

1)

E - Step
 \tilde{x}_i

$$P(k | x_i) = \frac{P(x_i | k) P(k)}{P(x_i)}$$

 $k=1$

$$P(x_i | k=1) \sim N(\mu_1 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \Sigma_1 = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix})$$

$$= \frac{1}{2\pi \sqrt{|\Sigma|}} e^{-\frac{1}{2}(x_i - \mu_1)^T \Sigma_1^{-1} (x_i - \mu_1)}$$

$$|\Sigma| = 2 \times 2 - 1 \times 1 = 3$$

$$\Sigma_1^{-1} = \frac{1}{3} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$$

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

$$x_i - \mu_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix} - \begin{bmatrix} 2 \\ 2 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}$$

$$= \frac{1}{2\pi \sqrt{3}} e^{-\frac{1}{2} \begin{bmatrix} -1 & 0 \end{bmatrix} \begin{bmatrix} 2/3 & -1/3 \\ -1/3 & 2/3 \end{bmatrix} \begin{bmatrix} -1 \\ 0 \end{bmatrix}}$$

$$= \frac{1}{2\pi \sqrt{3}} e^{-1/3}$$

$$\approx 0.0658$$

$$\begin{aligned} P_{\text{posterior}} &= P(x_i | k=1) \times P_{\text{prior}} \\ &= 0.1317 \times 0.5 \\ &= 0.0658 \end{aligned}$$

$$k=2$$

$$P(x_1 | k=2) \sim N(\mu_2 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix})$$

$$= \frac{1}{2\pi \sqrt{|\Sigma_2|}} e^{-\frac{1}{2}(x_1 - \mu_2)^T \Sigma_2^{-1} (x_1 - \mu_2)}$$

$$|\Sigma_2| = 2 \times 2 = 4$$

$$\Sigma_2^{-1} = \frac{1}{4} \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix}$$

$$x_1 - \mu_2 = \begin{bmatrix} 1 \\ 2 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$= \frac{1}{2\pi \sqrt{4}} e^{-\frac{1}{2} \begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix}}$$

$$= \frac{1}{2\pi \sqrt{4}} e^{-9/4}$$

$$\approx 0.0228$$

$$\begin{aligned} \text{Posterior}_2 &= P(x_1 | k=2) \times \text{Prior}_2 \\ &= 0.0228 \times 0.5 \\ &= 0.0114 \end{aligned}$$

$$\begin{aligned} P(k=1 | x_1) &= \frac{0.033}{0.033 + 0.0114} \\ &\approx 0.7434 \end{aligned}$$

$$\begin{aligned} P(k=2 | x_1) &= \frac{0.0114}{0.033 + 0.0114} \\ &\approx 0.2568 \end{aligned}$$

x_2

$$P(K | x_2) = \frac{P(x_2 | K) P(K)}{P(x_2)}$$

$K=1$

$$P(x_2 | K=1) \sim N(\mu_1 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \Sigma_1 = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix})$$

$$= \frac{1}{2\pi \sqrt{|\Sigma|}} e^{-\frac{1}{2} (x_2 - \mu_1)^T \Sigma_1^{-1} (x_2 - \mu_1)}$$

$$|\Sigma| = 2 \times 2 - 1 \times 1 = 3$$

$$\Sigma_1^{-1} = \frac{1}{3} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$$

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

$$x_2 - \mu_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 2 \\ 2 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

$$= \frac{1}{2\pi \sqrt{3}} e^{-\frac{1}{2} \begin{bmatrix} -1 & -1 \end{bmatrix} \begin{bmatrix} 2/3 & -1/3 \\ -1/3 & 2/3 \end{bmatrix} \begin{bmatrix} -1 \\ -1 \end{bmatrix}}$$

$$= \frac{1}{2\pi \sqrt{3}} e^{-7/3}$$

$$\approx 0.00891$$

$$\begin{aligned} \text{Posterior}_1 &= P(x | K=1) \times \text{Prior}_1 \\ &= 9.7716 \times 0.5 \\ &= 0.00446 \end{aligned}$$

$$K=2$$

$$P(x_2 | K=2) \sim N(\mu_2 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix})$$

$$= \frac{1}{2\pi \sqrt{|\Sigma_2|}} e^{-\frac{1}{2}(x_2 - \mu_2)^T \Sigma_2^{-1} (x_2 - \mu_2)}$$

$$|\Sigma_2| = 2 \times 2 = 4$$

$$\Sigma_2^{-1} = \frac{1}{4} \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix}$$

$$x_2 - \mu_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$= \frac{1}{2\pi \sqrt{4}} e^{-\frac{1}{2} \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix}}$$

$$= \frac{1}{2\pi \sqrt{4}} e^{-1/2}$$

$$\approx 0.0483$$

$$\begin{aligned} \text{Posterior}_2 &= P(x | K=2) \times \text{Prior}_2 \\ &= 0.0483 \times 0.5 \\ &= 0.0241 \end{aligned}$$

$$\begin{aligned} P(K=1 | x_2) &= \frac{0.00446}{0.0241 + 0.00446} \\ &= 0.1562 \end{aligned}$$

$$\begin{aligned} P(K=2 | x_2) &= \frac{0.0241}{0.0241 + 0.00446} \\ &= 0.8438 \end{aligned}$$

x_3

$$P(K | x_3) = \frac{P(x_3 | K) P(K)}{P(x_3)}$$

$K=1$

$$P(x_3 | K=1) \sim N(\mu_1 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \Sigma_1 = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix})$$

$$= \frac{1}{2\pi \sqrt{|\Sigma_1|}} e^{-\frac{1}{2} (x_3 - \mu_1)^T \Sigma_1^{-1} (x_3 - \mu_1)}$$

$$|\Sigma_1| = 2 \times 2 - 1 \times 1 = 3$$

$$\Sigma_1^{-1} = \frac{1}{3} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$$

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

$$x_3 - \mu_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} - \begin{bmatrix} 2 \\ 2 \end{bmatrix} = \begin{bmatrix} -1 \\ -2 \end{bmatrix}$$

$$= \frac{1}{2\pi \sqrt{3}} e^{-\frac{1}{2} \begin{bmatrix} -1 & -2 \end{bmatrix} \begin{bmatrix} 2/3 & -1/3 \\ -1/3 & 2/3 \end{bmatrix} \begin{bmatrix} -1 \\ -2 \end{bmatrix}}$$

$$= \frac{1}{2\pi \sqrt{3}} e^{-1}$$

$$\approx 0.0338$$

$$\begin{aligned} P_{\text{posterior}_1} &= P(x | K=1) \times P_{\text{prior}_1} \\ &= 0.0338 \times 0.5 \\ &= 0.0169 \end{aligned}$$

$$K=2$$

$$P(x_3 | K=2) \sim N(\mu_2 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix})$$

$$= \frac{1}{2\pi \sqrt{|\Sigma_2|}} e^{-\frac{1}{2}(x_3 - \mu_2)^T \Sigma_2^{-1} (x_3 - \mu_2)}$$

$$|\Sigma_2| = 2 \times 2 = 4$$

$$\Sigma_2^{-1} = \frac{1}{4} \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix}$$

$$x_3 - \mu_2 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$= \frac{1}{2\pi \sqrt{4}} e^{-\frac{1}{2} \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}}$$

$$= \frac{1}{2\pi \sqrt{4}} e^{-1/4}$$

$$\approx 0.062$$

$$P_{\text{prior}} \sigma_2 = P(K=2) \times P_{\text{prior}} \sigma_2$$

$$= 0.062 \times 0.5$$

$$= 0.031$$

$$P(K=1 | x_3) = \frac{0.0169}{0.031 + 0.0169}$$

$$\approx 0.3528$$

$$P(K=2 | x_3) = \frac{0.031}{0.031 + 0.0169}$$

$$\approx 0.6472$$

M - Step

$$\mu_c = \frac{\sum_{i=1}^N P(c|x_i) \cdot x_i}{\sum_{i=1}^N P(c|x_i)}$$

$$\begin{aligned} P_{prior_c} &= P(c=k) \\ &= \frac{\sum_{j=1}^K P(c=j|x_i)}{\sum_{j=1}^K P(c=j|x_i)} = N \end{aligned}$$

$$\sum_c \mu_c = \frac{\sum_{k=1}^N P(c|k) \cdot (x_k^i - \mu_c^i)(x_k^j - \mu_c^j)}{\sum_{k=1}^N P(c|k)}$$

c=1

$$\begin{aligned} \sum_{i=1}^N P(c|x_i) &= 0.7434 + 0.1562 + 0.3538 \\ &= 1.2534 \end{aligned}$$

$$\mu_1 = \frac{0.7434 \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 0.1562 \begin{bmatrix} -1 \\ 1 \end{bmatrix} + 0.3538 \begin{bmatrix} 1 \\ 0 \end{bmatrix}}{1.2534}$$

$$= \begin{bmatrix} 0.7508 \\ 1.3108 \end{bmatrix}$$

$$\sum_{i=1}^{(1,1)} = \frac{0.7434(1 - 0.7508)^2 + 0.1562(-1 - 0.7508)^2 + 0.3538(1 - 0.7538)^2}{1.2534}$$

$$= 0.4359$$

$$\sum_{i=1}^{(2,1)} = \sum_{i=1}^{(1,2)}$$

$$= \frac{0.7434(1 - 0.7508)(2 - 1.3108) + 0.1562(-1 - 0.7508)(1 - 1.3108) + 0.3538(1 - 0.7508)(0 - 1.3108)}{1.2534}$$

$$= 0.0775$$

$$\sum_i^{(2,2)} = \frac{0.7434(2 - 1.3108)^2 + 0.1562(1 - 1.3108)^2 + 0.3538(0 - 1.3108)^2}{1.2534}$$

$$= 0.7788$$

Therefore, $\Sigma_i = \begin{bmatrix} 0.4359 & 0.0775 \\ 0.0775 & 0.7788 \end{bmatrix}$

$$\pi_i = \frac{1.2534}{3}$$

$$= 0.4178$$

$C = 2$

$$\sum_{i=1}^2 P(C|x_i) = 0.2568 + 0.8438 + 0.6472$$

$$= 1.7478$$

$$N_2 = \frac{0.2568 \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 0.8438 \begin{bmatrix} -1 \\ 1 \end{bmatrix} + 0.6472 \begin{bmatrix} 1 \\ 0 \end{bmatrix}}{1.7478}$$

$$= \begin{bmatrix} 0.0344 \\ 0.7766 \end{bmatrix}$$

$$\sum_i^{(1,1)} = \frac{0.2568(1 - 0.0344)^2 + 0.8438(-1 - 0.0344)^2 + 0.6472(1 - 0.0344)^2}{1.7478}$$

$$= 0.9988$$

$$\sum^{(2,1)} = \sum^{(1,2)}$$

$$= \frac{0.2568(1 - 0.0344)(2 - 0.7766) + 0.8438(-1 - 0.0344)(1 - 0.7766) + 0.6472(1 - 0.0344)(0 - 0.7766)}{1.7478}$$

$$= -0.2157$$

$$\sum_{i=1}^3 (x_i^{(2,2)})^2 = \frac{0.2568(2 - 0.7766)^2 + 0.8438(1 - 0.7766)^2 + 0.6472(0 - 0.7766)^2}{1.7478}$$

$$= 0.6517$$

Therefore, $\Sigma_2 = \begin{bmatrix} 0.9988 & -0.2157 \\ -0.2157 & 0.4673 \end{bmatrix}$

$$\tilde{\pi}_2 = \frac{1.7478}{3}$$

$$= 0.5826$$

To summarize, the new parameters are:

$$\mu_1 = \begin{bmatrix} 0.7508 \\ 1.3108 \end{bmatrix}$$

$$\mu_2 = \begin{bmatrix} 0.6344 \\ 0.7766 \end{bmatrix}$$

$$\Sigma_1 = \begin{bmatrix} 0.4359 & 0.0775 \\ 0.0775 & 0.7788 \end{bmatrix}$$

$$\Sigma_2 = \begin{bmatrix} 0.9988 & -0.2157 \\ -0.2157 & 0.4673 \end{bmatrix}$$

$$\tilde{\pi}_1 = 0.4178$$

$$\tilde{\pi}_2 = 0.5826 //$$

2.c) Since we are considering a hard assignment under MAP assumption, we must assign each observation to the highest posterior probability:

$$P(k=1|x_1) = \frac{0.033}{0.033 + 0.0114}$$

$$= 0.7434$$

$$P(k=1|x_2) = \frac{0.00446}{0.0241 + 0.00446}$$

$$= 0.1562$$

$$P(k=1|x_3) = \frac{0.0169}{0.031 + 0.0169}$$

$$= 0.3528$$

$$P(k=2|x_1) = \frac{0.0114}{0.033 + 0.0114}$$

$$= 0.2568$$

$$P(k=2|x_2) = \frac{0.0241}{0.0241 + 0.00446}$$

$$= 0.8438$$

$$P(k=2|x_3) = \frac{0.031}{0.031 + 0.0169}$$

$$= 0.6472$$

Therefore, x_1 is assigned to cluster 1 and x_2 and x_3 are assigned to cluster 2. //

2. b) Silhouette of c_2 (largest cluster) using Euclidean distance:

$$S(c_2) = \frac{S(x_2) + S(x_3)}{2}$$

$$\begin{aligned} a(x_2) &= d(x_2, x_3) \\ &= \sqrt{(-1-1)^2 + (1-0)^2} \\ &= \sqrt{5} \end{aligned}$$

$$\begin{aligned} b(x_2) &= d(x_2, x_1) \\ &= \sqrt{(-1-1)^2 + (1-2)^2} \\ &= \sqrt{5} \end{aligned}$$

$$a(x_2) \leq b(x_2)$$

\Downarrow

$$S(x_2) = 1 - \frac{a(x_2)}{b(x_2)}$$

$$= 1 - 1$$

$$= 0$$

$$\begin{aligned} a(x_3) &= d(x_3, x_2) \\ &= a(x_2) \\ &= \sqrt{5} \end{aligned}$$

$$\begin{aligned} b(x_3) &= d(x_3, x_1) \\ &= \sqrt{(1-1)^2 + (0-2)^2} \\ &= 2 \end{aligned}$$

$$a(x_3) > b(x_3)$$

\Downarrow

$$S(x_3) = \frac{b(x_3)}{a(x_3)} - 1$$

$$= \frac{2}{\sqrt{5}} - 1$$

$$\approx -0.1056$$

$$= \frac{0 - 0.1056}{2}$$

$$= -0.0528$$

Therefore, the silhouette of the largest cluster is $S(c_2) = -0.0528_{//}$

II. Programming and critical analysis

The code for the questions is in the Appendix of this document.

1) Seed: 0

Silhouette: 0.11362027575179438

Purity: 0.7671957671957672

Seed: 1

Silhouette: 0.1140355420137708

Purity: 0.7632275132275133

Seed: 2

Silhouette: 0.11362027575179438

Purity: 0.7671957671957672

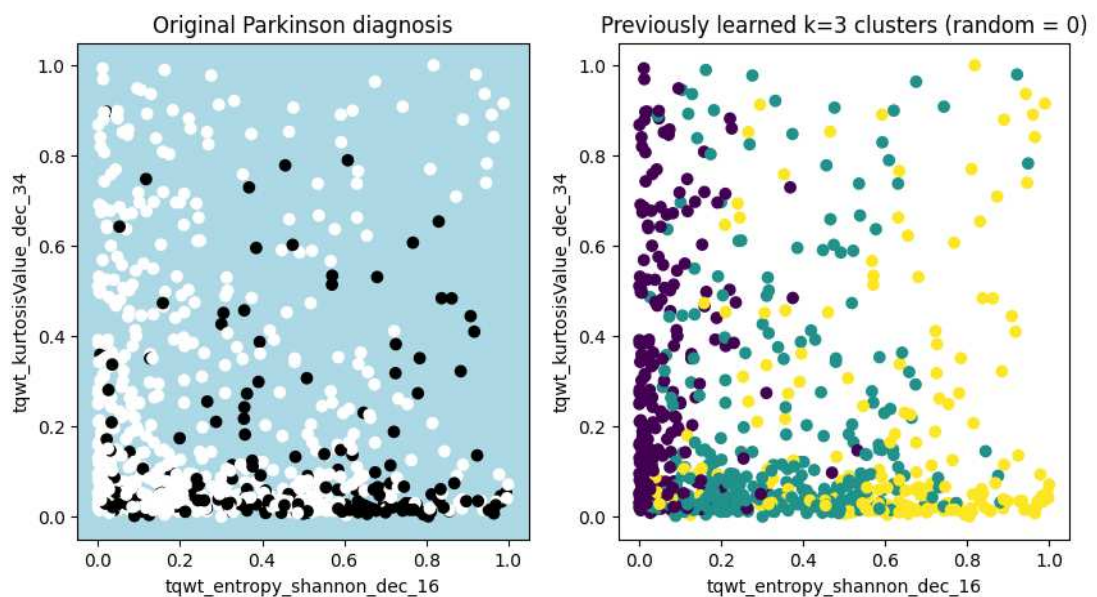
2) In the previous question, we observed that the results of the silhouette and purity were non-deterministic.

The `cluster.KMeans()` function has a parameter called `random_state`. This parameter determines random number generation for centroid initialization and we use an integer in this parameter to make the randomness deterministic.

Different `random_state` values will not prevent the algorithm from converging to the same final point since we can see that seed 0 and seed 2 results are the same.

However, the results for seed 1 are different from the others due to the randomness in the initialization of the centroids and that is what is causing the non-determinism.

3)



4) Number of components: 31

III. APPENDIX

Code used in question 1) of **II. Programming and critical analysis:**

Loads the data:

```
import numpy as np
import pandas as pd
from scipy.io.arff import loadarff
from sklearn.metrics import silhouette_score

# Load the data
data = loadarff('pd_speech.arff')
df = pd.DataFrame(data[0])
df['class'] = df['class'].str.decode('utf-8')

X = df.drop('class', axis=1)
Y = df['class']
```

Calculates the silhouettes and purities:

```
from sklearn import datasets, metrics, cluster, mixture
from sklearn.preprocessing import MinMaxScaler

#normalize the data with MinMaxScaler
scaler = MinMaxScaler()
X = scaler.fit_transform(X)
X = pd.DataFrame(X, columns=df.columns[:-1])

#compute purity
def purity_score(y_true, y_pred):
    # compute contingency/confusion matrix
    confusion_matrix = metrics.cluster.contingency_matrix(y_true, y_pred)
    return np.sum(np.amax(confusion_matrix, axis=0)) / np.sum(confusion_matrix)

seeds = [0, 1, 2]

#apply k-means clustering fully unsupervised (without targets) on the normalized data with
k = 3 and three different seeds (using random  $\epsilon \in \{0,1,2\}$ ).
for seed in seeds:
    #parameterize clustering
    kmeans = cluster.KMeans(n_clusters=3, random_state=seed)

    #learn model
    kmeans_model = kmeans.fit(X)

    # check the produced clusters
    y_pred = kmeans_model.labels_
```

```

print("Seed: ", seed)

print("Silhouette:", metrics.silhouette_score(X, y_pred, metric='euclidean'))

print("Purity:", purity_score(Y, y_pred))

print("\n")

```

Code used in question 3) of **II. Programming and critical analysis** (this code is the continuation of the code used in question 1)):

```

import matplotlib.pyplot as plt

#select the 2 features with the highest variance
Xnew = X.var().sort_values(ascending=False).head(2)

colors = [Y, y_pred]

#plot side by side
fig, ax = plt.subplots(1, 2, figsize=(10, 5))
for i in range(2):
    ax[i].scatter(X[Xnew.index[0]], X[Xnew.index[1]], c=colors[i])
    ax[i].set_xlabel(Xnew.index[0])
    ax[i].set_ylabel(Xnew.index[1])
    if i == 0:
        #background color
        ax[i].set_facecolor('lightblue')
        ax[i].set_title('Original Parkinson diagnosis')
    else:
        ax[i].set_title('Previously learned k=3 clusters (random = 0)')
plt.show()

```

Code used in question 4) of **II. Programming and critical analysis** (this code is the continuation of the code used in question 1) and 3)):

```

from sklearn.decomposition import PCA

#apply PCA
pca = PCA(n_components=0.8)
pca.fit(X)

#number of components
print("Number of components:", pca.n_components_)

```