

Tut 6: LDA (Bayes' Classifier)

June 2024

The general mathematical formulation of a generative model:

$$\begin{aligned} h_D(\mathbf{x}) &= \operatorname{argmax}_{j \in \{1, \dots, K\}} \mathbb{P}(\mathbf{X} = \mathbf{x} | Y = j) \mathbb{P}(Y = j) \\ &= \operatorname{argmax}_{j \in \{1, \dots, K\}} [\ln \mathbb{P}(\mathbf{X} = \mathbf{x} | Y = j) + \ln \mathbb{P}(Y = j)] \end{aligned} \quad (6.1)$$

QDA (only works for numeric inputs which follows the normal distribution):

$$\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = j) \approx \frac{1}{(2\pi)^{p/2} \sqrt{|\mathbf{C}_j|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_j)^T \mathbf{C}_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j) \right\}.$$

LDA (only works for numeric inputs which follows the normal distribution):

$$\begin{aligned} \mathbb{P}(\mathbf{X} = \mathbf{x} | Y = j) &\approx \frac{1}{(2\pi)^{p/2} \sqrt{|\mathbf{C}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_j)^T \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu}_j) \right\}. \\ \Rightarrow h_D(\mathbf{x}) &= \operatorname{argmax}_{j \in \{1, \dots, K\}} \left\{ \ln \mathbb{P}(Y = j) + \tilde{\boldsymbol{\mu}}_j^T \mathbf{C}^{-1} \left[\mathbf{x} - \frac{1}{2} \tilde{\boldsymbol{\mu}}_j \right] \right\}. \end{aligned}$$

Naive Bayes:

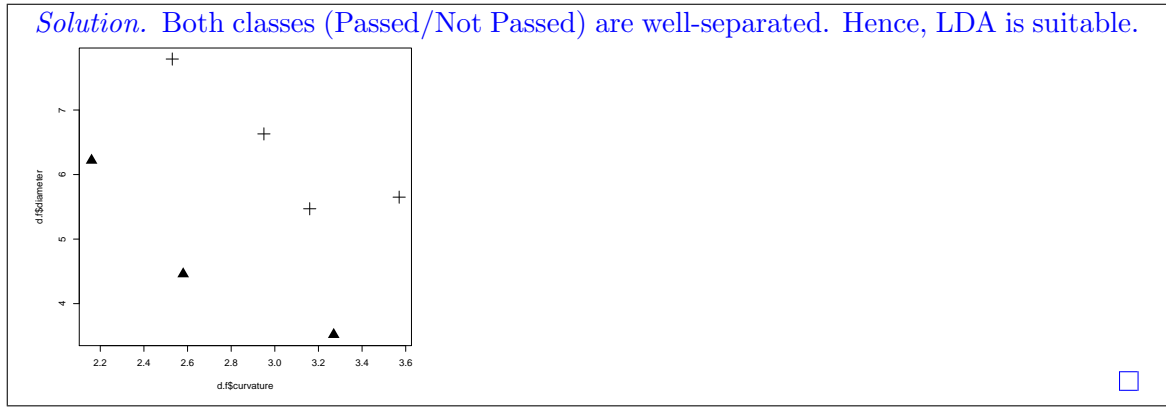
$$\mathbb{P}(\mathbf{X} = \mathbf{x} | Y = j) \approx \prod_{i=1}^p \mathbb{P}(X_i = x_i | Y = j)$$

1. Factory XYZ produces very expensive and high quality golf balls that their qualities are measured in term of curvature and diameter. Result of quality control by experts is given in the table below:

Curvature	Diameter	Result
2.95	6.63	Passed
2.53	7.79	Passed
3.57	5.65	Passed
3.16	5.47	Passed
2.58	4.46	Not Passed
2.16	6.22	Not Passed
3.27	3.52	Not Passed

As a consultant to the factory, you get a task to set up the criteria for automatic quality control using LDA model. Then, the manager of the factory also wants to test your criteria upon a new type of golf ball which have curvature 2.81 and diameter 5.46.

- (a) Plot the data with axes of curvature and diameter. Comment on the plot.



- (b) Write the data into matrix form by separating into “Passed” and “Not Passed”.

Solution. Let 1 = Passed and 2 = Not Passed.

$$X_1 = \begin{bmatrix} 2.95 & 6.63 \\ 2.53 & 7.79 \\ 3.57 & 5.65 \\ 3.16 & 5.47 \end{bmatrix} \quad X_2 = \begin{bmatrix} 2.58 & 4.46 \\ 2.16 & 6.22 \\ 3.27 & 3.52 \end{bmatrix}$$

□

- (c) Compute the prior probability for both classes.

Solution. $\hat{\pi}_1 = \frac{n_1}{n} = \frac{4}{7}, \quad \hat{\pi}_2 = \frac{n_2}{n} = \frac{3}{7}$

□

- (d) Compute the mean vectors for both classes.

Solution. The mean vectors are

$$\hat{\mu}_1 = \left[\frac{2.95 + 2.53 + 3.57 + 3.16}{4}, \frac{6.63 + 7.79 + 5.65 + 5.47}{4} \right] = [3.0525, 6.3850]$$

$$\hat{\mu}_2 = \left[\frac{2.58 + 2.16 + 3.27}{3}, \frac{4.46 + 6.22 + 3.52}{3} \right] = [2.67, 4.7333]$$

□

- (e) Compute the group covariance matrix.

Solution.

$$X_1 - \hat{\mu}_1 = \begin{bmatrix} 2.95 & 6.63 \\ 2.53 & 7.79 \\ 3.57 & 5.65 \\ 3.16 & 5.47 \end{bmatrix} - [3.0525, 6.3850] = \begin{bmatrix} -0.1025 & 0.2450 \\ -0.5225 & 1.4050 \\ 0.5175 & -0.7350 \\ 0.1075 & -0.9150 \end{bmatrix}$$

$$\Rightarrow C_1 = (X_1 - \hat{\mu})^T (X_1 - \hat{\mu}) = \begin{bmatrix} 0.562875 & -1.23795 \\ -1.237950 & 3.41150 \end{bmatrix}$$

$$X_2 - \hat{\mu}_2 = \begin{bmatrix} 2.58 & 4.46 \\ 2.16 & 6.22 \\ 3.27 & 3.52 \end{bmatrix} - [2.67, 4.7333] = \begin{bmatrix} -0.0900 & -0.2733 \\ -0.5100 & 1.4867 \\ 0.6000 & -1.2133 \end{bmatrix}$$

$$C_2 = (X_2 - \hat{\mu})^T (X_2 - \hat{\mu}) = \begin{bmatrix} 0.6282 & -1.461600 \\ -1.4616 & 3.757067 \end{bmatrix}$$

The group covariance matrix

$$C = \frac{1}{7-2}(C_1 + C_2) = \begin{bmatrix} 0.238215 & -0.539910 \\ -0.539910 & 1.433713 \end{bmatrix}$$

□

- (f) Write down the discriminant functions for both classes.

Solution. The discriminant functions are $\delta_j(X) = \ln(\pi_j) - \frac{1}{2}\mu_j C^{-1}\mu_j^T + \mu_j C^{-1}\mathbf{x}^T$, $j = 1, 2$:

$$\begin{aligned} \mu_1 C^{-1} &= [3.0525 \quad 6.3850] \begin{bmatrix} 0.238215 & -0.539910 \\ -0.539910 & 1.433713 \end{bmatrix}^{-1} = [156.38334 \quad 63.34455] \\ \Rightarrow \delta_1(X) &= \ln \frac{4}{7} - \frac{1}{2}\mu_1 C^{-1} \begin{bmatrix} 3.0525 \\ 6.3850 \end{bmatrix} + \mu_1 C^{-1} \mathbf{x}^T = [156.38334 \quad 63.34455] \mathbf{x}^T - 441.4672 \\ \mu_2 C^{-1} &= [2.6700 \quad 4.7333] \begin{bmatrix} 0.238215 & -0.539910 \\ -0.539910 & 1.433713 \end{bmatrix}^{-1} = [127.59686 \quad 51.35205] \\ \delta_2(X) &= \ln \frac{3}{7} - \frac{1}{2}\mu_2 C^{-1} \begin{bmatrix} 2.6700 \\ 4.7333 \end{bmatrix} + \mu_2 C^{-1} \mathbf{x}^T = [127.59686 \quad 51.35205] \mathbf{x}^T - 292.7214 \end{aligned}$$

□

- (g) Transform all the given data into discriminant functions.

Solution. For the first data, $\mathbf{x}_1 = [2.95, 6.63]$,

$$\delta_1(\mathbf{x}_1) = [156.38334 \quad 63.34455] [2.95 \quad 6.63]^T - 441.4672 = 439.8380$$

$$\delta_2(\mathbf{x}_1) = [127.59686 \quad 51.35205] [2.95 \quad 6.63]^T - 292.7214 = 424.1534, \text{ etc.}$$

X_1	X_2	$\delta_1(X)$	$\delta_2(X)$	Class
2.95	6.63	439.8380	424.1534	1
2.53	7.79	447.6367	430.1311	1
3.57	5.65	474.7180	452.9385	1
3.16	5.47	399.1988	391.3804	1
2.58	4.46	244.5185	265.5086	2
2.16	6.22	290.3239	302.2976	2
3.27	3.52	292.8791	305.2795	2

□

- (h) Locate the new golf ball in the plot as well as the functions to classify it.

Solution. $x_{new} = [2.81, 5.46]$

$$\delta_1(x_{new}) = [156.38334 \quad 63.34455] \begin{bmatrix} 2.81 \\ 5.46 \end{bmatrix} - 441.4672 = 343.8312$$

$$\delta_2(x_{new}) = [127.59686 \quad 51.35205] \begin{bmatrix} 2.81 \\ 5.46 \end{bmatrix} - 292.7214 = 346.208$$

Hence, the new golf ball should be classified into class 2, i.e. “not passed”.

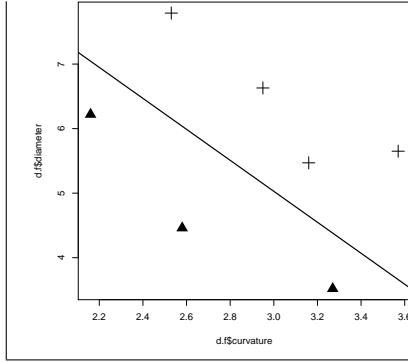
□

- (i) Plot the discriminant line into plot of $\delta_1(X)$ versus $\delta_2(X)$.

Solution. The discriminant functions fail when $\delta_1(X) = \delta_2(X)$, i.e.

$$(156.38334 - 127.59686)x_1 + (63.34455 - 51.35205)x_2 - (441.4672 - 292.7214) = 0 \\ \Rightarrow x_2 = 12.40324 - 2.4x_1$$

which can be used to plot a line using `abline(12.40324, -2.4)`.



2. (Final Assessment Jan 2021 Q5(a)) The data in Table 5.1 contains size measurements for two penguin species, i.e. Adelie and Gentoo, observed on three islands in the Palmer Archipelago, Antarctica.

Table 5.1: Palmer Archipelago Penguin data

Flipper length (mm)	Body mass (gram)	Species
196	4400	Adelie
188	3050	Adelie
219	5250	Gentoo
193	4200	Adelie
208	4200	Gentoo
215	5000	Gentoo
197	4775	Adelie

Construct a linear discriminant analysis (LDA) model for the given data in Table 5.1 by following the following steps.

- (i) Write down the general mathematical formula of the LDA model. (2 marks)

Solution. The general mathematical formulation is

[2 marks]

$$\begin{aligned}
 h_D(\mathbf{x}) &= \operatorname{argmax}_{j \in \{1, \dots, K\}} \frac{1}{(2\pi)^{p/2} \sqrt{|\mathbf{C}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \vec{\mu}_j)^T \mathbf{C}^{-1} (\mathbf{x} - \vec{\mu}_j) \right\} \mathbb{P}(Y = j) \\
 &= \operatorname{argmax}_{j \in \{1, \dots, K\}} \left\{ \ln \mathbb{P}(Y = j) + \vec{\mu}_j^T \mathbf{C}^{-1} \left[\mathbf{x} - \frac{1}{2} \vec{\mu}_j \right] \right\}.
 \end{aligned}$$

□

- (ii) Write down the estimates of all the parameters of discriminant functions of the LDA model in part (i). (11 marks)

Solution. The prior probability estimates for both classes are

$$\hat{\mathbb{P}}(Y = \text{Adelie}) = \frac{4}{7}, \quad \hat{\mathbb{P}}(Y = \text{Gentoo}) = \frac{3}{7}. \quad [2 \text{ marks}]$$

The mean vectors for both classes are

$$\begin{aligned}
 \vec{\mu}_{\text{Adelie}} &= \frac{1}{4}((196, 4400) + (188, 3050) + (193, 4200) + (197, 4775)) \\
 &= (193.50, 4106.25) \\
 \vec{\mu}_{\text{Gentoo}} &= \frac{1}{3}((219, 5250) + (208, 4200) + (215, 5000)) \\
 &= (214, 4816.667)
 \end{aligned} \quad [3 \text{ marks}]$$

We now estimate the “unscaled” covariance matrix estimate for Adelie:

$$X_{\text{Adelie}} - \vec{\mu}_{\text{Adelie}} = \begin{bmatrix} 196 & 4400 \\ 188 & 3050 \\ 193 & 4200 \\ 197 & 4775 \end{bmatrix} - [193.50, 4106.25] = \begin{bmatrix} 2.5 & 293.5 \\ -5.5 & -1056.5 \\ -0.5 & 93.5 \\ 3.5 & 668.5 \end{bmatrix}$$

$$C_{\text{Adelie}} = (X_{\text{Adelie}} - \vec{\mu}_{\text{Adelie}})^T (X_{\text{Adelie}} - \vec{\mu}_{\text{Adelie}}) = \begin{bmatrix} 49.0 & 8837.5 \\ 8837.5 & 1657969.0 \end{bmatrix} \quad [1.5 \text{ marks}]$$

We now estimate the “unscaled” covariance matrix estimate for Gentoo:

$$X_{\text{Gentoo}} - \vec{\mu}_{\text{Gentoo}} = \begin{bmatrix} 219 & 5250 \\ 208 & 4200 \\ 215 & 5000 \end{bmatrix} - [214, 4816.667] = \begin{bmatrix} 5 & 433.333 \\ -6 & -616.667 \\ 1 & 183.333 \end{bmatrix}$$

$$C_{\text{Gentoo}} = (X_{\text{Gentoo}} - \vec{\mu}_{\text{Gentoo}})^T (X_{\text{Gentoo}} - \vec{\mu}_{\text{Gentoo}}) = \begin{bmatrix} 62 & 6050.0 \\ 6050 & 601666.7 \end{bmatrix} \quad [1.5 \text{ marks}]$$

The group covariance matrix estimate is

$$C = \frac{1}{7-2} (C_{\text{Adelie}} + C_{\text{Gentoo}}) = \begin{bmatrix} 22.2 & 2977.5 \\ 2977.5 & 451927.1 \end{bmatrix} \quad [1 \text{ mark}]$$

Note that using scientific calculator or linear algebra, we can obtain

$$\mu_{\text{Adelie}}^T C^{-1} = \begin{bmatrix} 64.4419433 \\ -0.4154865 \end{bmatrix}, \quad \mu_{\text{Gentoo}}^T C^{-1} = \begin{bmatrix} 70.5666664 \\ -0.454267 \end{bmatrix}. \quad [1 \text{ mark}]$$

The discriminant functions (of the LDA model) for both classes are [1 mark]

$$\delta_{\text{Adelie}}(\mathbf{x}) = \ln \frac{4}{7} + \begin{bmatrix} 64.4419433 \\ -0.4154865 \end{bmatrix} \left(\mathbf{x} - \begin{bmatrix} 96.750 \\ 2053.125 \end{bmatrix} \right),$$

$$\delta_{\text{Gentoo}}(\mathbf{x}) = \ln \frac{3}{7} + \begin{bmatrix} 70.5666664 \\ -0.454267 \end{bmatrix} \left(\mathbf{x} - \begin{bmatrix} 107 \\ 2408.334 \end{bmatrix} \right)$$

□

- (iii) Use the discriminant functions of the LDA model in part (ii) to determine the species of a penguin with a flipper length of 218 mm and a body mass of 4590 gram. (4 marks)

Solution.

$$\delta_{\text{Adelie}}\left(\begin{bmatrix} 218 \\ 4590 \end{bmatrix}\right) = \ln \frac{4}{7} + \begin{bmatrix} 64.4419433 \\ -0.4154865 \end{bmatrix} \left(\begin{bmatrix} 218 \\ 4590 \end{bmatrix} - \begin{bmatrix} 96.750 \\ 2053.125 \end{bmatrix} \right),$$

$$= -0.5596 - (-6759.548) = 6758.989$$

$$\delta_{\text{Gentoo}}\left(\begin{bmatrix} 218 \\ 4590 \end{bmatrix}\right) = \ln \frac{3}{7} + \begin{bmatrix} 70.5666664 \\ -0.454267 \end{bmatrix} \left(\begin{bmatrix} 218 \\ 4590 \end{bmatrix} - \begin{bmatrix} 107 \\ 2408.334 \end{bmatrix} \right)$$

$$= -0.8473 - (-6841.841) = 6840.994 \quad [1.5+1.5=3 \text{ marks}]$$

Since $6840.994 > 6758.989$, the LDA model predicts the species to be Gentoo. [1 mark]

□

- (iv) (Not part of the final) Try to write down the QDA model.

Solution. Note that there are two inputs Flipper length and Body mass, therefore, $p = 2$ and $p/2 = 1$.

Part of the QDA model associated with Adelie is

$$\mathbb{P}(Y = \text{Adelie} | X_1 = x_1, X_2 = x_2) \propto P(Y = \text{Adelie}) \times \frac{1}{(2\pi)\sqrt{|C_{\text{Adelie}}|}} \exp \left\{ -\frac{1}{2} \begin{bmatrix} x_1 - 193.50 \\ x_2 - 4106.25 \end{bmatrix} C_{\text{Adelie}}^{-1} \begin{bmatrix} x_1 - 193.50 & x_2 - 4106.25 \end{bmatrix} \right\}.$$

Try to write down the QDA model associated with Gentoo and then the complete model.

□

3. (Final Exam Jan 2024 Sem, Q4(a)) Given the training data with three numeric features “bill length” (unit: mm), “bill depth” (unit: mm), “flipper length” (unit: mm) and the label “species” in Table 4.1.

Table 4.1: Training data of the penguin data with three different labels of penguins — Adelie, Chinstrap and Gentoo.

Obs.	bill length	bill depth	flipper length	species
A	41.1	19.1	188	Adelie
B	35.9	19.2	189	Adelie
C	36.0	17.9	190	Adelie
D	43.4	14.4	218	Gentoo
E	50.0	15.2	218	Gentoo
F	44.5	14.7	214	Gentoo
G	50.6	19.4	193	Chinstrap
H	45.7	17.0	195	Chinstrap

- (a) Write down the min-max scaling for all the features in Table 4.1 which transform Table 4.1 to Table 4.2.

Table 4.2: Scaled training data from Table 4.1

Obs.	bill length	bill depth	flipper length	species
A	0.3537	0.94	0.0000	Adelie
B	0.0000	0.96	0.0333	Adelie
C	0.0068	0.70	0.0667	Adelie
D	0.5102	0.00	1.0000	Gentoo
E	0.9592	0.16	1.0000	Gentoo
F	0.5850	0.06	0.8667	Gentoo
G	1.0000	1.00	0.1667	Chinstrap
H	0.6667	0.52	0.2333	Chinstrap

(3 marks)

Solution. $S_1(x) = \frac{x - 35.9}{14.7}$ [1 mark]

$S_2(x) = \frac{x - 14.4}{5}$ [1 mark]

$S_3(x) = \frac{x - 188}{30}$ [1 mark]

Average: 2.33 / 3 marks in Jan 2024; 14.81% below 1.5 marks.

□

(b) Given the results listing of LDA for Table 4.2:

```
lda(species ~ ., data = D.train.s)
```

Prior probabilities of groups:

Adelie	Chinstrap	Gentoo
0.375	0.250	0.375

Group means:

	bill_length	bill_depth	flipper_length
Adelie	0.1202	0.8667	0.0333
Chinstrap	0.8333	0.7600	0.2000
Gentoo	0.6848	0.0733	0.9556

By using Table 4.2, suppose the unscaled group covariance matrix for the species Adelie is

$$\begin{bmatrix} 0.0818 & 0.0248 & -0.0116 \\ 0.0248 & 0.0419 & -0.0080 \\ -0.0116 & -0.0080 & 0.0022 \end{bmatrix},$$

the unscaled group covariance matrix for the species Gentoo is

$$\begin{bmatrix} 0.1157 & 0.0379 & 0.0133 \\ 0.0379 & 0.0131 & 0.0018 \\ 0.0133 & 0.0018 & 0.0119 \end{bmatrix},$$

find the unscaled group covariance matrix for the species Chinstrap and the estimated common covariance matrix \hat{C} . (5 marks)

Solution. The unscaled group covariance for the species Chinstrap is

$$\begin{aligned} & \begin{bmatrix} 1.0000 - 0.83335 & 0.6667 - 0.83335 \\ 1.00 - 0.76 & 0.52 - 0.76 \\ 0.1667 - 0.2 & 0.2333 - 0.2 \end{bmatrix} \begin{bmatrix} 1.0000 - 0.83335 & 1.00 - 0.76 & 0.1667 - 0.2 \\ 0.6667 - 0.83335 & 0.52 - 0.76 & 0.2333 - 0.2 \end{bmatrix} \\ &= \begin{bmatrix} 0.16665 & -0.16665 \\ 0.24000 & -0.24000 \\ -0.03330 & 0.03330 \end{bmatrix} \begin{bmatrix} 0.16665 & 0.24 & -0.0333 \\ -0.16665 & -0.24 & 0.0333 \end{bmatrix} \\ &= \begin{bmatrix} 0.0555 & 0.0800 & -0.0111 \\ 0.0800 & 0.1152 & -0.0160 \\ -0.0111 & -0.0160 & 0.0022 \end{bmatrix} \end{aligned}$$

.....[3 marks]

The estimated common covariance matrix

$$\begin{aligned} \hat{C} &= \frac{1}{8-3} \begin{bmatrix} 0.0818 + 0.1157 + 0.0555 & 0.0248 + 0.0379 + 0.0800 & -0.0116 + 0.0133 - 0.0111 \\ 0.0248 + 0.0379 + 0.0800 & 0.0419 + 0.0131 + 0.1152 & -0.0080 + 0.0018 - 0.0160 \\ -0.0116 + 0.0133 - 0.0111 & -0.0080 + 0.0018 - 0.0160 & 0.0022 + 0.0119 + 0.0022 \end{bmatrix} \\ &= \begin{bmatrix} 0.05050 & 0.02854 & -0.00188 \\ 0.02854 & 0.03404 & -0.00444 \\ -0.00188 & -0.00444 & 0.00326 \end{bmatrix} \end{aligned}$$

.....[2 marks]

Average: 1.91 / 5 marks in Jan 2024; 51.85% below 2.5 marks. □

(c) Suppose the inverse matrix of the estimated common covariance matrix \hat{C} is

$$\begin{bmatrix} 39.33 & -36.50 & -27.03 \\ -36.50 & 69.60 & 73.74 \\ -27.03 & 73.74 & 391.59 \end{bmatrix}.$$

By finding the posterior probabilities or otherwise, predict the species of penguin with a bill length in 51.3 mm, a bill depth in 19.2 mm and a flipper length in 193 mm. (7 marks)

Solution. To perform prediction, one needs to scale the features using functions from part (i):

$$\mathbf{x}^* = \frac{51.3 - 35.9}{14.7} = 1.0476, \quad \frac{19.2 - 14.4}{5} = 0.96, \quad \frac{193 - 188}{30} = 0.1667 \quad [1 \text{ mark}]$$

Suppose we are estimating the posterior probabilities for j (being one of the penguin species), we will be using the formula

$$P(Y = j|\mathbf{x}) \propto P(Y = j) \exp(-\frac{1}{2}(\mathbf{x} - \mu_j)\hat{C}^{-1}(\mathbf{x} - \mu_j)^T). \quad [1 \text{ mark}]$$

Let $A = \text{Adelie}$, $C = \text{Chinstrap}$ and $G = \text{Gentoo}$.

$$\begin{aligned} P(Y = A|\mathbf{x}^*) &\propto 0.375 \exp(-\frac{1}{2} \begin{bmatrix} 1.0476 - 0.1202 \\ 0.96 - 0.8667 \\ 0.1667 - 0.0333 \end{bmatrix}^T \hat{C}^{-1} \begin{bmatrix} 0.9274 \\ 0.0933 \\ 0.1334 \end{bmatrix}^T) \\ &= 0.375 \exp(-\frac{30.2321}{2}) \end{aligned} \quad [2 \text{ marks}]$$

$$\begin{aligned} P(Y = C|\mathbf{x}^*) &\propto 0.250 \exp(\begin{bmatrix} 1.0476 - 0.8333 \\ 0.96 - 0.7600 \\ 0.1667 - 0.2000 \end{bmatrix}^T \hat{C}^{-1} \begin{bmatrix} 0.2143 \\ 0.2000 \\ -0.0333 \end{bmatrix}^T) \\ &= 0.250 \exp(-\frac{1.299226}{2}) \end{aligned} \quad [1 \text{ mark}]$$

$$\begin{aligned} P(Y = G|\mathbf{x}^*) &\propto 0.375 \exp(\begin{bmatrix} 1.0476 - 0.6848 \\ 0.96 - 0.0733 \\ 0.1667 - 0.9556 \end{bmatrix}^T \hat{C}^{-1} \begin{bmatrix} 0.3628 \\ 0.8867 \\ -0.7889 \end{bmatrix}^T) \\ &= 0.375 \exp(-\frac{192.4342}{2}) \end{aligned} \quad [1 \text{ mark}]$$

Since $P(Y = G|\mathbf{x}^*)$ and $P(Y = A|\mathbf{x}^*)$ are very small, the penguin species is predicted to be Chinstrap. [1 mark]

Average: 0.87 / 7 marks in Jan 2024; 88.89% below 3.5 marks. □