

I. Pen-and-paper

1)

Homework IV

1.

$$X_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad X_2 = \begin{bmatrix} 0 \\ 2 \end{bmatrix} \quad X_3 = \begin{bmatrix} 3 \\ -1 \end{bmatrix} \quad \pi_1 = 0,5 \quad \pi_2 = 0,5$$

$$\mu_1 = \begin{bmatrix} 2 \\ -1 \end{bmatrix} \quad \mu_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \Sigma_1 = \begin{bmatrix} 4 & 1 \\ 1 & 4 \end{bmatrix} \quad \Sigma_2 = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

$$P(X_1 | C_1) = \frac{1}{2\pi \sqrt{|\Sigma_1|}} e^{-\frac{1}{2}((1,0) - (2,-1)) \begin{bmatrix} 4 & 1 \\ 1 & 4 \end{bmatrix}^{-1} ((1,0) - (2,-1))}$$

$$= \frac{1}{2\pi \times \sqrt{15}} e^{-\frac{1}{30}((1,0) - (2,-1)) \begin{bmatrix} 4 & 1 \\ 1 & 4 \end{bmatrix}^{-1} ((1,0) - (2,-1))}$$

$$= \frac{1}{2\pi \sqrt{15}} e^{-1/3} \approx 0,029$$

$$P(X_1 | C_2) = \frac{1}{2\pi \sqrt{|\Sigma_2|}} e^{-\frac{1}{2}((1,0) - (1,1)) \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}^{-1} ((1,0) - (1,1))}$$

$$= \frac{1}{2\pi \times 2} e^{-2/8} \approx 0,062$$

$$\delta_{11} = \frac{0,029 \times 0,5}{0,5 \times 0,029 + 0,5 \times 0,062} \approx 0,315$$

$$\delta_{21} = \frac{0,062 \times 0,5}{0,5 \times 0,062 + 0,5 \times 0,029} \approx 0,681$$

$$P(X_2 | C_1) = \frac{1}{2\pi \sqrt{|\Sigma_1|}} e^{-\frac{1}{2}((0,2) - (2,1)) \begin{bmatrix} 4 & 1 \\ 1 & 4 \end{bmatrix}^{-1} ((0,2) - (2,1))}$$

$$= \frac{1}{2\pi \sqrt{15}} e^{-1/30 \times 64} \approx 0,005$$

$$P(x_2 | C_2) = \frac{1}{2\pi \times 2} e^{-1/2 \left((1025 - 0.13) \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}^{-1} \left(\begin{bmatrix} 0 \\ 2 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right)^T \right)}$$

$$= \frac{1}{4\pi} e^{-1/8 \times 6} = 0,048$$

$$\delta_{12} = \frac{0,005 \times 0,5}{0,5 \times (0,005 + 0,048)} = 0,094$$

$$\delta_{22} = \frac{0,048 \times 0,5}{0,5 \times (0,005 + 0,048)} = 0,906$$

$$P(x_3 | C_1) = \frac{1}{2\pi \sqrt{15}} e^{-1/2 \left((123 - 13) \begin{bmatrix} 1 & 1 \\ 0 & 2 \end{bmatrix}^{-1} \left(\begin{bmatrix} 3 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right)^T \right)}$$

$$= 0,035$$

$$P(x_3 | C_2) = \frac{1}{2\pi \times 2} e^{-1/2 \left((123 - 13) \begin{bmatrix} 1 & 1 \\ 0 & 2 \end{bmatrix}^{-1} \left(\begin{bmatrix} 3 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right)^T \right)}$$

$$= 0,011$$

$$\delta_{13} = \frac{0,035 \times 0,5}{0,5 \times (0,035 + 0,011)} = 0,766$$

$$\delta_{23} = \frac{0,011 \times 0,5}{0,5 \times (0,011 + 0,035)} = 0,234$$

$$N_1 = 0,315 + 0,094 + 0,766 = 1,175$$

$$N_2 = 0,681 + 0,906 + 0,234 = 1,821$$

$$\mu_1 = \frac{1}{1,175} \times \left(0,315 \times \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 0,094 \times \begin{bmatrix} 0 \\ 2 \end{bmatrix} + 0,766 \times \begin{bmatrix} 3 \\ 1 \end{bmatrix} \right)$$

$$= \begin{bmatrix} 2,220 \\ 0,450 \end{bmatrix}$$

$$\mu_2 = \frac{1}{1,821} \times \left(0,681 \times \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 0,906 \times \begin{bmatrix} 0 \\ 2 \end{bmatrix} + 0,234 \times \begin{bmatrix} 3 \\ 1 \end{bmatrix} \right) = \begin{bmatrix} 0,739 \\ 0,866 \end{bmatrix}$$

$$\begin{aligned}
 \Sigma_1 &= \frac{1}{1,175} \left(0,315 \times \left(\begin{bmatrix} -1 \\ 0 \end{bmatrix} - \begin{bmatrix} 2,220 \\ -0,430 \end{bmatrix} \right) \left(\begin{bmatrix} 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 2,220 & -0,430 \end{bmatrix} \right) \right. \\
 &\quad + 0,054 \times \left(\begin{bmatrix} 0 \\ 2 \end{bmatrix} - \begin{bmatrix} 2,220 \\ -0,430 \end{bmatrix} \right) \left(\begin{bmatrix} 0 & 2 \end{bmatrix} \cdot \begin{bmatrix} 2,220 & -0,430 \end{bmatrix} \right) \\
 &\quad \left. + 0,766 \times \left(\begin{bmatrix} 3 \\ -1 \end{bmatrix} - \begin{bmatrix} 2,220 \\ -0,430 \end{bmatrix} \right) \left(\begin{bmatrix} 3 & -1 \end{bmatrix} \cdot \begin{bmatrix} 2,220 & -0,430 \end{bmatrix} \right) \right) \\
 &= \frac{1}{1,175} \times \left(0,315 \times \begin{bmatrix} -1,220 \\ 0,430 \end{bmatrix} \begin{bmatrix} -1,220 & 0,430 \end{bmatrix} + \right. \\
 &\quad 0,054 \times \begin{bmatrix} -2,220 \\ 2,430 \end{bmatrix} \begin{bmatrix} -2,220 & 2,430 \end{bmatrix} + \\
 &\quad \left. 0,766 \times \begin{bmatrix} 0,780 \\ -0,510 \end{bmatrix} \begin{bmatrix} 0,780 & -0,510 \end{bmatrix} \right) \\
 &= \frac{1}{1,175} \left(0,315 \times \begin{bmatrix} 1,488 & -0,238 \\ -0,538 & 0,210 \end{bmatrix} + 0,054 \times \begin{bmatrix} 4,928 & -5,538 \\ -5,538 & 6,200 \end{bmatrix} \right. \\
 &\quad \left. + 0,766 \times \begin{bmatrix} 0,608 & -0,338 \\ -0,338 & 0,260 \end{bmatrix} \right) \\
 &= \begin{bmatrix} 1,130 & -0,861 \\ -0,861 & 0,728 \end{bmatrix} \\
 \Sigma_2 &= \frac{1}{1,821} \times \left(0,681 \times \begin{bmatrix} 0,241 \\ -0,866 \end{bmatrix} \begin{bmatrix} 0,241 & -0,866 \end{bmatrix} + \right. \\
 &\quad 0,506 \times \begin{bmatrix} -0,759 \\ 1,134 \end{bmatrix} \begin{bmatrix} -0,759 & 1,134 \end{bmatrix} + \\
 &\quad \left. 0,234 \times \begin{bmatrix} 2,241 \\ -1,866 \end{bmatrix} \begin{bmatrix} 2,241 & -1,866 \end{bmatrix} \right) \\
 &= \frac{1}{1,821} \times \left(0,681 \times \begin{bmatrix} 0,058 & -0,705 \\ -0,205 & 0,750 \end{bmatrix} + 0,506 \times \begin{bmatrix} 0,576 & -0,861 \\ -0,861 & 1,226 \end{bmatrix} \right. \\
 &\quad \left. + 0,234 \times \begin{bmatrix} 5,022 & -4,182 \\ -4,182 & 3,482 \end{bmatrix} \right) \\
 &= \begin{bmatrix} 0,954 & -1,044 \\ -1,044 & 1,368 \end{bmatrix}
 \end{aligned}$$

$$\bar{\pi}_1 = \frac{1,179}{1,179 + 1,821} = 0,393$$

$$\bar{\pi}_2 = \frac{1,821}{1,821 + 1,179} = 0,607$$

2nd Epoch

$$P(x_1 | C_1) = \frac{1}{\sqrt{2\pi} \times \sqrt{0,125}} e^{-\frac{1}{2} \left(\begin{bmatrix} 1 & 0 \end{bmatrix} - \begin{bmatrix} 2,22 & -0,49 \end{bmatrix} \begin{bmatrix} 1,180 & -0,861 \\ -0,861 & 0,728 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 0 \end{bmatrix} - \begin{bmatrix} 1,180 & -0,861 \\ -0,861 & 0,728 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1,180 & -0,861 \\ -0,861 & 0,728 \end{bmatrix} \right)}$$

$$= 0,116$$

$$P(x_1 | C_2) = \frac{1}{\sqrt{2\pi} \times \sqrt{0,215}} e^{-\frac{1}{2} \left(\begin{bmatrix} 1 & 0 \end{bmatrix} - \begin{bmatrix} 0,755 & 0,866 \end{bmatrix} \begin{bmatrix} 0,654 & -1,044 \\ -1,044 & 1,368 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 0 \end{bmatrix} - \begin{bmatrix} 0,755 & 0,866 \end{bmatrix} \begin{bmatrix} 0,654 & -1,044 \\ -1,044 & 1,368 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 0,654 & -1,044 \\ -1,044 & 1,368 \end{bmatrix} \right)}$$

$$= 0,149$$

$$\delta_{11} = \frac{0,116 \times 0,393}{0,116 \times 0,393 + 0,149 \times 0,607} = 0,335$$

$$\delta_{21} = \frac{0,149 \times 0,607}{0,149 \times 0,607 + 0,116 \times 0,393} = 0,665$$

$$P(x_2 | C_1) = \frac{1}{\sqrt{2\pi} \times \sqrt{0,115}} e^{-\frac{1}{2} \left(\begin{bmatrix} 0,1 & 7 \end{bmatrix} - \begin{bmatrix} 2,22 & -0,49 \end{bmatrix} \begin{bmatrix} 1,180 & -0,861 \\ -0,861 & 0,728 \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ 2 \end{bmatrix} - \begin{bmatrix} 1,180 & -0,861 \\ -0,861 & 0,728 \end{bmatrix}^{-1} \begin{bmatrix} 1,180 & -0,861 \\ -0,861 & 0,728 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \end{bmatrix} \begin{bmatrix} 1,180 & -0,861 \\ -0,861 & 0,728 \end{bmatrix} \right)}$$

$$= 0,001$$

$$P(x_2 | C_2) = \frac{1}{\sqrt{2\pi} \times \sqrt{0,215}} e^{-\frac{1}{2} \left(\begin{bmatrix} 0,1 & 7 \end{bmatrix} - \begin{bmatrix} 0,755 & 0,866 \end{bmatrix} \begin{bmatrix} 0,654 & -1,044 \\ -1,044 & 1,368 \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ 2 \end{bmatrix} - \begin{bmatrix} 0,755 & 0,866 \end{bmatrix} \begin{bmatrix} 0,654 & -1,044 \\ -1,044 & 1,368 \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ 2 \end{bmatrix} \begin{bmatrix} 0,654 & -1,044 \\ -1,044 & 1,368 \end{bmatrix} \right)}$$

$$= 0,207$$

$$\delta_{12} = \frac{0,001 \times 0,393}{0,001 \times 0,393 + 0,207 \times 0,607} = 0,003$$

$$\delta_{22} = \frac{0,207 \times 0,607}{0,207 \times 0,607 + 0,001 \times 0,393} = 0,997$$

$$P(x_3 | C_1) = \frac{1}{2\pi \sqrt{0.123}} e^{-1/2 \left(\begin{bmatrix} 3 & 1 \end{bmatrix} - \begin{bmatrix} 2.17 & 0.45 \end{bmatrix} \right)^T \begin{bmatrix} 1.100 & -0.861 \\ -0.861 & 0.728 \end{bmatrix}^{-1} \left(\begin{bmatrix} 3 \\ 1 \end{bmatrix} - \begin{bmatrix} 2.17 \\ 0.45 \end{bmatrix} \right)}$$

$$\approx 0.344$$

$$P(x_3 | C_2) = \frac{1}{2\pi \sqrt{0.215}} e^{-1/2 \left(\begin{bmatrix} 3 & 1 \end{bmatrix} - \begin{bmatrix} 0.755 & 0.866 \end{bmatrix} \right)^T \begin{bmatrix} 0.954 & -1.044 \\ -1.044 & 1.368 \end{bmatrix}^{-1} \left(\begin{bmatrix} 3 \\ 1 \end{bmatrix} - \begin{bmatrix} 0.755 \\ 0.866 \end{bmatrix} \right)}$$

$$\approx 0.012$$

$$\delta_{13} = \frac{0.344 \times 0.337}{0.344 \times 0.337 + 0.012 \times 0.607} \approx 0.549$$

$$\delta_{23} = \frac{0.012 \times 0.607}{0.012 \times 0.607 + 0.344 \times 0.337} \approx 0.051$$

$$N_1 = 0.335 + 0.003 + 0.549 = 1.287$$

$$N_2 = 0.665 + 0.957 + 0.051 = 1.713$$

$$\mu_1 = \frac{1}{1.287} \left(0.335 \times \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 0.003 \times \begin{bmatrix} 0 \\ 2 \end{bmatrix} + 0.549 \times \begin{bmatrix} 3 \\ 1 \end{bmatrix} \right)$$

$$= \begin{bmatrix} 2.472 \\ -0.733 \end{bmatrix}$$

$$\mu_2 = \frac{1}{1.713} \left(0.665 \times \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 0.957 \times \begin{bmatrix} 0 \\ 2 \end{bmatrix} + 0.051 \times \begin{bmatrix} 3 \\ 1 \end{bmatrix} \right)$$

$$= \begin{bmatrix} 0.418 \\ 1.134 \end{bmatrix}$$

$$\begin{aligned} \Sigma_1 &= \frac{1}{1,287} \left(0,335 \times \begin{bmatrix} -1,472 \\ 0,733 \end{bmatrix} \begin{bmatrix} -1,472 & 0,733 \end{bmatrix} + \right. \\ &\quad 0,003 \times \begin{bmatrix} -2,472 \\ 2,733 \end{bmatrix} \begin{bmatrix} -2,472 & 2,733 \end{bmatrix} + \\ &\quad \left. 0,949 \times \begin{bmatrix} 0,528 \\ -0,267 \end{bmatrix} \begin{bmatrix} 0,528 & -0,267 \end{bmatrix} \right) \\ &= \frac{1}{1,287} \left(0,335 \times \begin{bmatrix} 2,167 & -1,075 \\ -1,075 & 0,537 \end{bmatrix} + 0,003 \times \begin{bmatrix} 6,111 & -6,756 \\ -6,756 & 7,469 \end{bmatrix} \right. \\ &\quad \left. + 0,949 \times \begin{bmatrix} 0,279 & -0,141 \\ -0,141 & 0,141 \end{bmatrix} \right) \end{aligned}$$

$$= \begin{bmatrix} 0,784 & -0,400 \\ -0,400 & 0,210 \end{bmatrix}$$

$$\begin{aligned} \Sigma_2 &= \frac{1}{1,713} \left(0,665 \begin{bmatrix} 0,522 \\ -1,134 \end{bmatrix} \begin{bmatrix} 0,522 & -1,134 \end{bmatrix} + \right. \\ &\quad 0,357 \begin{bmatrix} -0,478 \\ 0,866 \end{bmatrix} \begin{bmatrix} -0,478 & 0,866 \end{bmatrix} + \\ &\quad \left. 0,051 \begin{bmatrix} 2,522 \\ -2,134 \end{bmatrix} \begin{bmatrix} 2,522 & -2,134 \end{bmatrix} \right) \end{aligned}$$

$$= \begin{bmatrix} 0,628 & -0,631 \\ -0,631 & 1,071 \end{bmatrix}$$

$$\overline{\Pi}_1 = \frac{1,287}{1,713 + 1,287} = 0,429$$

$$\overline{\Pi}_2 = \frac{1,713}{1,713 + 1,287} = 0,571$$

2)

a)

2. a)

X_1 :

$$P(X_1 | C_1) = \frac{1}{\sqrt{2\pi} \sqrt{0,005}} e^{-\frac{1}{2} \left(\left(\begin{bmatrix} 1 & 0 \end{bmatrix} - \begin{bmatrix} 2,472 & -0,733 \end{bmatrix} \right) \begin{bmatrix} 0,784 & -0,4 \\ -0,4 & 0,21 \end{bmatrix}^{-1} \left(\begin{bmatrix} 1 \\ 0 \end{bmatrix} - \begin{bmatrix} 2,472 \\ -0,733 \end{bmatrix} \right) \right)}$$

$$= 0,571$$

$$P(X_1 | C_2) = \frac{1}{\sqrt{2\pi} \sqrt{0,060}} e^{-\frac{1}{2} \left(\left(\begin{bmatrix} 1 & 0 \end{bmatrix} - \begin{bmatrix} 0,478 & 1,134 \end{bmatrix} \right) \begin{bmatrix} 0,428 & -0,631 \\ -0,631 & 1,071 \end{bmatrix}^{-1} \left(\begin{bmatrix} 1 \\ 0 \end{bmatrix} - \begin{bmatrix} 0,478 \\ 1,134 \end{bmatrix} \right) \right)}$$

$$= 0,112$$

$$\delta_{11} = \frac{0,428 \times 0,571}{0,428 \times 0,571 + 0,571 \times 0,112} = 0,793$$

$$\delta_{21} = \frac{0,571 \times 0,112}{0,571 \times 0,112 + 0,428 \times 0,571} = 0,207$$

X_2 :

$$P(X_2 | C_1) = \frac{1}{\sqrt{2\pi} \sqrt{0,005}} e^{-\frac{1}{2} \left(\left(\begin{bmatrix} 0 & 2 \end{bmatrix} - \begin{bmatrix} 2,472 & -0,733 \end{bmatrix} \right) \begin{bmatrix} 0,784 & -0,4 \\ -0,4 & 0,21 \end{bmatrix}^{-1} \left(\begin{bmatrix} 0 \\ 2 \end{bmatrix} - \begin{bmatrix} 2,472 \\ -0,733 \end{bmatrix} \right) \right)}$$

$$= 0,000$$

$$P(X_2 | C_2) = \frac{1}{\sqrt{2\pi} \sqrt{0,060}} e^{-\frac{1}{2} \left(\left(\begin{bmatrix} 0 & 2 \end{bmatrix} - \begin{bmatrix} 0,478 & 1,134 \end{bmatrix} \right) \begin{bmatrix} 0,428 & -0,631 \\ -0,631 & 1,071 \end{bmatrix}^{-1} \left(\begin{bmatrix} 0 \\ 2 \end{bmatrix} - \begin{bmatrix} 0,478 \\ 1,134 \end{bmatrix} \right) \right)}$$

$$= 0,004$$

$$\delta_{12} = \frac{0,000 \times 0,428}{0,000 \times 0,428 + 0,004 \times 0,571} = 0$$

$$\delta_{22} = \frac{0,004 \times 0,571}{0,000 \times 0,428 + 0,004 \times 0,571} = 1$$

$u_1 \in C_1$

$u_2 \in C_2$

$$P(x_3 | C_1) = \frac{1}{2\pi\sqrt{0,005}} e^{-\frac{1}{2} \left(\begin{bmatrix} 1 \\ 2 \end{bmatrix} - \begin{bmatrix} 1,172 \\ -0,783 \end{bmatrix} \right)^T \begin{bmatrix} 0,181 & -0,1 \\ -0,1 & 0,21 \end{bmatrix}^{-1} \left(\begin{bmatrix} 1 \\ 2 \end{bmatrix} - \begin{bmatrix} 1,172 \\ -0,783 \end{bmatrix} \right)}$$

$$= 1,955$$

$$P(x_3 | C_2) = \frac{1}{2\pi\sqrt{0,060}} e^{-\frac{1}{2} \left(\begin{bmatrix} 1 \\ 2 \end{bmatrix} - \begin{bmatrix} 0,448 \\ 1,134 \end{bmatrix} \right)^T \begin{bmatrix} 0,418 & -0,421 \\ -0,421 & 1,071 \end{bmatrix}^{-1} \left(\begin{bmatrix} 1 \\ 2 \end{bmatrix} - \begin{bmatrix} 0,448 \\ 1,134 \end{bmatrix} \right)}$$

$$= 0,000$$

$$\delta_{13} = \frac{0,423 \times 1,955}{0,423 \times 1,955 + 0,571 \times 0} = 1 \quad \left\{ \begin{array}{l} x_3 \in C_1 \end{array} \right.$$

$$\delta_{23} = \frac{0,571 \times 0}{0,571 \times 0 + 0,423 \times 1,955} = 0$$

b)

b)

$$\|x_1, x_3\|_2 = \sqrt{(1-3)^2 + (0-(-1))^2} = \sqrt{4+1} = \sqrt{5}$$

$$\|x_1, x_2\|_2 = \sqrt{(1-0)^2 + (0-2)^2} = \sqrt{1+4} = \sqrt{5}$$

$$S(x_1) = 1 - \frac{a(x_1)}{b(x_1)} = 1 - \frac{\|x_1, x_3\|_2}{\|x_1, x_2\|_2} = 1 - \frac{\sqrt{5}}{\sqrt{5}} = 0$$

$$\|x_2, x_3\|_2 = \sqrt{(0-3)^2 + (2-(-1))^2} = \sqrt{9+9} = \sqrt{18} = 3\sqrt{2}$$

$$S(x_3) = 1 - \frac{a(x_3)}{b(x_3)} = 1 - \frac{\sqrt{5}}{3\sqrt{2}} = 1 - 0,473 = 0,527$$

$$S(C_1) = \frac{S(x_1) + S(x_3)}{2} = \frac{0 + 0,527}{2} \approx 0,263$$

II. Programming and critical analysis

Here we have an overview of all the imports used and how we imported the dataset, which are common to all questions. Therefore, we will be omitting this code from the beginning of each task to avoid redundancy

```
import pandas as pd, matplotlib.pyplot as plt
import seaborn as sns

from sklearn.preprocessing import MinMaxScaler, StandardScaler
from sklearn.cluster import KMeans
from sklearn.decomposition import PCA

# Load the data
df = pd.read_csv('accounts.csv')
X = df.drop('deposit', axis=1)
X = pd.get_dummies(X, drop_first=True)
X = X.iloc[:, :8].drop_duplicates().dropna()
Y = pd.get_dummies(df['deposit'], drop_first=True)
```

1)

a) Normalizing the data and applying k-means for each k value

```
# Normalize the data using MinMaxScaler
scaler = MinMaxScaler()
X_normalized = scaler.fit_transform(X)

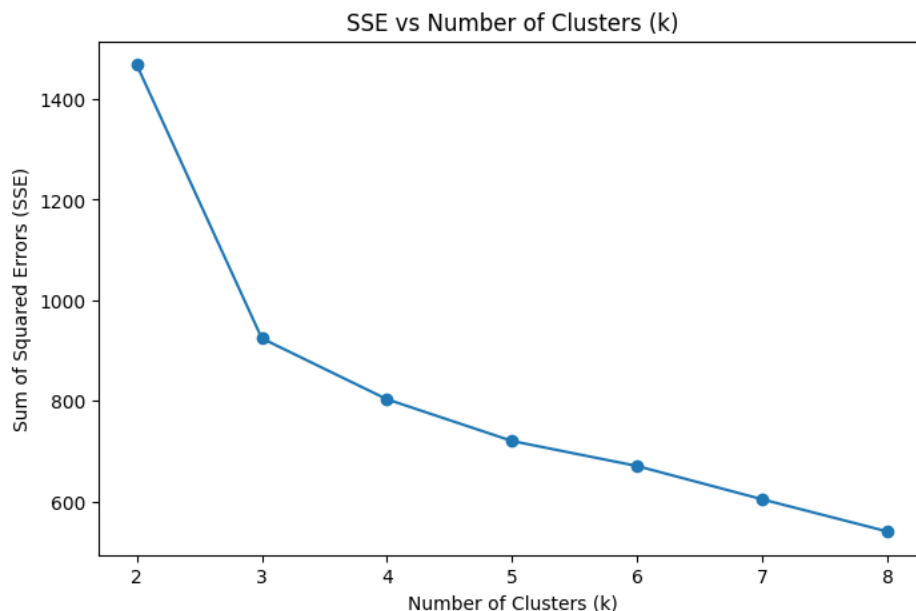
sse = []
k_values = range(2, 9)

# Apply k-means clustering for each k value
for k in k_values:
    kmeans = KMeans(n_clusters=k, max_iter=500, random_state=42)
    kmeans.fit_predict(X_normalized)
    sse.append(kmeans.inertia_)
```

Plotting the result

```
# Plot the SSE for each k
plt.figure(figsize=(8, 5))
plt.plot(k_values, sse, marker='o')
plt.title('SSE vs Number of Clusters (k)')
plt.xlabel('Number of Clusters (k)')
plt.ylabel('Sum of Squared Errors (SSE)')
plt.show()
```

Output:



- b) According to the plot, and using the knee/elbow finding method we come to the conclusion that the ideal number of clusters is 3. This method consists of finding the point where the decrease in SSE becomes less accentuated with the increase of clusters.

For all values before this point, the model might suffer from underfitting and for values above, the increase in the number of clusters and complexity does not provide substantial improvements to the model. Not only does it not justify, it also contributes to overfitting the data.

- c) K-modes is an adaptation of k-means used to better handle categorical features. It uses the mode (most frequent value) instead of means to represent the centroid of the cluster and uses distance based in dissimilarity (like Hamming Distance) to group the data.

Given that our dataset's features are predominantly categorical (10 categorical out of 17), in theory, k-modes would be a better clustering approach, considering the explanation above.

2)

- a) Normalizing the data with StandardScaler, applying the PCA and getting the explained variance ratio for each component

```
scaler = StandardScaler()
X_normalized = scaler.fit_transform(X)

# Apply PCA
pca = PCA(n_components=2)
pca.fit(X_normalized)

variance_explained = pca.explained_variance_ratio_

print("The first component explains", round(variance_explained[0], 5), "variability in the data set.")
print("The second component explains", round(variance_explained[1], 5), "variability in the data set.")
print("The top 2 components explain", round(variance_explained[1] + variance_explained[0], 5), "variability in the data set.")
```

The first component explains 0.19568 variability in the data set.

The second component explains 0.14462 variability in the data set.

The top 2 components explain 0.3403 variability in the dataset.

- b) Applying k-means clustering and PCA

```
kmeans = KMeans(n_clusters=3, random_state=42)
labels = kmeans.fit_predict(X_normalized)

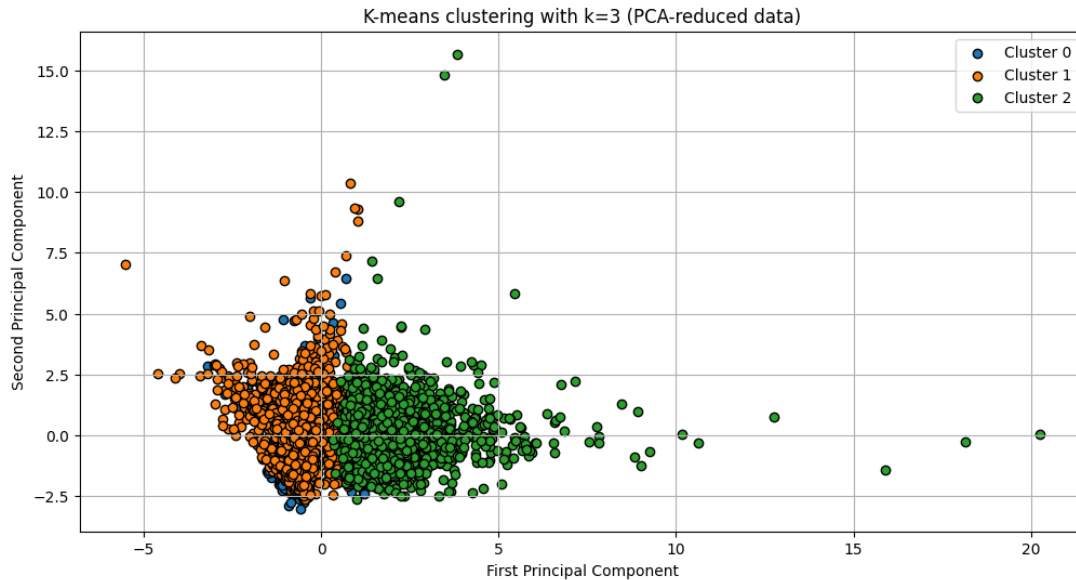
X_projected = pca.fit_transform(X_normalized)
```

Plotting the results:

```
plt.figure(figsize=(12, 6))
for cluster in range(3):
    plt.scatter(X_projected[labels == cluster, 0], X_projected[labels == cluster, 1], label=f'Cluster {cluster}', edgecolors='k')

plt.title('K-means clustering with k=3 (PCA-reduced data)')
plt.xlabel('First Principal Component')
plt.ylabel('Second Principal Component')
plt.grid()
plt.legend()
plt.show()
```

0 output:



We can only clearly separate two of the three clusters (the yellow one and the magenta one). The third one (dark blue) is very hard to distinguish from the other two as it underlies both of them.

This happens because we are losing a lot of information on the data variance by only considering the top two components to draw the plot. The fact that there is significant overlap between some clusters suggests that using only these two components in a 2D space isn't enough to represent a lot of the variance between the different observations of the dataset, that might become more apparent in higher dimensional spaces.

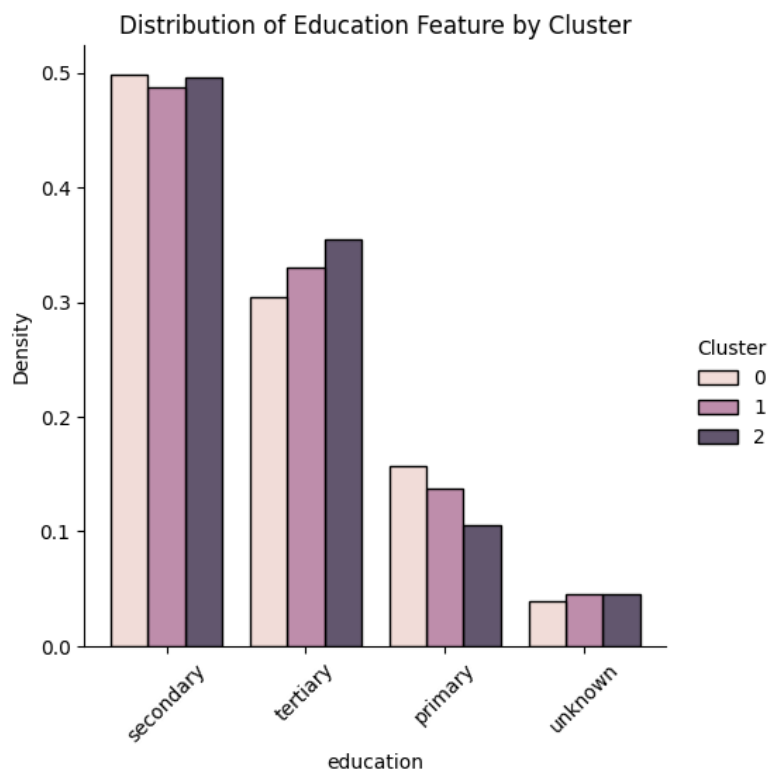
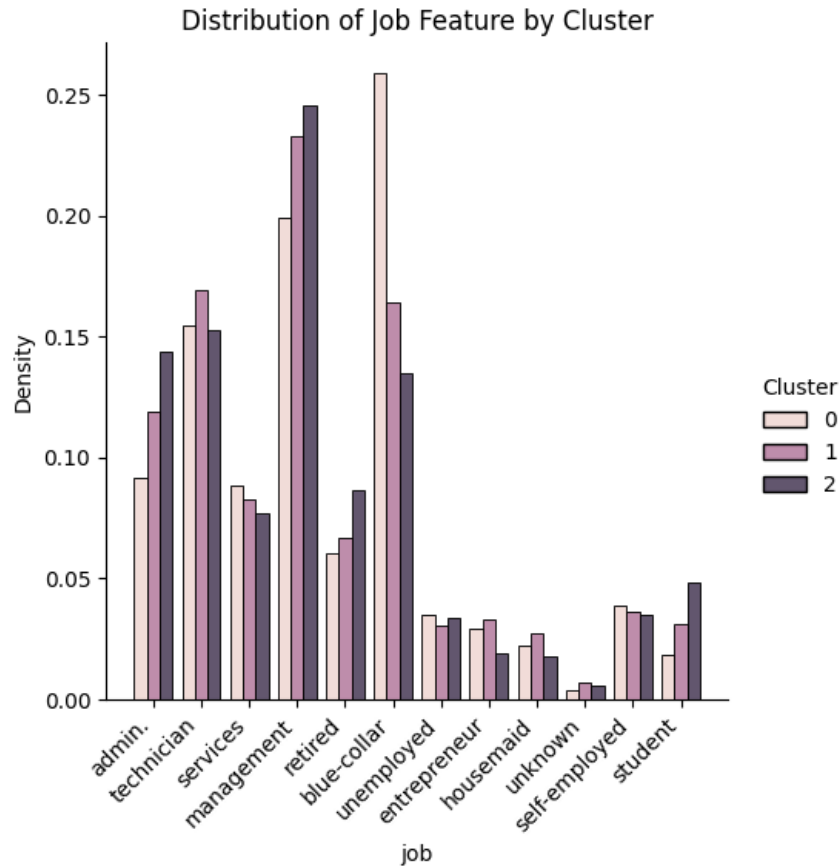
c) Plotting both graphs

```
df['Cluster'] = labels

# Plot the 'job' feature distribution across clusters
sns.displot(df, x="job", hue="Cluster", multiple="dodge", stat='density', shrink=0.8, common_norm=False)
plt.title('Distribution of Job Feature by Cluster')
plt.xticks(rotation=45, ha = 'right')
plt.show()

# Plot the 'education' feature distribution across clusters
sns.displot(df, x="education", hue="Cluster", multiple="dodge", stat='density', shrink=0.8, common_norm=False)
plt.title('Distribution of Education Feature by Cluster')
plt.xticks(rotation=45)
plt.show()
```


Output:



The main differences in the obtained plots can be described as follows:

- Jobs
 - Cluster 0: Jobs like “blue-collar”, “services”, “unemployed” and “self-employed” appear mostly in this cluster. They can be classified as less qualified positions and working-class jobs.
 - Cluster 1: Jobs like “technician”, “entrepreneur”, “housemaid” and “unknown” appear mostly in this cluster. We can’t find a significant relationship between these jobs, which makes us believe this cluster is better described by other features.
 - Cluster 2: Jobs like “admin”, “management”, “retired” and “student” appear mostly in this cluster. They describe either people who are not working or highly qualified jobs
- Education
 - Cluster 0: Mostly represented by “primary” and “secondary” levels of education, which translates well to the conclusions we made for the job distribution.
 - Cluster 1: Mostly represented by “unknown” level of education, which can also be a good indicator that this cluster might be better represented by other features.
 - Cluster 2: Mostly represented by “tertiary” level of education, which also translates well to the conclusions we made for the job distribution.

END