

**ENG EC 414 (Ishwar) Introduction to Machine Learning**

**HW 4**

© 2015 – Spring 2022 Prakash Ishwar

**Issued:** Fri 4 Mar 2022      **Due:** 10:55pm Boston time Fri 18 Mar 2022 in [Gradescope](#) + [Blackboard](#).  
**Required reading:** Slides on LDA + OLS/ridge regression + your notes from lectures & Discussions 5, 6.

---

**Important:** Before you proceed, please read the documents pertaining to *Homework formatting and submission guidelines* and the *HW-grading policies* in the Homeworks section of Blackboard, especially guidelines for submitting [reports in Gradescope](#), and [code in Blackboard](#). In particular, for computer assignments *you are prohibited from using any online code or built-in MATLAB functions except as indicated in the problem or skeleton code (when provided)*.

In order to receive full credit, all work should be supported by a concise explanation that is clear, relevant, specific, logical, and correct. In particular, for each part, you must clearly outline the key steps and provide proper justification for your calculations.

**Note:** Problem difficulty = number of coffee cups ☕

**Problem 4.1** [15pts] (*LDA could perform poorly*) Consider a dataset  $\mathcal{D}$  of 40 points in  $\mathbb{R}^2$  with half the points in the positive class  $\mathcal{D}_+$  and the remaining in the negative class  $\mathcal{D}_-$ . Specifically, let

$$\mathcal{D}_- = \left\{ 9 \text{ points at } \begin{pmatrix} -1 \\ 0.5 \end{pmatrix}, \text{ one point at } \begin{pmatrix} -11 \\ 0.5 \end{pmatrix}, 9 \text{ points at } \begin{pmatrix} -1 \\ -1.5 \end{pmatrix}, \text{ one point at } \begin{pmatrix} -11 \\ -1.5 \end{pmatrix} \right\}$$

$$\mathcal{D}_+ = \left\{ 9 \text{ points at } \begin{pmatrix} 1 \\ 1.5 \end{pmatrix}, \text{ one point at } \begin{pmatrix} 11 \\ 1.5 \end{pmatrix}, 9 \text{ points at } \begin{pmatrix} 1 \\ -0.5 \end{pmatrix}, \text{ one point at } \begin{pmatrix} 11 \\ -0.5 \end{pmatrix} \right\}$$

- (a) Compute the mean vector and covariance matrix of each class and the average of the two class covariance matrices, i.e., compute  $\hat{\mu}_{x-}, \hat{\mu}_{x+}, \hat{S}_{x-}, \hat{S}_{x+}, \hat{S}_{x\text{-avg}}$ .
- (b) Compute  $\mathbf{w}_{\text{LDA}}, b_{\text{LDA}}$ , and the CCR of the LDA classifier.
- (c) Compute a linear binary classifier with the highest CCR and compare it with the LDA classifier.

**Problem 4.2** [20pts] (*Linear Least Squares Regression*) Consider the following training set  $\mathcal{D} = \{(x_1, y_1) = (-1, -1), (x_2, y_2) = (-1/2, -1/8), (x_3, y_3) = (0, 0), (x_4, y_4) = (1/2, 1/8), (x_5, y_5) = (1, 1)\}$ . Hand-compute the following:

- (a) [6pts] *Ordinary Least Squares*:  $(w_{OLS}, b_{OLS}) = \arg \min_{w,b} \sum_{j=1}^5 (y_j - wx_j - b)^2$ .
- (b) [14pts] *Polynomial Least Squares*:  $(\mathbf{w}_{PLS}, b_{PLS}) = \arg \min_{\mathbf{w},b} \sum_{j=1}^5 (y_j - \mathbf{w}^\top \boldsymbol{\phi}_j - b)^2$ . Where  $\boldsymbol{\phi}_j = ((x_j)^3, (x_j)^2)^\top$  and  $\mathbf{w} = (w_3, w_2)^\top$ .

**Problem 4.3** [43pts] (*LDA*) In this problem we will implement LDA on a synthetic dataset and also develop geometric intuition for the eigenvalues and eigenvectors of the empirical feature covariance matrices. You are provided skeleton code to assist you in implementing LDA.

- (a) [13pts] (*Synthetic dataset generation*) Write a Matlab function to generate as its output a 2-class labeled dataset consisting of 2D feature vectors each drawn independently from one of two Gaussian distributions (one distribution for each class). The function's inputs are (1) a specified number of examples  $n_1, n_2$  for each distribution (class), (2) specified  $2 \times 1$  mean vectors  $\mu_1, \mu_2$  for each class, and (3) a single common  $2 \times 2$  covariance matrix for both classes specified via two real, nonnegative eigenvalues  $\lambda_1, \lambda_2 \geq 0$ , and a single orientation variable  $\theta$  that defines their corresponding orthonormal eigenvectors  $\mathbf{u}_1 = (\cos \theta, \sin \theta)^\top$  and  $\mathbf{u}_2 = (\sin \theta, -\cos \theta)^\top$ .

Use this function to generate  $n_1 = 50$  class 1 examples and  $n_2 = 100$  class 2 examples with  $\mu_1 = (1, 2)^\top, \mu_2 = (3, 2)^\top$  and the following four choices of the  $\lambda$ 's and  $\theta$ : (i)  $\lambda_1 = 1, \lambda_2 = 0.25, \theta = 0$ , (ii)  $\lambda_1 = 1, \lambda_2 = 0.25, \theta = \pi/6$ , (iii)  $\lambda_1 = 1, \lambda_2 = 0.25, \theta = \pi/3$ , (iv)  $\lambda_1 = 0.25, \lambda_2 = 1, \theta = \pi/6$ . For each of these four choices, create a scatter plot with class 1 examples shown as solid blue circles and class 2 examples shown as solid red triangles. Discuss how the eigenvalues and eigenvectors affect the geometry of the dataset.

**Note:** For the rest of this problem, we will work with the dataset created by the **second** choice, i.e., choice number (ii), out of the four choices in part (a).

- (b) [17pts] Write a Matlab function which takes as inputs a 2D labeled dataset for binary classification and a unit direction (specified by an angle  $\phi$  in radians) and returns the following three outputs (1) the squared Euclidean distance between the class means of the projections of the feature vectors along the direction  $(\cos \phi, \sin \phi)^\top$ , i.e.,  $(\hat{\mu}_{2z}(\phi) - \hat{\mu}_{1z}(\phi))^2$ , where  $\hat{\mu}_{1z}(\phi)$  and  $\hat{\mu}_{2z}(\phi)$  are, respectively, the class 1 and class 2 empirical means along direction  $\phi$ , (2) the average within-class variance of the projections of the feature vectors along direction  $(\cos \phi, \sin \phi)^\top$ , i.e.,  $\frac{n_1}{n} \hat{\sigma}_{1z}^2(\phi) + \frac{n_2}{n} \hat{\sigma}_{2z}^2(\phi)$ , where  $\hat{\sigma}_{1z}^2(\phi)$  and  $\hat{\sigma}_{2z}^2(\phi)$  are, respectively, the class 1 and class 2 empirical variances along direction  $\phi$ , and (3) their ratio, i.e., Signal-to-Noise Ratio (SNR).

Use this function to plot all three quantities as a function of  $\phi$ , for  $\phi$  ranging from 0 to  $\pi$  in steps of  $\pi/48$ , for the dataset of part(a)(ii). Determine (i) the value of  $\phi$  which maximizes the squared distance between the class means of the projected values (ii) the value of  $\phi$  which minimizes the average within-class variance, and (iii) the value of  $\phi$  which maximizes the SNR. (iv) Also use Matlab's `ksdensity` function with default settings to plot estimates of the class 1 and class 2 probability densities of the projections of the feature vectors for  $\phi = 0, \pi/6, \pi/3$ . Discuss your findings.

- (c) [6pts] Write a Matlab function which takes as input a 2D labeled dataset for binary classification and outputs the LDA vector  $\mathbf{w}_{LDA} = \hat{S}_{x,\text{avg}}^{-1}(\hat{\mu}_{2x} - \hat{\mu}_{1x})$  where  $\hat{S}_{x,\text{avg}} = \frac{n_1}{n} \hat{S}_{1x} + \frac{n_2}{n} \hat{S}_{2x}$  and  $\hat{S}_{1x}$  and  $\hat{S}_{2x}$  are the empirical  $2 \times 2$  covariance matrices of classes 1 and 2 respectively.

Use this function to compute  $\mathbf{w}_{LDA}$  for the dataset of part(a)(ii). Compare it with the difference between the class 2 and class 1 empirical mean vectors  $(\hat{\mu}_{2x} - \hat{\mu}_{1x})$ . Overlay both these vectors on the scatter plot of the dataset with the vectors represented as arrows starting at the location of  $\hat{\mu}_{1x}$ . Discuss what you observe. Also compare the direction of  $\mathbf{w}_{LDA}$  with the value of  $\phi$  from part (b) which maximizes the SNR.

- (d) [7pts] Consider the following separating hyperplane decision rule for binary classification based on thresholding a linear (affine) function of the feature vector:  $h_{\mathbf{w}, b}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}^\top \mathbf{x} + b \leq 0 \\ 2 & \text{else} \end{cases}$ . Write a Matlab function which takes as inputs a 2D labeled dataset for binary classification,  $\mathbf{w}$ , and  $b$  as inputs and outputs the value of CCR that results from applying  $h_{\mathbf{w}, b}(\mathbf{x})$  to classify all feature vectors in the dataset into class 1 or class 2.

↑  
NOT  
TOO  
BAD

Use this function to create a plot of CCR as a function of  $b$  for the dataset of part(a)(ii) and  $\mathbf{w} = \mathbf{w}_{LDA}$ . Compute the value of the offset parameter  $b$  which maximizes the CCR and the resulting CCR value.

**Problem 4.4** [22pts] (*ridge regression*) In this problem we will implement ridge regression and apply it to a real-world 8-dimensional (8 features) prostate cancer dataset contained in the file `prostateStnd.mat`. In this dataset, 8 medically relevant features named `lcavol`, `lweight`, `age`, `lbph`, `svi`, `lcp`, `gleason`, and `pgg45` are used to estimate `lpsa` (log prostate specific antigen). The training and test data are provided as `(xtrain, ytrain)` and `(xtest, ytest)` respectively. The first 8 features correspond to the first 8 entries of `names`. The ninth entry of `names` (the last one) is the label to be predicted whose values are in `(ytest, ytrain)`.

- offset = mean  
Scaling = Std  
reprint b*
- (a) [4pts] As a first step, write Matlab code to normalize the **training dataset** so that post-normalization, each of the 8 features and the label in the normalized training dataset has zero mean and unit variance. This requires determining a pair of offset and scaling parameters, one pair for each feature and one pair for the label. These parameters must be computed only from the training dataset, but they must be applied to both the training and test datasets, i.e., we normalize both the training and test data, but we are only allowed to normalize the test data using parameters derived from the training data. In other words we must apply identical operations to training and test data. Only the training data will be actually normalized by the operation. If the test data is statistically similar to the training data, it too will be approximately normalized. Report the mean and variance of each feature and the label before normalization.
  - (b) [6pts] Next, write Matlab code to use the normalized data to train a ridge regression model for each of the following values of the quadratic regularization penalty parameter  $\lambda$ :  $\{e^{-5}, e^{-4}, e^{-3}, \dots, e^{10}\}$ .
  - (c) [9pts] In a single figure, plot the ridge regression coefficient of each feature (8 in total) as a function of  $\ln \lambda$  (8 curves in total) for  $\ln \lambda$  ranging from -5 to 10 in steps of 1. Use suitable colors and/or markers to distinguish between the 8 curves and label them appropriately in a legend. Discuss what happens to the coefficients as  $\lambda$  becomes larger.
  - (d) [3pts] In another figure, plot the mean-squared-error (MSE) of both the training and the test data as a function of  $\ln \lambda$ . Discuss your observations.

**Code-submission via Blackboard:** Create two dot-m files. One named `<yourBUemailID>.hw4_3.m` for Problem 4.3 and one named `<yourBUemailID>.hw4_4.m` for Problem 4.4. All local functions and scripts pertaining to one problem should appear within the single dot-m file for that problem. When run, your scripts should be able to display in the command window whatever you are asked to compute and report in each part. One skeleton code is provided for your reference: `skeleton.hw4_3.m`. Reach-out to the TA via Piazza and office hours for questions related to coding. When submitting code, please include a ‘Readme’ text file in your source code directory describing approximate running times of different parts, any additional comments (as needed) explaining how to use your code, and any dependencies between different parts. Please do not include into the directory, any data files that are already provided. Write your code under the assumption that all data files are in the same directory as your source code.

**Problem 4.1** [15pts] (LDA could perform poorly) Consider a dataset  $\mathcal{D}$  of 40 points in  $\mathbb{R}^2$  with half the points in the positive class  $\mathcal{D}_+$  and the remaining in the negative class  $\mathcal{D}_-$ . Specifically, let

$$\mathcal{D}_- = \left\{ 9 \text{ points at } \begin{pmatrix} -1 \\ 0.5 \end{pmatrix}, \text{ one point at } \begin{pmatrix} -11 \\ 0.5 \end{pmatrix}, 9 \text{ points at } \begin{pmatrix} -1 \\ -1.5 \end{pmatrix}, \text{ one point at } \begin{pmatrix} -11 \\ -1.5 \end{pmatrix} \right\}$$

$$\mathcal{D}_+ = \left\{ 9 \text{ points at } \begin{pmatrix} 1 \\ 1.5 \end{pmatrix}, \text{ one point at } \begin{pmatrix} 11 \\ 1.5 \end{pmatrix}, 9 \text{ points at } \begin{pmatrix} 1 \\ -0.5 \end{pmatrix}, \text{ one point at } \begin{pmatrix} 11 \\ -0.5 \end{pmatrix} \right\}$$

- (a) Compute the mean vector and covariance matrix of each class and the average of the two class covariance matrices, i.e., compute  $\hat{\mu}_{x-}, \hat{\mu}_{x+}, \hat{S}_{x-}, \hat{S}_{x+}, \hat{S}_{x-\text{avg}}$ .
- (b) Compute  $\mathbf{w}_{\text{LDA}}, b_{\text{LDA}}$ , and the CCR of the LDA classifier.
- (c) Compute a linear binary classifier with the highest CCR and compare it with the LDA classifier.

$$a) \quad \mathbf{x}_- = \begin{vmatrix} 9(-1) & -11 & 9(-1) & -11 \\ 9(1/2) & 1/2 & 9(-3/2) & -3/2 \end{vmatrix}$$

$$\bar{x}_1 = \begin{vmatrix} -1 \\ 1/2 \end{vmatrix} - \begin{vmatrix} -2 \\ -1/2 \end{vmatrix}$$

$$\bar{x}_1 = \begin{vmatrix} 1 \\ 1 \end{vmatrix}$$

$$\bar{x}_1 \cdot \bar{x}_1^\top = \begin{vmatrix} 1 \\ 1 \end{vmatrix} \cdot \begin{vmatrix} 1 \\ 1 \end{vmatrix} = \begin{vmatrix} 1 & 1 \\ 1 & 1 \end{vmatrix}$$

$$q(\bar{x}_1 \cdot \bar{x}_1^\top) = \begin{vmatrix} q & q \\ q & q \end{vmatrix}$$

$$\bar{x}_2 = \begin{vmatrix} -11 \\ 1/2 \end{vmatrix} - \begin{vmatrix} -2 \\ -1/2 \end{vmatrix}$$

$$\bar{x}_2 = \begin{vmatrix} -9 \\ 1 \end{vmatrix}$$

$$\bar{x}_2 \cdot \bar{x}_2^\top = \begin{vmatrix} -9 \\ 1 \end{vmatrix} \cdot \begin{vmatrix} -9 \\ 1 \end{vmatrix} = \begin{vmatrix} 81 & -9 \\ -9 & 1 \end{vmatrix}$$

$$\bar{x}_3 = \begin{vmatrix} -1 \\ -3/2 \end{vmatrix} - \begin{vmatrix} -2 \\ -1/2 \end{vmatrix}$$

$$\bar{x}_3 = \begin{vmatrix} 1 \\ -1 \end{vmatrix}$$

$$\bar{x}_3 \cdot \bar{x}_3^\top = \begin{vmatrix} 1 \\ -1 \end{vmatrix} \cdot \begin{vmatrix} 1 \\ -1 \end{vmatrix} = \begin{vmatrix} 1 & -1 \\ -1 & 1 \end{vmatrix}$$

$$q(\bar{x}_3 \cdot \bar{x}_3^\top) = \begin{vmatrix} q & -q \\ -q & q \end{vmatrix}$$

remember  $\hat{\mu} = \frac{1}{n} \sum_{j: y_j=1} x_j$

$$\hat{\mu}_- = \frac{1}{20} \begin{vmatrix} -40 \\ -10 \end{vmatrix} = \begin{vmatrix} -2 \\ -1/2 \end{vmatrix}$$

$$\hat{\mu}_- = \begin{vmatrix} -2 \\ -1/2 \end{vmatrix}$$

$$\mathbf{x}_+ = \begin{vmatrix} 9(1) & 11 & 9(1) & 11 \\ 9(3/2) & 3/2 & 9(-1/2) & -1/2 \end{vmatrix}$$

$$\hat{\mu}_+ = \frac{1}{20} \begin{vmatrix} 40 \\ 10 \end{vmatrix} = \begin{vmatrix} 2 \\ 1/2 \end{vmatrix}$$

$$\hat{\mu}_+ = \begin{vmatrix} 2 \\ 1/2 \end{vmatrix}$$

$$S = \frac{1}{n} \sum_{j=1}^n \bar{x}_j \cdot \bar{x}_j^\top$$

$$\bar{x}_4 = \begin{vmatrix} -1 \\ 1 \\ 3/2 \end{vmatrix} - \begin{vmatrix} -2 \\ -1 \\ 1/2 \end{vmatrix}$$

$$\bar{x}_4 = \begin{vmatrix} -9 \\ 1 \end{vmatrix}$$

$$\bar{x}_4 \cdot \bar{x}_4^T = \begin{vmatrix} -9 \\ 1 \end{vmatrix} \cdot \begin{vmatrix} -9 & 1 \\ 1 & 1 \end{vmatrix} = \begin{vmatrix} 81 & -9 \\ -9 & 1 \end{vmatrix}$$

$$\sum_{j=1}^2 \bar{x}_j \cdot \bar{x}_j^T = \begin{vmatrix} 9 & 1 \\ 1 & 1 \end{vmatrix} + \begin{vmatrix} 81 & -9 \\ -9 & 1 \end{vmatrix}$$

$$+ \begin{vmatrix} 9 & -1 \\ -1 & 1 \end{vmatrix} + \begin{vmatrix} 81 & 1 \\ 1 & 1 \end{vmatrix}$$

$$= \begin{vmatrix} 9+81+9+81 & 0 \\ 0 & 20 \end{vmatrix}$$

$$= \begin{vmatrix} 180 & 0 \\ 0 & 20 \end{vmatrix}$$

$$\hat{S}_- = \frac{1}{20} \begin{vmatrix} 180 & 0 \\ 0 & 20 \end{vmatrix} = \begin{vmatrix} 9 & 0 \\ 0 & 1 \end{vmatrix}$$

$$\hat{S}_- = \begin{vmatrix} 9 & 0 \\ 0 & 1 \end{vmatrix}$$

$$x_+ = \begin{vmatrix} 9(1) & 11 & 9(1) & 11 \\ 9(3/2) & 1/2 & 9(-1/2) & -1/2 \end{vmatrix}$$

$$\bar{x}_1 = \begin{vmatrix} 1 \\ 3/2 \end{vmatrix} - \begin{vmatrix} 2 \\ 1/2 \end{vmatrix}$$

$$\bar{x}_1 = \begin{vmatrix} -1 \\ 1 \end{vmatrix}$$

$$\bar{x}_1 \cdot \bar{x}_1^T = \begin{vmatrix} -1 \\ 1 \end{vmatrix} \cdot \begin{vmatrix} -1 & 1 \\ 1 & 1 \end{vmatrix} = \begin{vmatrix} 1 & -1 \\ -1 & 1 \end{vmatrix}$$

$$q(\bar{x}_1 \cdot \bar{x}_1^T) = \begin{vmatrix} 9 & -9 \\ -9 & 1 \end{vmatrix}$$

$$\bar{x}_2 = \begin{vmatrix} 11 \\ 3/2 \end{vmatrix} - \begin{vmatrix} 2 \\ 1/2 \end{vmatrix} = \begin{vmatrix} 9 \\ 1 \end{vmatrix}$$

$$\bar{x}_2 \cdot \bar{x}_2^T = \begin{vmatrix} 9 \\ 1 \end{vmatrix} \cdot \begin{vmatrix} 9 & 1 \\ 1 & 1 \end{vmatrix} = \begin{vmatrix} 81 & 9 \\ 9 & 1 \end{vmatrix}$$

$$\bar{x}_3 = \begin{vmatrix} 1 \\ -1/2 \end{vmatrix} - \begin{vmatrix} 2 \\ 1/2 \end{vmatrix} = \begin{vmatrix} -1 \\ -1 \end{vmatrix}$$

$$\bar{x}_3 \cdot \bar{x}_3^T = \begin{vmatrix} -1 \\ -1 \end{vmatrix} \cdot \begin{vmatrix} -1 & -1 \\ -1 & -1 \end{vmatrix} = \begin{vmatrix} 1 & 1 \\ 1 & 1 \end{vmatrix}$$

$$q(\bar{x}_3 \cdot \bar{x}_3^T) = \begin{vmatrix} 9 & 9 \\ 9 & 9 \end{vmatrix}$$

$$\bar{x}_4 = \begin{vmatrix} 11 \\ -1/2 \end{vmatrix} - \begin{vmatrix} 2 \\ 1/2 \end{vmatrix} = \begin{vmatrix} 9 \\ -1 \end{vmatrix}$$

$$\bar{x}_4 \cdot \bar{x}_4^T = \begin{vmatrix} 9 \\ -1 \end{vmatrix} \cdot \begin{vmatrix} 9 & -1 \\ -1 & 1 \end{vmatrix} = \begin{vmatrix} 81 & -9 \\ -9 & 1 \end{vmatrix}$$

$$\sum_{j=1}^2 \bar{x}_j \cdot \bar{x}_j^T = \begin{vmatrix} 9 & 1 \\ 1 & 1 \end{vmatrix} + \begin{vmatrix} 81 & -9 \\ -9 & 1 \end{vmatrix}$$

$$+ \begin{vmatrix} 9 & -1 \\ -1 & 1 \end{vmatrix} + \begin{vmatrix} 81 & 1 \\ 1 & 1 \end{vmatrix}$$

$$= \begin{vmatrix} 9+81+9+81 & 0 \\ 0 & 20 \end{vmatrix}$$

$$= \begin{vmatrix} 180 & 0 \\ 0 & 20 \end{vmatrix}$$

$$\hat{S}_+ = \frac{1}{20} \begin{vmatrix} 180 & 0 \\ 0 & 20 \end{vmatrix} = \begin{vmatrix} 9 & 0 \\ 0 & 1 \end{vmatrix}$$

$$\hat{S}_+ = \begin{vmatrix} 9 & 0 \\ 0 & 1 \end{vmatrix}$$

$$\hat{P}_- = 20/40 = 1/2$$

$$\hat{P}_+ = 1/2$$

$$\hat{S}_{x \cdot \text{avg}} = 1/2 \begin{vmatrix} 9 & 0 \\ 0 & 1 \end{vmatrix} + 1/2 \begin{vmatrix} 9 & 0 \\ 0 & 1 \end{vmatrix} \quad \text{Need to project } w^T \text{ to even point.}$$

$$\hat{S}_{x \cdot \text{avg}} = \begin{vmatrix} 9 & 0 & 1 \\ 0 & 1 & 1 \end{vmatrix}$$

$$z_i = w^T x_i$$

$$(a) + z_A = |4/9| 1 \begin{vmatrix} -1 \\ 0.5 \end{vmatrix} = -4/9 + 1/2 =$$

$$(a) + z_B = 1/8 = 0.05$$

$$+ z_B = |4/9| 1 \begin{vmatrix} -1 \\ 1/2 \end{vmatrix} = -11(4/9) + 1/2$$

$$b) (\hat{S}_{x \cdot \text{avg}})^{-1} = \frac{1}{9} \begin{vmatrix} 1 & 0 \\ 0 & 9 \end{vmatrix} \quad + z_B = -4.38$$

$$= \begin{vmatrix} 1/9 & 0 \\ 0 & 1 \end{vmatrix} \quad (a) + z_C = |4/9| 1 \begin{vmatrix} -1 \\ -3/2 \end{vmatrix} = 4/9 + -3/2$$

$$(a) + z_C = -1.05$$

$$\hat{\mu}_{x+} - \hat{\mu}_{x-} = \begin{vmatrix} 2 \\ 1/2 \end{vmatrix} - \begin{vmatrix} -2 \\ -1/2 \end{vmatrix} \quad + z_D = |4/9| 1 \begin{vmatrix} -1 \\ -3/2 \end{vmatrix} = 4/9(-1) + -3/2$$

$$= \begin{vmatrix} 4 \\ 1 \end{vmatrix} \quad + z_D = -6.38$$

$$(a) - z_E = |4/9| 1 \begin{vmatrix} 1 \\ 3/2 \end{vmatrix} = 4/9 + 3/2$$

$$(a) - z_E = 1.94$$

$$W_{LDA} = \begin{vmatrix} 1/9 & 0 \\ 0 & 1 \end{vmatrix} \cdot \begin{vmatrix} 4 \\ 1 \end{vmatrix} - z_F = |4/9| 1 \begin{vmatrix} 1 \\ 3/2 \end{vmatrix} = 4/9(1) + 3/2$$

$$2 \times 2 \quad 2 \times 1 \quad - z_F = 6.38$$

$$W_{LDA} = \begin{vmatrix} 4/9 \\ 1 \end{vmatrix} \quad (a) - z_G = |4/9| 1 \begin{vmatrix} 1 \\ -1/2 \end{vmatrix} = 4/9 + -1/2$$

$$(a) - z_G = -0.05$$

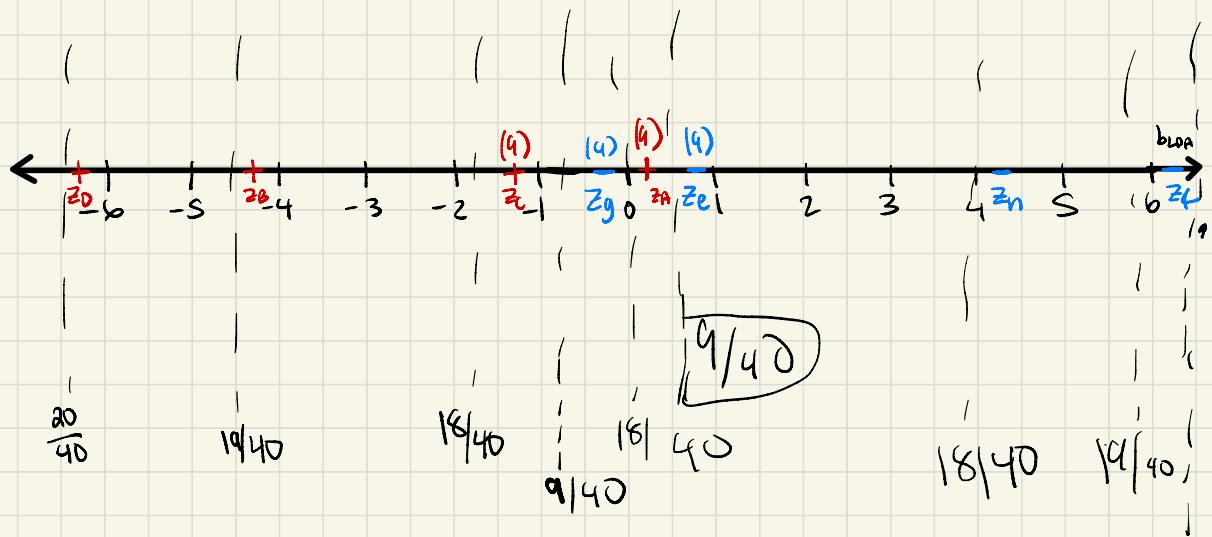
$$z_h = 4|a| \begin{pmatrix} 1 \\ -1/2 \end{pmatrix} = 4|a|(1) + -1/2$$

$$z_n = 4 \cdot 3\bar{8}$$

$$\bar{z} \leftarrow \begin{smallmatrix} + \\ \rightarrow \end{smallmatrix} \quad z \leftarrow \begin{smallmatrix} + \\ \rightarrow \end{smallmatrix}$$

$$1 \leftarrow \begin{smallmatrix} + \\ \rightarrow \end{smallmatrix}$$

$$\bar{z} \leftarrow \begin{smallmatrix} + \\ \rightarrow \end{smallmatrix} \quad z \leftarrow \begin{smallmatrix} + \\ \rightarrow \end{smallmatrix}$$



Best CCR = 1/2 (half right / half wrong)

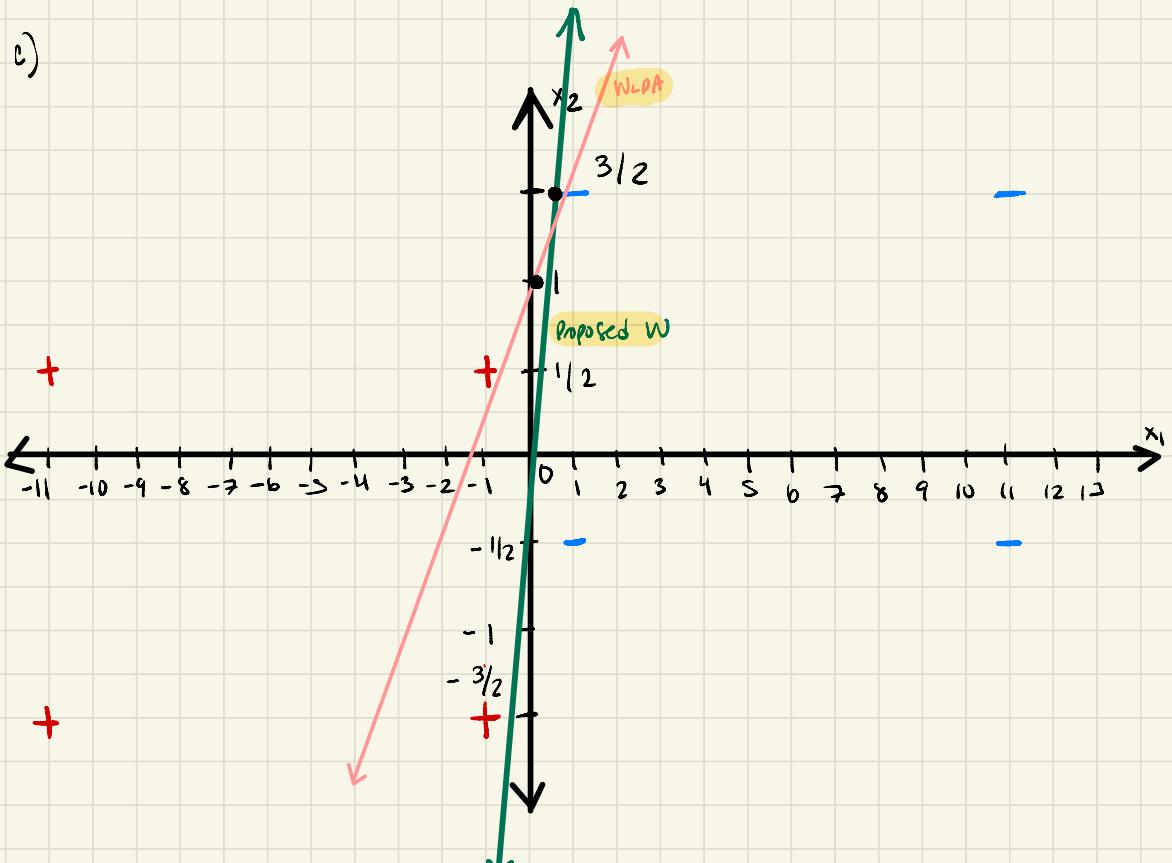
$$b_{LoA} = \begin{cases} b_{LoA} > b \cdot 38 \\ b_{LoA} < -b \cdot 38 \end{cases}$$

6 else

$$b_{LoA} = [-8, 8]$$

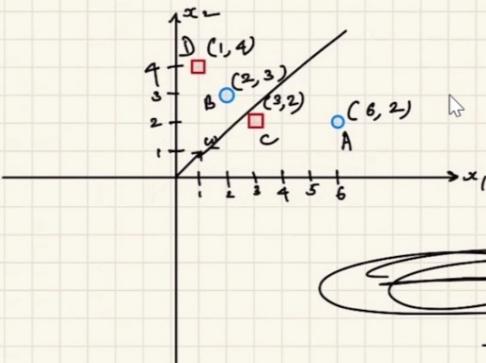
$$20/40$$

c)



$$W = \begin{bmatrix} 1/2 \\ 3/2 \end{bmatrix}$$

all points are correct  $CLR = 1$



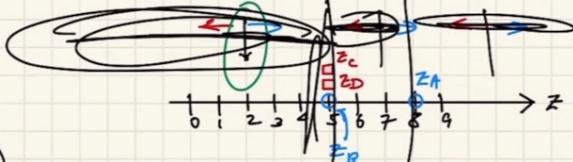
$$\underline{w} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$z_A = \underline{w}^T \underline{x}_A = (1 \ 1) \begin{pmatrix} 6 \\ 2 \end{pmatrix} = 8$$

$$z_B = \underline{w}^T \underline{x}_B = (1 \ 1) \begin{pmatrix} 2 \\ 3 \end{pmatrix} = 5$$

$$z_C = \underline{w}^T \underline{x}_C = (1 \ 1) \begin{pmatrix} 3 \\ 2 \end{pmatrix} = 5$$

$$z_D = \underline{w}^T \underline{x}_D = (1 \ 1) \begin{pmatrix} 1 \\ 4 \end{pmatrix} = 5$$



$$h(x) = \text{sign}(\underline{w}^T \underline{x} + b)$$

$$= \begin{cases} + & \text{if } \underline{w}^T \underline{x} + b > 0 \\ - & \text{if } \underline{w}^T \underline{x} + b < 0 \end{cases}$$

$$CCR = \frac{1}{2} = \frac{1}{4}$$

$$\left( \frac{3}{4} \right)$$

$$\frac{2}{4} = \frac{1}{2}$$

**Problem 4.2** [20pts] (Linear Least Squares Regression) Consider the following training set  $\mathcal{D} = \{(x_1, y_1) = (-1, -1), (x_2, y_2) = (-1/2, -1/8), (x_3, y_3) = (0, 0), (x_4, y_4) = (1/2, 1/8), (x_5, y_5) = (1, 1)\}$ . Hand-compute the following:

(a) [6pts] Ordinary Least Squares:  $(w_{OLS}, b_{OLS}) = \arg \min_{w,b} \sum_{j=1}^5 (y_j - w x_j - b)^2$ .

(b) [14pts] Polynomial Least Squares:  $(w_{PLS}, b_{PLS}) = \arg \min_{w,b} \sum_{j=1}^5 (y_j - \mathbf{w}^\top \phi_j - b)^2$ . Where  $\phi_j = ((x_j)^3, (x_j)^2)^\top$  and  $\mathbf{w} = (w_3, w_2)^\top$ .

$$a) \quad \mathbf{x} = \begin{vmatrix} -1 & -1/2 & 0 & 1/2 & 1 \end{vmatrix}$$

$$\mathbf{y} = \begin{vmatrix} -1 & -1/8 & 0 & 1/8 & 1 \end{vmatrix}$$

$$\hat{\mu}_x = -1 + -1/2 + 1/2 + 1$$

$$\hat{\mu}_x = \begin{vmatrix} 0 \end{vmatrix}$$

$$\hat{\mu}_y = \begin{vmatrix} -1 & -1/8 & 0 & 1/8 & 1 \end{vmatrix} / 5$$

$$\hat{\mu}_y = 0$$

$$\bar{x} = x$$

$$\bar{y} = y$$

$$\hat{S}_x = \frac{1}{n} \bar{x} \cdot \bar{x}^\top$$

$$\hat{S}_x = \frac{1}{5} \begin{vmatrix} -1 & -1/2 & 0 & 1/2 & 1 \end{vmatrix} \begin{vmatrix} -1 \\ -1/2 \\ 0 \\ 1/2 \\ 1 \end{vmatrix}^\top$$

$$\hat{S}_x = 1/2 \quad \begin{vmatrix} -1 \\ -1/2 \\ 0 \\ 1/2 \\ 1 \end{vmatrix}^\top$$

$$\hat{S}_{xy} = \frac{1}{n} \bar{x} \bar{y}^\top$$

$$\hat{S}_{xy} = \frac{1}{5} \begin{vmatrix} -1 & -1/2 & 0 & 1/2 & 1 \end{vmatrix} \begin{vmatrix} -1 \\ -1/8 \\ 0 \\ 1/8 \\ 1 \end{vmatrix}^\top$$

$$\hat{S}_{xy} = 17/40$$

$$w_{OLS} = (\hat{S}_x)^{-1} \cdot \hat{S}_{xy}$$

$$w_{OLS} = 2 \cdot \frac{17}{40} = \frac{34}{40}$$

$$w_{OLS} = 34/40$$

$$b_{OLS} = \hat{\mu}_y - (w_{OLS}^\top) \hat{\mu}_x$$

$$b_{OLS} = 0$$

$$b) \quad w^\top \phi_j + b = w_3 (x_j)^3 + w_2 (x_j)^2 + b$$

$$\phi_1 = (-1)^3, (-1)^2 \} = (-1, 1)$$

$$\phi_1 = (-1, 1)$$

$$\phi_2 = ((-1/2)^3, (-1/2)^2)$$

$$\phi_2 = (-1/8, 1/4)$$

$$\phi_5 = (1, 1)$$

$$(w_{PLS}, b_{PLS}) = \arg \min_{w,b} \sum_{j=1}^5 (y_j - w^\top \phi_j - b)^2$$

Need to find  $w_3$

$$x = \begin{vmatrix} -1^3 & -1/2^3 & 0^3 & 1/2^3 & 1^3 \end{vmatrix}$$

$$x = \begin{vmatrix} -1 & -1/8 & 0 & 1/8 & 1 \end{vmatrix}$$

$$\hat{\mu}_x = 0$$

$$\hat{S}_x = \frac{1}{S} \begin{vmatrix} -1 & -1/8 & 0 & 1/8 & 1 \end{vmatrix} \begin{pmatrix} -1 \\ -1/8 \\ 0 \\ 1/8 \\ 1 \end{pmatrix}$$

$$\hat{S}_x = (1 + -1/8 + -1/8 + 1/8(1/8) + 1)^{1/2}$$

$$\hat{S}_x = (2 + 2/16)^{1/2} = 17/40$$

$$\hat{S}_{xy} = \frac{1}{S} \begin{vmatrix} -1 & -1/8 & 0 & 1/8 & 1 \end{vmatrix} \begin{pmatrix} -1 \\ -1/8 \\ 0 \\ 1/8 \\ 1 \end{pmatrix}$$

$$\hat{S}_{xy} = 17/40$$

$$w_3 = \frac{40}{17} \cdot \frac{17}{40} = 1$$

$$w_3 = 1$$

Need to find  $w_2$

$$x = \begin{vmatrix} -1^2 & -1/2^2 & 0^2 & 1/2^2 & 1^2 \end{vmatrix}$$

$$x = \begin{vmatrix} 1 & 1/4 & 0 & 1/4 & 1 \end{vmatrix}$$

$$\hat{\mu}_x = 1/5 (2^2/4) = 1/2$$

$$\bar{x} = \hat{x} - \hat{\mu}_x = \begin{vmatrix} 1/2 & -1/4 & -1/2 & -1/2 & 1/2 \end{vmatrix}$$

$$\hat{S}_x = \frac{1}{S} \begin{vmatrix} 1/2 & -1/4 & -1/2 & -1/2 & 1/2 \end{vmatrix} \begin{pmatrix} 1/2 \\ -1/4 \\ -1/2 \\ -1/2 \\ 1/2 \end{pmatrix}$$

$$\hat{S}_x = 17/40$$

$$\hat{S}_{xy} = \frac{1}{S} \begin{vmatrix} 1/2 & -1/4 & -1/2 & -1/2 & 1/2 \end{vmatrix} \begin{pmatrix} -1 \\ -1/8 \\ 0 \\ 1/8 \end{pmatrix}$$

$$\hat{S}_{xy} = 0$$

$$w_2 = 0$$

$$\Phi_1 : (-1 - |w_3 w_2| \begin{pmatrix} -1 \\ 1 \end{pmatrix} - b)^2$$

$$(-1 + w_3 - w_2 - b)^2$$

$$(\cancel{+} \cancel{-b})^2$$

$$(-b)^2$$

$$\Phi_2 : (-1/8 - |w_3 w_2| \begin{pmatrix} -1/8 \\ -1/4 \end{pmatrix} - b)^2$$

$$(-\cancel{1/8} + \cancel{1/8} w_3 + w_2/4 - b)^2$$

$$(-b)^2$$

$$\Phi_3 : (-1 w_3 w_2 | \begin{pmatrix} -1/8 \\ -1/4 \end{pmatrix} - b)^2$$

$$(w_3/8 + w_2/4 - b)^2$$

$$(1/8 - b)^2$$

$$\theta_4: \left( \|w_3 w_2\| \left\| \begin{pmatrix} \|w_8\| \\ -b \end{pmatrix} \right\|^2 \right)$$

$$\left( \cancel{\|w_8\|} - \cancel{w_3 w_2} - w_4 - b \right)^2$$

$$(-b)^2$$

$$\theta_5: \left( 1 - \|w_3 w_2\| \left\| \begin{pmatrix} -b \end{pmatrix} \right\|^2 \right)$$

$$\left( \cancel{1} - \cancel{w_3 w_2} - b \right)^2$$

$$(-b)^2$$

$$\arg \min_{M \in \mathbb{R}} \left[ (-b)^2 + (-b)^2 + (\|w_8\| - b)^2 + (-b)^2 + (-b)^2 \right]$$

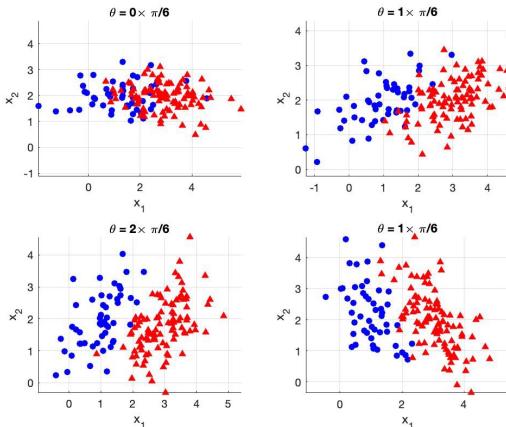
$$\frac{\|w_8 - Sb\|}{S} = 0 \rightarrow \|w_8\| = Sb$$

$$b = \|w_8\|$$

**Problem 4.3** [43pts] (*LDA*) In this problem we will implement LDA on a synthetic dataset and also develop geometric intuition for the eigenvalues and eigenvectors of the empirical feature covariance matrices. You are provided skeleton code to assist you in implementing LDA.

- (a) [13pts] (*Synthetic dataset generation*) Write a Matlab function to generate as its output a 2-class labeled dataset consisting of 2D feature vectors each drawn independently from one of two Gaussian distributions (one distribution for each class). The function's inputs are (1) a specified number of examples  $n_1, n_2$  for each distribution (class), (2) specified  $2 \times 1$  mean vectors  $\mu_1, \mu_2$  for each class, and (3) a single common  $2 \times 2$  covariance matrix for both classes specified via two real, nonnegative eigenvalues  $\lambda_1, \lambda_2 \geq 0$ , and a single orientation variable  $\theta$  that defines their corresponding orthonormal eigenvectors  $\mathbf{u}_1 = (\cos \theta, \sin \theta)^T$  and  $\mathbf{u}_2 = (\sin \theta, -\cos \theta)^T$ .

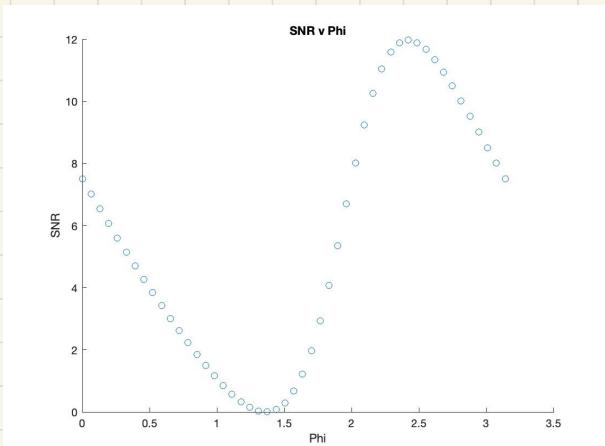
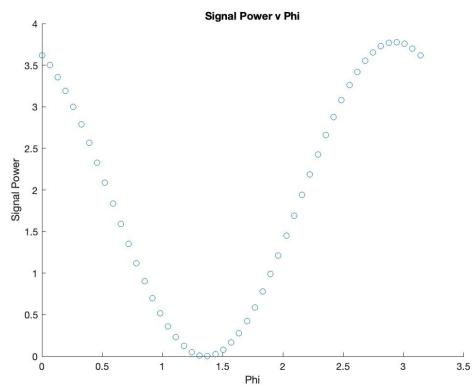
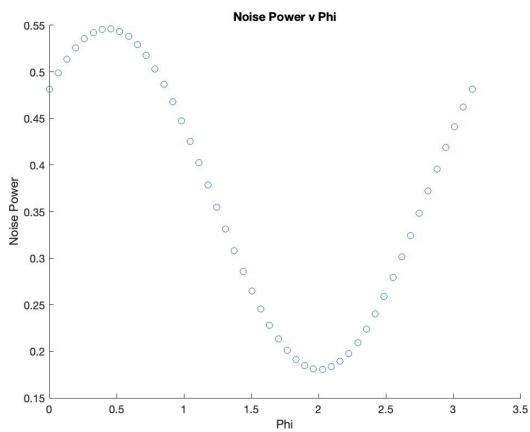
Use this function to generate  $n_1 = 50$  class 1 examples and  $n_2 = 100$  class 2 examples with  $\mu_1 = (1, 2)^T, \mu_2 = (3, 2)^T$  and the following four choices of the  $\lambda$ 's and  $\theta$ : (i)  $\lambda_1 = 1, \lambda_2 = 0.25, \theta = 0$ , (ii)  $\lambda_1 = 1, \lambda_2 = 0.25, \theta = \pi/6$ , (iii)  $\lambda_1 = 1, \lambda_2 = 0.25, \theta = \pi/3$ , (iv)  $\lambda_1 = 0.25, \lambda_2 = 1, \theta = \pi/6$ . For each of these four choices, create a scatter plot with class 1 examples shown as solid blue circles and class 2 examples shown as solid red triangles. Discuss how the eigenvalues and eigenvectors affect the geometry of the dataset.



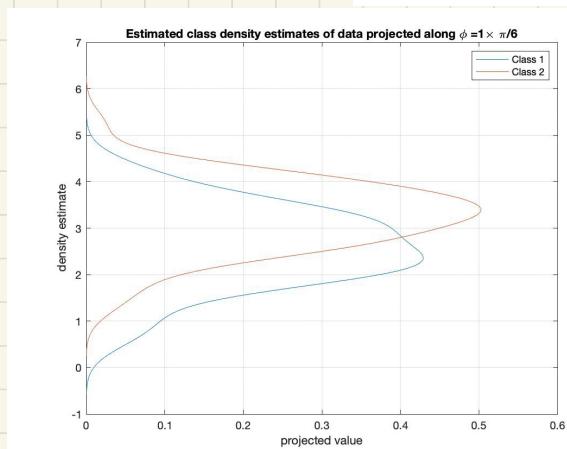
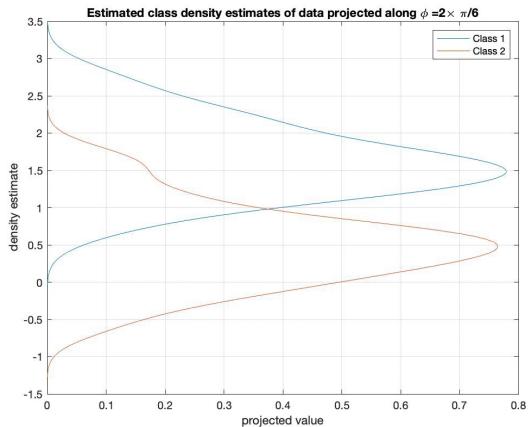
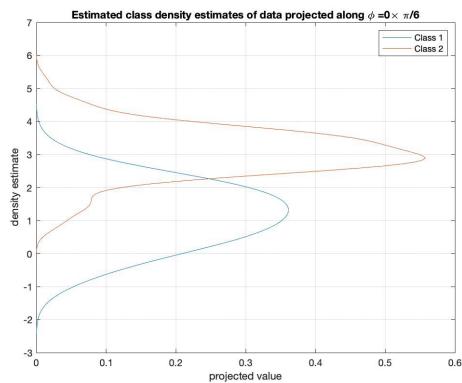
Eigenvalues change the orientation of the points! It adjusts the slope and placement of data.

- (b) [17pts] Write a Matlab function which takes as inputs a 2D labeled dataset for binary classification and a unit direction (specified by an angle  $\phi$  in radians) and returns the following three outputs (1) the squared Euclidean distance between the class means of the projections of the feature vectors along the direction  $(\cos \phi, \sin \phi)^\top$ , i.e.,  $(\hat{\mu}_{2z}(\phi) - \hat{\mu}_{1z}(\phi))^2$ , where  $\hat{\mu}_{1z}(\phi)$  and  $\hat{\mu}_{2z}(\phi)$  are, respectively, the class 1 and class 2 empirical means along direction  $\phi$ , (2) the average within-class variance of the projections of the feature vectors along direction  $(\cos \phi, \sin \phi)^\top$ , i.e.,  $\frac{n_1}{n} \hat{\sigma}_{1z}^2(\phi) + \frac{n_2}{n} \hat{\sigma}_{2z}^2(\phi)$ , where  $\hat{\sigma}_{1z}^2(\phi)$  and  $\hat{\sigma}_{2z}^2(\phi)$  are, respectively, the class 1 and class 2 empirical variances along direction  $\phi$ , and (3) their ratio, i.e., Signal-to-Noise Ratio (SNR).

Use this function to plot all three quantities as a function of  $\phi$ , for  $\phi$  ranging from 0 to  $\pi$  in steps of  $\pi/48$ , for the dataset of part(a)(ii). Determine (i) the value of  $\phi$  which maximizes the squared distance between the class means of the projected values (ii) the value of  $\phi$  which minimizes the average within-class variance, and (iii) the value of  $\phi$  which maximizes the SNR. (iv) Also use Matlab's `ksdensity` function with default settings to plot estimates of the class 1 and class 2 probability densities of the projections of the feature vectors for  $\phi = 0, \pi/6, \pi/3$ . Discuss your findings.



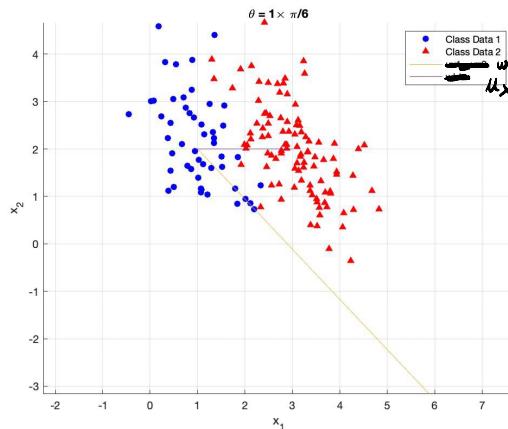
Noise, and signal are highly dependent on phi as they keep the shape of the cos/sine waves.



The density plots are adjusted profoundly according to  $\phi$ . Projecting points to a line is sensitive and can easily be influenced with the slightest change.

- (c) [6pts] Write a Matlab function which takes as input a 2D labeled dataset for binary classification and outputs the LDA vector  $\mathbf{w}_{LDA} = \hat{\mathbf{S}}_{x,\text{avg}}^{-1}(\hat{\mu}_{2x} - \hat{\mu}_{1x})$  where  $\hat{\mathbf{S}}_{x,\text{avg}} = \frac{n_1}{n}\hat{\mathbf{S}}_{1x} + \frac{n_2}{n}\hat{\mathbf{S}}_{2x}$  and  $\hat{\mathbf{S}}_{1x}$  and  $\hat{\mathbf{S}}_{2x}$  are the empirical  $2 \times 2$  covariance matrices of classes 1 and 2 respectively.

Use this function to compute  $\mathbf{w}_{LDA}$  for the dataset of part(a)(ii). Compare it with the difference between the class 2 and class 1 empirical mean vectors ( $\hat{\mu}_{2x} - \hat{\mu}_{1x}$ ). Overlay both these vectors on the scatter plot of the dataset with the vectors represented as arrows starting at the location of  $\hat{\mu}_{1x}$ . Discuss what you observe. Also compare the direction of  $\mathbf{w}_{LDA}$  with the value of  $\phi$  from part (b) which maximizes the SNR.

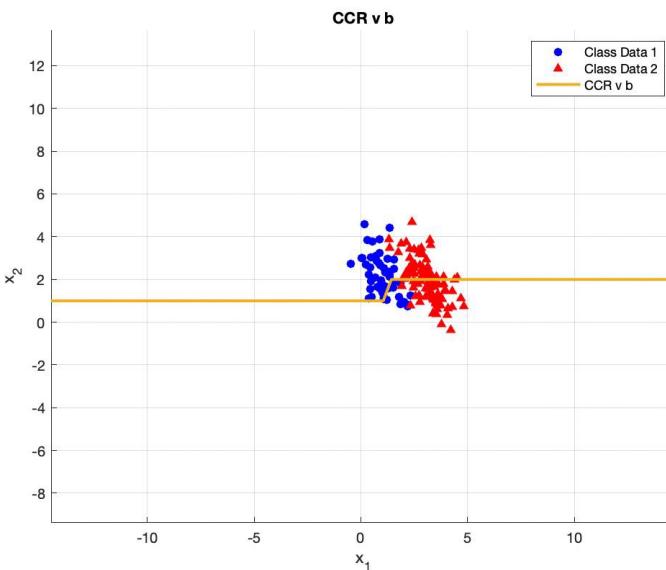


(Type  
in  
code)

The direction of  $w_{LDA}$  is very similar to the value of  $\phi$  which maximizes the SNR.  $w_{LDA}$  has a negative slope.

- (d) [7pts] Consider the following separating hyperplane decision rule for binary classification based on thresholding a linear (affine) function of the feature vector:  $h_{\mathbf{w},b}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}^\top \mathbf{x} + b \leq 0 \\ 2 & \text{else} \end{cases}$ . Write a Matlab function which takes as inputs a 2D labeled dataset for binary classification,  $\mathbf{w}$ , and  $b$  as inputs and outputs the value of CCR that results from applying  $h_{\mathbf{w},b}(\mathbf{x})$  to classify all feature vectors in the dataset into class 1 or class 2.

Use this function to create a plot of CCR as a function of  $b$  for the dataset of part(a)(ii) and  $\mathbf{w} = \mathbf{w}_{LDA}$ . Compute the value of the offset parameter  $b$  which maximizes the CCR and the resulting CCR value.



B that maximizes CCR: 1.428

**Problem 4.4** [22pts] (*ridge regression*) In this problem we will implement ridge regression and apply it to a real-world 8-dimensional (8 features) prostate cancer dataset contained in the file `prostateStnd.mat`. In this dataset, 8 medically relevant features named `lcavol`, `lweight`, `age`, `lbph`, `svi`, `lcp`, `gleason`, and `pgg45` are used to estimate `lpsa` (log prostate specific antigen). The training and test data are provided as `(xtrain, ytrain)` and `(xtest, ytest)` respectively. The first 8 features correspond to the first 8 entries of `names`. The ninth entry of `names` (the last one) is the label to be predicted whose values are in `(ytest, ytrain)`.

- (a) [4pts] As a first step, write Matlab code to **normalize the training dataset** so that post-normalization, each of the 8 features and the label in the normalized training dataset has zero mean and unit variance. This requires determining a *pair* of offset and scaling parameters, one pair for each feature and one pair for the label. These parameters must be computed only from the training dataset, but they must be applied to both the training and test datasets, i.e., we normalize both the training and test data, but we are only allowed to normalize the test data using parameters derived from the training data. In other words we must apply identical operations to training and test data. Only the training data will be actually normalized by the operation. If the test data is statistically similar to the training data, it too will be approximately normalized. **Report** the mean and variance of each feature and the label *before normalization*.

#### The mean and variance of training data before normalization

The mean of feature column `lcavol` is: -0.031  
The variance of feature column `lcavol` is: 1.123

The mean of feature column `lweight` is: -0.054  
The variance of feature column `lweight` is: 0.931

The mean of feature column `age` is: 0.119  
The variance of feature column `age` is: 1.026

The mean of feature column `lbph` is: -0.020  
The variance of feature column `lbph` is: 1.028

The mean of feature column `svi` is: 0.018  
The variance of feature column `svi` is: 1.040

The mean of feature column `lcp` is: -0.025  
The variance of feature column `lcp` is: 1.014

The mean of feature column `gleason` is: -0.030  
The variance of feature column `gleason` is: 0.974

The mean of feature column `pgg45` is: 0.067  
The variance of feature column `pgg45` is: 1.091

The mean of label column `lpsa` is: 2.452  
The variance of label column `lpsa` is: 1.459

- (b) [6pts] Next, write Matlab code to use the normalized data to train a ridge regression model for each of the following values of the quadratic regularization penalty parameter  $\lambda$ :  $\{e^{-5}, e^{-4}, e^{-3}, \dots, e^{10}\}$ .
- (c) [9pts] In a single figure, plot the ridge regression coefficient of each feature (8 in total) as a function of  $\ln \lambda$  (8 curves in total) for  $\ln \lambda$  ranging from -5 to 10 in steps of 1. Use suitable colors and/or markers to distinguish between the 8 curves and label them appropriately in a legend. Discuss what happens to the coefficients as  $\lambda$  becomes larger.
- (d) [3pts] In another figure, plot the mean-squared-error (MSE) of both the training and test data as a function of  $\ln \lambda$ . Discuss your observations.

I HAVE CODE FOR 4.4b & c  
DID GET TO FINISH & TEST.  
DID NOT GET A CHANCE  
TO COMPLETE 4.4d.