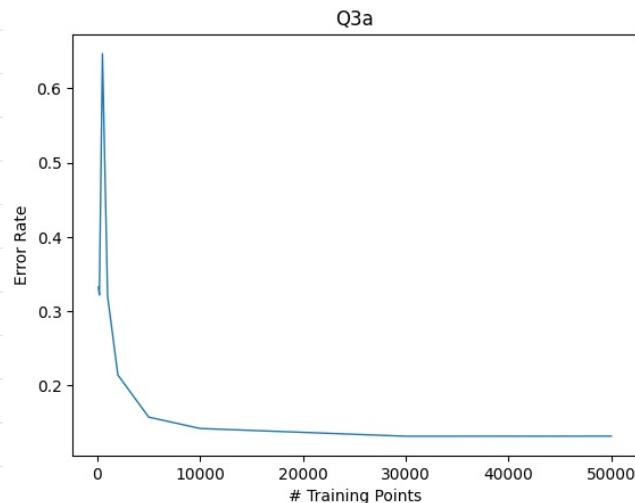
*Optional:* If you use the defaults, expect relatively low classification rates. We suggest using a Bag-Of-Words model. You are encouraged to explore alternative hand-crafted features, and are welcome to use any third-party library to implement them, as long as they do not use a separate model for their computation (no large language models, BERT, or word2vec!).

2. Looked at the covariance matrix with 0

The diagonal terms seem to be larger than the other terms. Indicating a strong signal



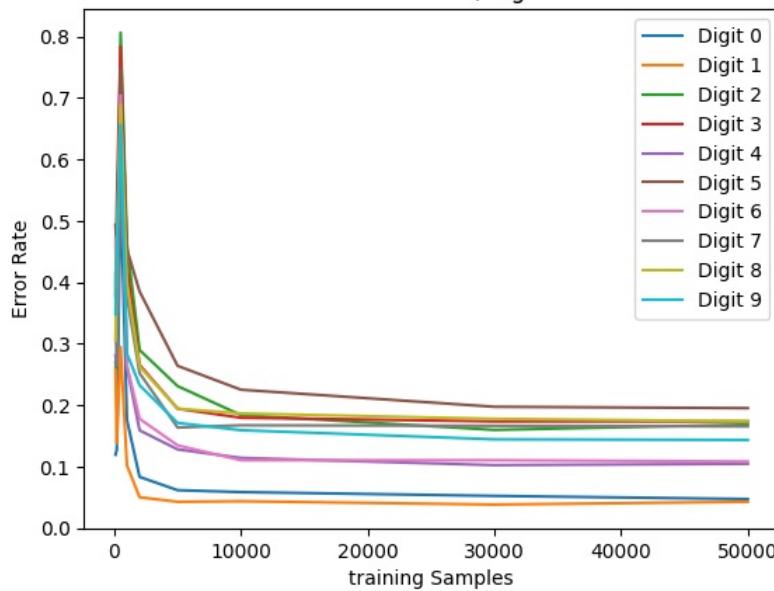
3.



c) QDA model was better at around 1000 training points. Then overfit.  
LDA never overfit but was not as good.

d)

LDA classification, digitwise



Digits that have a more common way of writing and are less complicated are easier to classify.

This is b/c they are less variance in what they look like.

4. kaggle username: Frank Jin (fj0228)

Mnist: 0.914

5. kaggle username: Frank Jin (fj0228)

SPAM: 0.793

## Submission Checklist

Please ensure you have completed the following before your final submission.

At the beginning of your writeup...

1. Have you copied and hand-signed the honor code specified in Question 1?
2. Have you listed all students (Names and ID numbers) that you collaborated with?

In your writeup for Question 8...

1. Have you included your **Kaggle Score** and **Kaggle Username** for **both** questions 8.4 and 8.5?

At the end of the writeup...

1. Have you provided a code appendix including all code you wrote in solving the homework?
2. Have you included featurize.py in your code appendix if you modified it?

## Executable Code Submission

1. Have you created an archive containing all “.py” files that you wrote or modified to generate your homework solutions (including featurize.py if you modified it)?
2. Have you removed all data and extraneous files from the archive?
3. Have you included a README file in your archive briefly describing how to run your code on the test data and reproduce your Kaggle results?

## Submissions

1. Have you submitted your test set predictions for both **MNIST** and **SPAM** to the appropriate Kaggle challenges?
2. Have you submitted your written solutions to the Gradescope assignment titled **HW3 Write-Up** and selected pages appropriately?
3. Have you submitted your executable code archive to the Gradescope assignment titled **HW3 Code**?

Congratulations! You have completed Homework 3.