

# Machine Learning

## Homework 2

系級：醫工三    學號：B12508025    姓名：盧雅筠

1. For each fold of leave-one-out cross-validation

a. Determination of  $k$ , the number of Eigenvalue Selected

i. Methodology: 為了在 Leave-One-Out Cross-Validation 的 103 Folds 中客觀且自動地決定最佳主成分數量  $k$ ，採用 Maximum Distance Method 來尋找 Scree Graph 上的 Elbow Point。Scree Graph 繪製了特徵值  $\lambda_i$  與其對應的主成分順序  $i$ 。曲線的 Elbow Point 代表了「邊際效益遞減」的關鍵位置——在此點之前的主成分包含了數據的主要結構與變異，而此點之後的主成分則多為雜訊。定義 Elbow Point 為曲線上距離「連接首尾特徵值之直線」垂直距離最遠的點。

ii. Mathematical Formulation: PCA 計算出的特徵值序列為  $\lambda_1, \lambda_2, \dots, \lambda_d$ ，將 Scree Graph 視為 2D 平面上的一組座標點集  $P$ ：

$$P_i = (x_i, y_i) = (i, \lambda_i), \quad \text{for } i = 1, \dots, d$$

定義一條連接第一個點  $P_1$  與最後一個點  $P_d$  的直線向量  $\vec{L}$ ：

$$\vec{L} = P_d - P_1 = (d - 1, \lambda_d - \lambda_1)$$

對於曲線上任意一點  $P_i$ ，計算從起點指向該點的向量  $\vec{V}_i$ ：

$$\vec{V}_i = P_i - P_1 = (i - 1, \lambda_i - \lambda_1)$$

為了找出轉折點，計算  $P_i$  到直線  $\vec{L}$  的 Perpendicular Euclidean Distance  $d(P_i, \vec{L})$ 。這可以通過向量的正射影與排斥 Rejection 原理求得：

$$d(P_i, \vec{L}) = \left\| \vec{V}_i - \frac{\vec{V}_i \cdot \vec{L}}{\|\vec{L}\|^2} \vec{L} \right\|$$

最後，最佳的維度  $k$  被定義為具有最大垂直距離的索引值：

$$k = \arg \max_i (dist_i)$$

iii. Code Explanation: 程式碼如 Figure 7，用 `find_elbow_point` 函式來執行上述運算，該函式被應用於 LOOCV 的每一次迭代中。

iv. Report of Selected  $k$  for each of the 103 folds: 結果如 Figure 1。



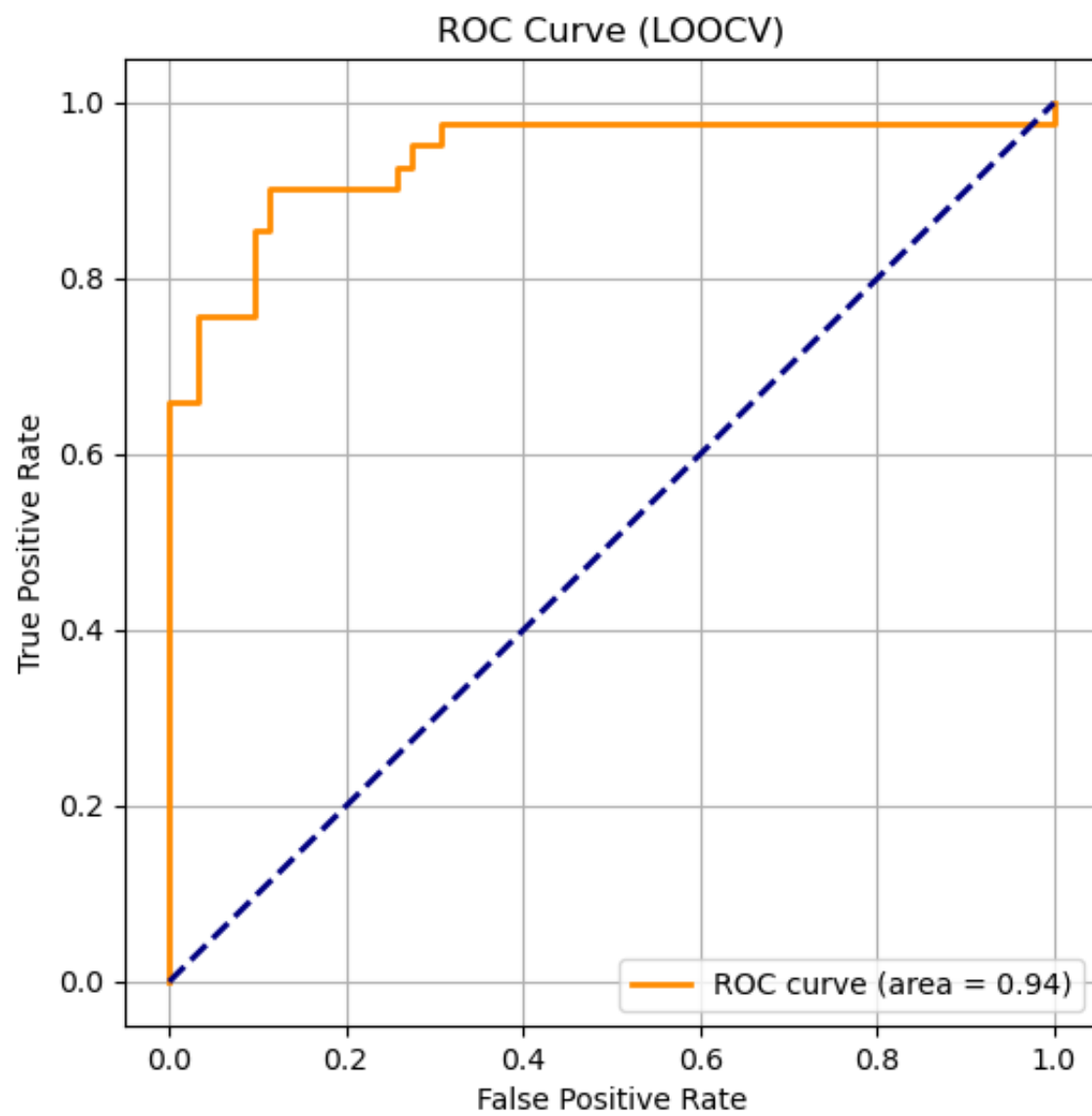


Figure 3

## 2. Use ALL 103 data

- a. Eigenvalue Analysis and Variance Explained: Scree Graph 如 Figure 4 展示了特徵值隨主成分順序的遞減趨勢。可以看到在前幾個主成分後，特徵值迅速下降，顯示數據的主要變異集中在少數幾個維度上。Proportion of Variance Explained 如 Figure 4 展示了前  $m$  個特徵向量所解釋的累積變異量比例。

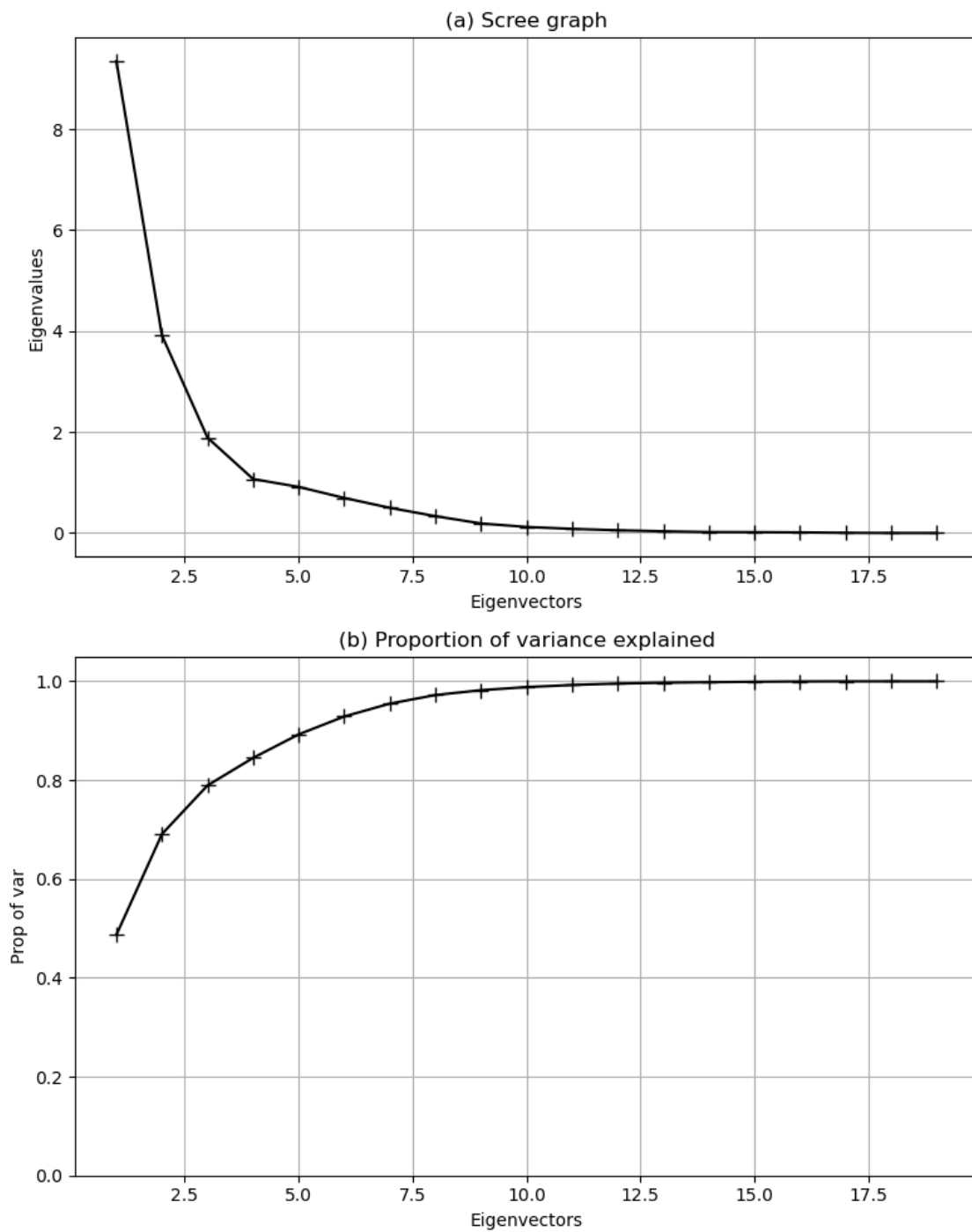


Figure 4

- b. The largest five eigenvalues of PCA results: 首先對數據進行特徵分解，並將特徵值按降序排列。Top 5 Eigenvalues 根據 PCA 分析結果，最大的五個特徵值如 Figure 5 所示。

```
-----  
PART 2: Use ALL 103 data  
-----  
Performance Report  
Top 5 Eigenvalues:  
1: 9.3448  
2: 3.9126  
3: 1.8877  
4: 1.0692  
5: 0.9131
```

Figure 5

- c. Bivariate Gaussian Bayesian Decision Model: 為了將原本的高維數據可視化並進行分類，將所有數據投影到由前兩個特徵向量（對應最大兩個特徵值）所張成的 2D 空間，分別標記為 PC1 與 PC2。在此 2D 子空間中，構建了一個 Bivariate Gaussian Bayesian Decision Model。該模型的視覺化結果如 Figure 6 所示。橫軸與縱軸分別為 PC1 與 PC2。Class 1 以紅色加號 (+) 表示；Class 0 以藍色圓圈 (o) 表示。可以看到兩類數據在 PC1 與 PC2 空間中有一定的分離度。而 Contour Map 中，實線代表 Class 1 的分佈；虛線代表 Class 0 的分佈。最後，Decision Boundary 為黑色粗線，代表  $\Delta = 0$  的軌跡。

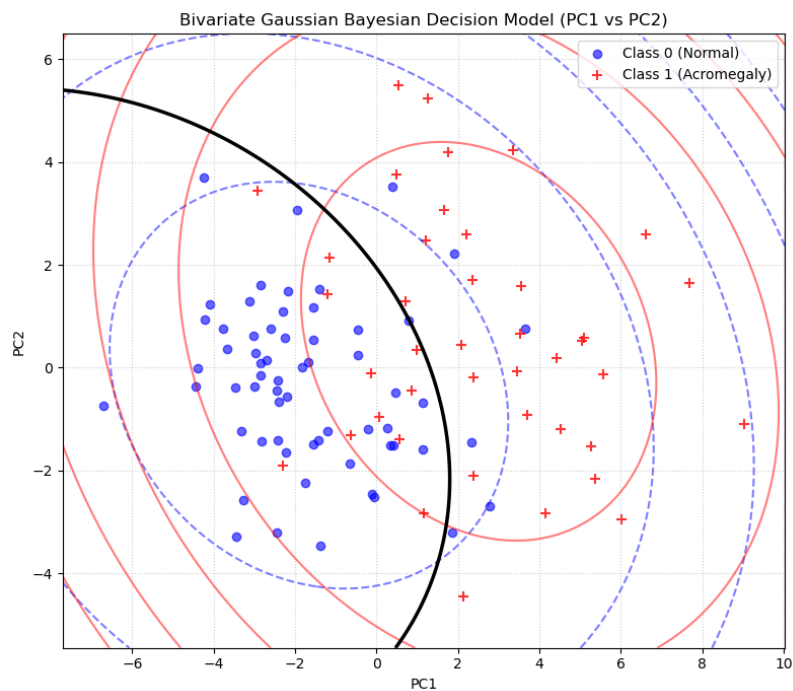


Figure 6

```

1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  from sklearn.metrics import roc_curve, auc, confusion_matrix
5  from sklearn.decomposition import PCA
6  from sklearn.preprocessing import StandardScaler
7
8  ''' Multivariate Gaussian Bayesian classifier '''
9  def Multivariate_Gaussian_Distribution_log_likelihood(X, mu, cov):
10     X = np.atleast_2d(X)
11     d = X.shape[1]
12
13     # Cholesky decomposition, check for positive definiteness
14     try:
15         L = np.linalg.cholesky(cov)
16     except np.linalg.LinAlgError:
17         # If not positive definite, add small ridge to diagonal
18         cov = cov + np.eye(d) * 1e-4
19         L = np.linalg.cholesky(cov)
20
21     Z = np.linalg.solve(L, (X - mu).T)
22     quad = np.sum(Z**2, axis=0)
23     logdet = 2.0 * np.sum(np.log(np.diag(L)))
24     log_likelihood = -0.5 * (d * np.log(2 * np.pi) + logdet + quad)
25     return log_likelihood
26
27 def Bayesian_decision_classifier(X, mu0, mu1, cov0, cov1, p0, p1):
28     ll0 = Multivariate_Gaussian_Distribution_log_likelihood(X, mu0, cov0) + np.log(p0)
29     ll1 = Multivariate_Gaussian_Distribution_log_likelihood(X, mu1, cov1) + np.log(p1)
30     delta = ll1 - ll0
31     m = np.maximum(ll0, ll1)
32     num1 = np.exp(ll1 - m)
33     den = np.exp(ll0 - m) + num1
34     posterior_1 = num1 / den
35     return delta, posterior_1
36
37 ''' Screen Graph of Elbow Method to find optimal k '''
38 def find_elbow_point(eigenvalues):
39     n_points = len(eigenvalues)
40     # if less than 2 points, return 1 as k
41     if n_points < 2:
42         return 1
43
44     all_coords = np.vstack((range(n_points), eigenvalues)).T # each row is [x, y] = [index, eigenvalue]
45     first_point = all_coords[0]
46     last_point = all_coords[-1]
47
48     vec_line = last_point - first_point # vector from first to last point
49     vec_line_norm = vec_line / np.sqrt(np.sum(vec_line**2))
50
51     vec_from_first = all_coords - first_point # vectors from first point to all points
52     scalar_product = np.sum(vec_from_first * vec_line_norm, axis=1) # projection length on line
53     vec_from_first_parallel = np.outer(scalar_product, vec_line_norm) # projection vectors on line(parallel)
54     vec_to_line = vec_from_first - vec_from_first_parallel # vectors from points to line(perpendicular)
55     dist_to_line = np.sqrt(np.sum(vec_to_line ** 2, axis=1)) # distances to line
56
57     # Index of the point with maximum distance to line
58     best_k = np.argmax(dist_to_line) + 1
59
60     return max(1, best_k)

```

```

62 ''' Custom PCA Implementation '''
63 class CustomPCA:
64     def __init__(self, n_components=None):
65         self.n_components = n_components
66         self.components_ = None
67         self.mean_ = None
68         self.explained_variance_ = None
69         self.explained_variance_ratio_ = None
70
71     def fit(self, X):
72         # 1. Centering
73         self.mean_ = np.mean(X, axis=0)
74         X_centered = X - self.mean_
75
76         # 2. Covariance Matrix ~ X^T X
77         cov_matrix = np.cov(X_centered, rowvar=False) # rowvar=False 表示每一行是一個特徵
78
79         # 3. Eigendecomposition: Cv = λv
80         eigenvalues, eigenvectors = np.linalg.eigh(cov_matrix)
81
82         # 4. Sort Eigenvalues and Eigenvectors, descending order
83         idx = np.argsort(eigenvalues)[::-1]
84         self.explained_variance_ = eigenvalues[idx]
85         sorted_vectors = eigenvectors[:, idx]
86
87         # 5. calculate explained variance ratio
88         total_var = np.sum(self.explained_variance_)
89         self.explained_variance_ratio_ = self.explained_variance_ / total_var
90
91         # 6. Select top n_components
92         if self.n_components is None:
93             n_comp = X.shape[1] # all components
94         else:
95             n_comp = min(self.n_components, X.shape[1])
96         self.components_ = sorted_vectors[:, :n_comp].T # W: (n_components, n_features)
97         return self
98
99     def transform(self, X):
100         X_centered = X - self.mean_
101         # Project data onto principal components: X_new = X_centered @ W.T
102         return np.dot(X_centered, self.components_.T)
103
104     def fit_transform(self, X):
105         self.fit(X)
106         return self.transform(X)
107
108
109 ''' Load Data '''
110 data_path = 'AcromegalyFeatureSet.xlsx'
111 data = pd.read_excel(data_path)
112 data.rename(columns=lambda s: s.strip() if isinstance(s, str) else s, inplace=True)
113
114 X_raw = data.drop(columns=['SeqNum', 'Gender', 'GroundTruth']).values
115 y = data['GroundTruth'].values
116 n_samples, n_features = X_raw.shape
117
118 print(f"Data Loaded: {n_samples} samples, {n_features} features.")
119
120 ''' Part1: Leave-One-Out Cross Validation with PCA and Bayesian Classifier '''
121 print("-" * 60)
122 print("PART 1: For each fold of leave-one-out cross-validation")
123 print("-" * 60)
124 posts1 = []
125 deltas = []
126 k_values = []
127 preds = []

```



```

129 for i in range(n_samples):
130     # 1. Split Data
131     X_train_raw = np.delete(X_raw, i, axis=0)
132     y_train = np.delete(y, i)
133     X_test_raw = X_raw[i].reshape(1, -1) # (d,) -> (1, d)
134
135     # 2. Standardization
136     scaler = StandardScaler() # standard to zero mean and unit variance
137     X_train_std = scaler.fit_transform(X_train_raw)
138     X_test_std = scaler.transform(X_test_raw)
139
140     # 3. PCA on Training Data
141     pca_full = CustomPCA(n_components=None) # 計算所有成分
142     pca_full.fit(X_train_std)
143
144     # 4. Find k using Screen Graph of Elbow Method
145     eigenvalues = pca_full.explained_variance_
146     k = find_elbow_point(eigenvalues)
147     k_values.append(k)
148     print(f"Fold {i+1}/{n_samples}: The numbers of eigenvalues selected = {k}")
149
150     # 5. Project Data
151     pca_k = CustomPCA(n_components=k)
152     X_train_pca = pca_k.fit_transform(X_train_std)
153     X_test_pca = pca_k.transform(X_test_std)
154
155     # 6. Train Bayesian Classifier
156     X0 = X_train_pca[y_train == 0]
157     X1 = X_train_pca[y_train == 1]
158
159     mu0 = np.mean(X0, axis=0)
160     mu1 = np.mean(X1, axis=0)
161
162     cov0 = np.cov(X0, rowvar=False) + np.eye(k) * 1e-4
163     cov1 = np.cov(X1, rowvar=False) + np.eye(k) * 1e-4
164
165     p0 = len(X0) / len(X_train_pca)
166     p1 = len(X1) / len(X_train_pca)
167
168     # 7. Test
169     delta, posterior1 = Bayesian_decision_classifier(X_test_pca, mu0, mu1, cov0, cov1, p0, p1)
170
171     # --- FIX START: 使用 .item() 解決 DeprecationWarning ---
172     posterior1_val = posterior1.item()
173     delta_val = delta.item()
174
175     posts1.append(posterior1_val)
176     deltas.append(delta_val)
177     preds.append(1 if delta_val >= 0 else 0)
178     # --- FIX END ---
179
180 # Performance
181 posts1 = np.array(posts1)
182 deltas = np.array(deltas)
183 preds = np.array(preds)
184 k_values = np.array(k_values)
185
186 fpr, tpr, _ = roc_curve(y, posts1)
187 roc_auc = auc(fpr, tpr)
188 cm = confusion_matrix(y, preds)
189 TN, FP, FN, TP = cm.ravel()
190
191 accuracy = (TP + TN) / (TP + TN + FP + FN)
192 sensitivity = TP / (TP + FN) if (TP + FN) > 0 else 0
193 specificity = TN / (TN + FP) if (TN + FP) > 0 else 0

```

```

195 print("-" * 30)
196 print("Performance Report")
197 print(f"Selected k (mean): {np.mean(k_values):.2f}")
198 print(f"Selected k (min/max): {np.min(k_values)} / {np.max(k_values)}")
199 print(f"Confusion Matrix:\n{cm}")
200 print(f"Accuracy: {accuracy:.3f}")
201 print(f"Sensitivity: {sensitivity:.3f}")
202 print(f"Specificity: {specificity:.3f}")
203 print(f"AUC: {roc_auc:.3f}")
204
205 # ROC curve plot
206 plt.figure(figsize=(6, 6))
207 plt.plot(fpr, tpr, color='darkorange', lw=2, label=f'ROC curve (area = {roc_auc:.2f})')
208 plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
209 plt.xlabel('False Positive Rate')
210 plt.ylabel('True Positive Rate')
211 plt.title('ROC Curve (LOOCV)')
212 plt.legend(loc="lower right")
213 plt.grid(True)
214 plt.savefig('homework2_roc.png')
215
216 ''' Part2: Full Data Analysis with PCA and Bayesian Classifier '''
217 print("\n" + "-" * 60)
218 print("PART 2: Use ALL 103 data")
219 print("-" * 60)
220
221 scaler_all = StandardScaler()
222 X_std_all = scaler_all.fit_transform(X_raw)
223
224 # 使用 CustomPCA 計算所有資料
225 pca_all = CustomPCA(n_components=None)
226 pca_all.fit(X_std_all)
227
228 eigenvalues_all = pca_all.explained_variance_
229 ratios = pca_all.explained_variance_ratio_
230 cumulative_ratio = np.cumsum(ratios)
231 print("Performance Report")
232 print("Top 5 Eigenvalues:")
233 for i, val in enumerate(eigenvalues_all[:5]):
234     print(f" {i+1}: {val:.4f}")
235
236 fig, (ax1, ax2) = plt.subplots(2, 1, figsize=(8, 10))
237
238 ax1.plot(range(1, len(eigenvalues_all)+1), eigenvalues_all, 'k--', markersize=8)
239 ax1.set_title('(a) Scree graph')
240 ax1.set_ylabel('Eigenvalues')
241 ax1.set_xlabel('Eigenvectors')
242 ax1.grid(True)
243
244 ax2.plot(range(1, len(cumulative_ratio)+1), cumulative_ratio, 'k--', markersize=8)
245 ax2.set_title('(b) Proportion of variance explained')
246 ax2.set_ylabel('Prop of var')
247 ax2.set_xlabel('Eigenvectors')
248 ax2.set_ylim([0, 1.05])
249 ax2.grid(True)
250
251 plt.tight_layout()
252 plt.savefig('homework2_pca_scree.png')
253 # plt.show()

```

```

255 # Plotting 2D Bivariate Gaussian Decision Boundary
256 pca_2d = PCA(n_components=2)
257 X_pca_2d = pca_2d.fit_transform(X_std_all)
258
259 X0_2d = X_pca_2d[y == 0]
260 X1_2d = X_pca_2d[y == 1]
261
262 mu0_2d = np.mean(X0_2d, axis=0)
263 mu1_2d = np.mean(X1_2d, axis=0)
264 cov0_2d = np.cov(X0_2d, rowvar=False) + np.eye(2) * 1e-4
265 cov1_2d = np.cov(X1_2d, rowvar=False) + np.eye(2) * 1e-4
266 p0_2d = len(X0_2d) / len(X_pca_2d)
267 p1_2d = len(X1_2d) / len(X_pca_2d)
268
269 x_min, x_max = X_pca_2d[:, 0].min() - 1, X_pca_2d[:, 0].max() + 1
270 y_min, y_max = X_pca_2d[:, 1].min() - 1, X_pca_2d[:, 1].max() + 1
271 xx, yy = np.meshgrid(np.linspace(x_min, x_max, 200),
272                      np.linspace(y_min, y_max, 200))
273 pts = np.c_[xx.ravel(), yy.ravel()]
274
275 LL0 = Multivariate_Gaussian_Distribution_log_likelihood(pts, mu0_2d, cov0_2d).reshape(xx.shape)
276 LL1 = Multivariate_Gaussian_Distribution_log_likelihood(pts, mu1_2d, cov1_2d).reshape(xx.shape)
277
278 boundary = (LL1 + np.log(p1_2d)) - (LL0 + np.log(p0_2d))
279
280 plt.figure(figsize=(8, 7))
281 plt.scatter(X0_2d[:, 0], X0_2d[:, 1], c='blue', marker='o', label='Class 0 (Normal)', alpha=0.6)
282 plt.scatter(X1_2d[:, 0], X1_2d[:, 1], c='red', marker='+', s=60, label='Class 1 (Acromegaly)', alpha=0.8)
283
284 plt.contour(xx, yy, LL0, levels=5, colors='blue', linestyle='dashed', alpha=0.5)
285 plt.contour(xx, yy, LL1, levels=5, colors='red', linestyle='solid', alpha=0.5)
286 plt.contour(xx, yy, boundary, levels=[0], colors='black', linewidths=2.5, linestyle='-')
287
288 plt.title('Bivariate Gaussian Bayesian Decision Model (PC1 vs PC2)')
289 plt.xlabel('PC1')
290 plt.ylabel('PC2')
291 plt.legend(loc='best')
292 plt.grid(True, linestyle=':', alpha=0.6)
293 plt.tight_layout()
294 plt.savefig('homework2_2d_boundary.png')

```

Figure 1