

Question 2

2(e)

Case 1: $P(C=1|x)=0.5$ Given: $P(C=1|x) = \frac{1}{1+e^{-z}}$, $z = w^T x + w_0$

Let $y = \frac{1}{1+e^{-z}} \Rightarrow y = \frac{1}{1+\frac{1}{e^z}}$
 $= \frac{1}{\frac{e^z+1}{e^z}} = \frac{e^z}{e^z+1} = y$

$\ln\left(\frac{y}{1-y}\right) = \ln\left(\frac{\frac{e^z}{e^z+1}}{1-\frac{e^z}{e^z+1}}\right)$ (log-odds ratio from slides)
 $= \ln\left(\frac{\frac{e^z}{e^z+1}}{\frac{(e^z+1)-e^z}{e^z+1}}\right) = \ln\left(\frac{\frac{e^z}{e^z+1}}{\frac{1}{e^z+1}}\right) = \ln(e^z)$

① $\ln\left(\frac{y}{1-y}\right) = z$

$\ln\left(\frac{0.5}{1-0.5}\right) = w_0 + w_1 x_1 + w_2 x_2$

$0 = w_0 + w_1 x_1 + w_2 x_2$

$w_2 x_2 = -w_0 - w_1 x_1$

$x_2 = \frac{-w_0 - w_1 x_1}{w_2}$, for $P(C=1|x)=0.5$

Case 2: $P(C=1|x)=0.6$

$\ln\left(\frac{0.6}{1-0.6}\right) = w_0 + w_1 x_1 + w_2 x_2$ (reuse eqn. ①)

$0.405465 = w_0 + w_1 x_1 + w_2 x_2$

$w_2 x_2 = -w_0 + 0.405465 - w_1 x_1$

$x_2 = \frac{-w_0 + 0.405465 - w_1 x_1}{w_2}$, for $P(C=1|x)=0.6$

Case 3: $P(C=1|x)=0.05$

$\ln\left(\frac{0.05}{1-0.05}\right) = w_0 + w_1 x_1 + w_2 x_2$ (reuse eqn. ①)

$-2.9444... = w_0 + w_1 x_1 + w_2 x_2$

$w_2 x_2 = -w_0 - 2.9444 - w_1 x_1$

$x_2 = \frac{-w_0 - 2.9444 - w_1 x_1}{w_2}$, for $P(C=1|x)=0.05$

Question 3

3(a)

Assume $p(y) = \text{prior}$
 $p(t|\phi) = \phi^t (1-\phi)^{1-t}$ (prior is Bernoulli)

Prove $\phi = \frac{1}{N} \sum_{n=1}^N \mathbb{I}[t^{(n)}=1]$

$$\ell(\phi, \mu_0, \mu_1, \Sigma) = -\ln \prod_{n=1}^N p(x^{(n)} | t^{(n)}, \mu_0, \mu_1, \Sigma) p(t^{(n)} | \phi) \quad (\text{slide 17})$$

(choose arbitrary class $k \in \{0, 1\}$, preserving generality, Assume $\Sigma_1 = \Sigma_2$
 (classes share same covar. matrix))

$$= -\ln \prod_{n=1}^N p(x^{(n)} | t^{(n)}, \mu_k, \Sigma) p(t^{(n)} | \phi) \quad (\text{Eqn from slide 6})$$

$$= -\ln \prod_{n=1}^N \left(\frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_k|^{\frac{1}{2}}} \cdot e^{-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k)} \right) \cdot \phi^t (1-\phi)^{1-t}$$

$$= -\left[\sum_{n=1}^N \ln(1) - \ln((2\pi)^{\frac{d}{2}} |\Sigma_k|^{\frac{1}{2}}) + \ln(e^{-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k)}) \right. \\ \left. + t \ln \phi + (1-t) \ln(1-\phi) \right]$$

$$= \sum_{n=1}^N \left\{ \frac{d}{2} \ln(2\pi) + \frac{1}{2} \ln |\Sigma_k| + (x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k) \right\} + t \ln \phi + (1-t) \ln(1-\phi) \quad \text{Eqn (1)}$$

Now equate derivative to 0

$$\frac{d}{d\phi} \ell(\phi, \mu_k, \Sigma_k) = 0 = \sum_{n=1}^N 0 + 0 + 0 + \frac{1}{\phi} t \ln \phi + \frac{1}{1-\phi} (1-t) \ln(1-\phi)$$

$$= \sum_{n=1}^N \frac{t^{(n)}}{\phi} + \frac{1}{\phi} (1-t^{(n)}) + 0 \cdot \ln(1-\phi) + (1-t^{(n)}) \cdot \frac{1}{(1-\phi)} (-1) \quad (\text{chain rule})$$

$$= \sum_{n=1}^N \left(\frac{t^{(n)}}{\phi} - \frac{(1-t^{(n)})}{(1-\phi)} \right) = \sum_{n=1}^N \frac{t^{(n)}(1-\phi) - \phi(1-t^{(n)})}{\phi(1-\phi)} = 0$$

$$0 = \sum_{n=1}^N (t^{(n)} - \phi) = \sum_{n=1}^N t^{(n)} - N\phi$$

$$\Rightarrow \phi = \frac{1}{N} \sum_{n=1}^N t^{(n)}$$

$$\Rightarrow \phi = \frac{1}{N} \sum_{n=1}^N \mathbb{I}[t^{(n)}=1] \quad (\text{since prior is Bernoulli, } t^{(n)} \in \{0, 1\})$$

$$\therefore \phi = \frac{1}{N} \sum_{n=1}^N \mathbb{I}[t^{(n)}=1]$$

3(b)

Given $\Sigma = \frac{1}{N} \sum_{n=1}^N (x^{(n)} - \mu_{t^{(n)}})^T (x^{(n)} - \mu_{t^{(n)}})$, where $x^{(n)}$ is row vector ①

Let X_0 be the data matrix for class 0, X_1 for class 1.

Prove $\Sigma = [y_0^T y_0 + y_1^T y_1] / N$, where $y_i = x_i - \vec{\mu}_i^T$
 where μ_i is column vector
 and $\vec{\mu}$ is column vector

$$y_0 = x_0 - \vec{\mu}_0^T$$

$$= \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix}_{nd} - \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1} \begin{bmatrix} \mu_{01} & \mu_{02} & \dots & \mu_{0d} \end{bmatrix}_{1 \times d}$$

$$= \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix}_{nd} - \begin{bmatrix} \mu_{01}^{(1)} & \mu_{02}^{(1)} & \dots & \mu_{0d}^{(1)} \\ \mu_{01}^{(2)} & \mu_{02}^{(2)} & \dots & \mu_{0d}^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ \mu_{01}^{(n)} & \mu_{02}^{(n)} & \dots & \mu_{0d}^{(n)} \end{bmatrix}_{nd}$$

$$y_0 = \begin{bmatrix} x^{(1)} - \mu_0 \\ x^{(2)} - \mu_0 \\ \vdots \\ x^{(n)} - \mu_0 \end{bmatrix}_{nd} \quad \text{②}$$

$$\Sigma = \frac{1}{N} \left[\sum_{n=1}^{N_0} (x^{(n)} - \mu_{t^{(n)}})^T (x^{(n)} - \mu_{t^{(n)}}) + \sum_{n=1}^{N_1} (x^{(n)} - \mu_{t^{(n)}})^T (x^{(n)} - \mu_{t^{(n)}}) \right] \quad \left(\begin{array}{l} \text{split given} \\ \text{eqn. ① by} \\ \text{class} \end{array} \right)$$

$$= \frac{1}{N} \left[\left[(x^{(1)} - \mu_0) \ (x^{(2)} - \mu_0) \dots (x^{(n)} - \mu_0) \right] \begin{bmatrix} x^{(1)} - \mu_0 \\ x^{(2)} - \mu_0 \\ \vdots \\ x^{(n)} - \mu_0 \end{bmatrix} + \sum_{n=1}^{N_1} \dots \right]$$

$$= \frac{1}{N} \left[y_0^T y_0 + \sum_{n=1}^{N_1} (x^{(n)} - \mu_{t^{(n)}})^T (x^{(n)} - \mu_{t^{(n)}}) \right] \quad \left(\begin{array}{l} \text{plug in eqn.} \\ \text{②} \end{array} \right)$$

$$\therefore \Sigma = \frac{1}{N} [y_0^T y_0 + y_1^T y_1] \quad \left(\begin{array}{l} \text{plug in eqn.} \\ \text{② but symmetrically} \\ \text{for } y_1 \text{ and } y_1^T \end{array} \right)$$

3(c)

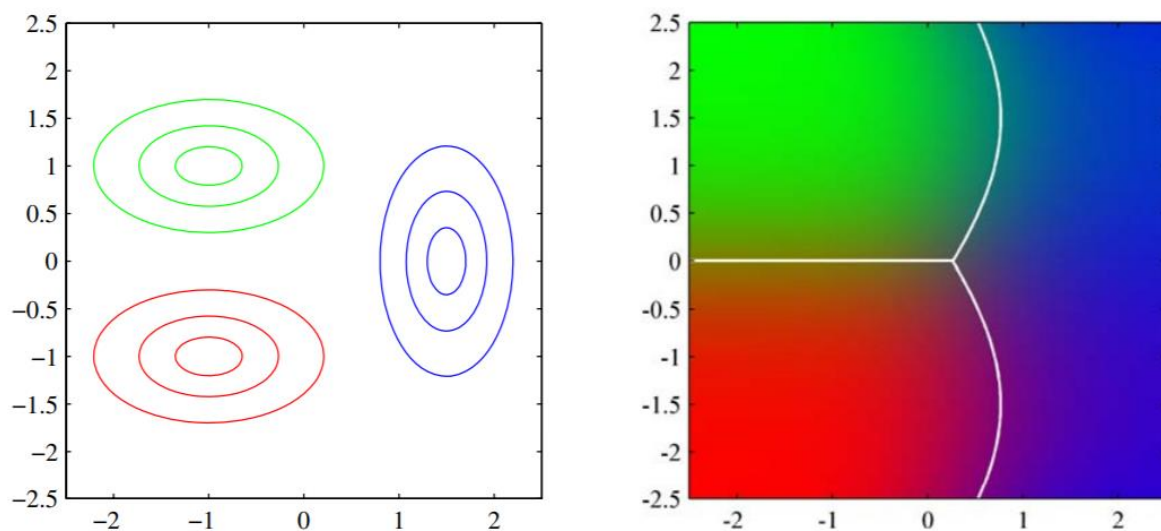
"I don't know"

Question 4

4(b)

Quadratic Discriminant Analysis, or QDA, does not assume that the covariance matrices are the same between different classes. In the case that the covariance matrices are the same, then the decision boundary between the classes which share the covariance matrix is linear. On the other hand, as we see in the Figure for 4(a), when the covariance matrices are different, QDA is free to learn quadratic boundaries. In particular, this is due to the fact that when you are calculating the decision boundaries with different covariance matrices, then the squared terms will not cancel out and you will receive quadratic decision boundaries.

The recommended text by *Bishop* demonstrates this in Figure 4.11:



Question 5

5(e)

To put it briefly, adding noise prevents the classifiers from learning parts of the whitespace surrounding the digits. This is related to the error message received about collinear variables because the classifications without noise experience a phenomenon called multicollinearity where one predictor variable can be linearly predicted from the others fairly accurately – leading to large changes in the results from small changes in the model or data. This is reflected as low accuracies in my printouts for models trained without noise.

5(f)

The similar accuracies for Naïve Bayes can be explained because Naïve Bayes relies on the conditional independence assumption that all the variables in the dataset are not correlated to each other. When this assumption holds, Naïve Bayes will converge quicker than discriminative models so you can get away with having less data to achieve similar accuracies. This is evident when comparing the Naïve Bayes accuracies between part e) and part f).

The performance for full Gaussian Bayes has plummeted because it does not converge as quickly and since we have only 90% of the data in part f) as we did in part e). The full Gaussian Bayes model has not yet converged and it is reflected by a drop in the accuracies.

5(g)

To put it simply and briefly, Gaussian Naïve Bayes predicts the probability of each class based on the feature vectors. A Gaussian is estimated for each feature, for each class. The class likelihoods are then used along with the priors and evidence to predict the posteriors for each class. The class with the largest posterior gets returned.

The background noise is no longer visible in the mean vectors because it got averaged out.