

Introduction to Machine Learning - Project

Talla DIAGNE

12 novembre 2024

Table des matières

Introduction	1
1 Dataset exploration and preprocessing	1
2 Principal Component Analysis (PCA)	2
3 Training and evaluation	3
Results	6
Conclusion	7

Introduction to Machine Learning Assignment

Medical Image Diagnosis Enhancement Using PCAA and KNN

Talla DIAGNE

Nous attestons que ce travail est original, que nous citons en référence toutes les sources utilisées et qu'il ne comporte pas de plagiat.

Introduction

Medical image analysis plays a crucial role in diagnosing diseases and disorders. In this report, we present an analysis of medical images using Principal Component Analysis (PCA) for dimensionality reduction and K-Nearest Neighbors (KNN) for classification. The goal is to develop a model that can accurately classify medical images into different categories, aiding in diagnosis and treatment planning.

1 Dataset exploration and preprocessing

Understanding the dataset is the first step in any machine learning task. Our dataset consists of medical images (MRI brain scans). We associate each image with a label indicating the corresponding medical condition or anomaly. We explored the dataset to gain insights into its structure and characteristics. This exploration involved determining the number of images, image dimensions, and the distribution of labels across the dataset. The dataset contains images from various medical conditions, including glioma, meningioma, pituitary tumors, and normal scans.

```
1 Train dataset :
2 Number of images: 5712
3 Image dimensions: (512, 512, 3)
4 Label distribution: {'glioma': 1321, 'meningioma': 1339, 'notumor': 1595, 'pituitary': 1457}
```

Before applying PCA and KNN, we preprocess the images to ensure they are suitable for analysis. The preprocessing steps include :

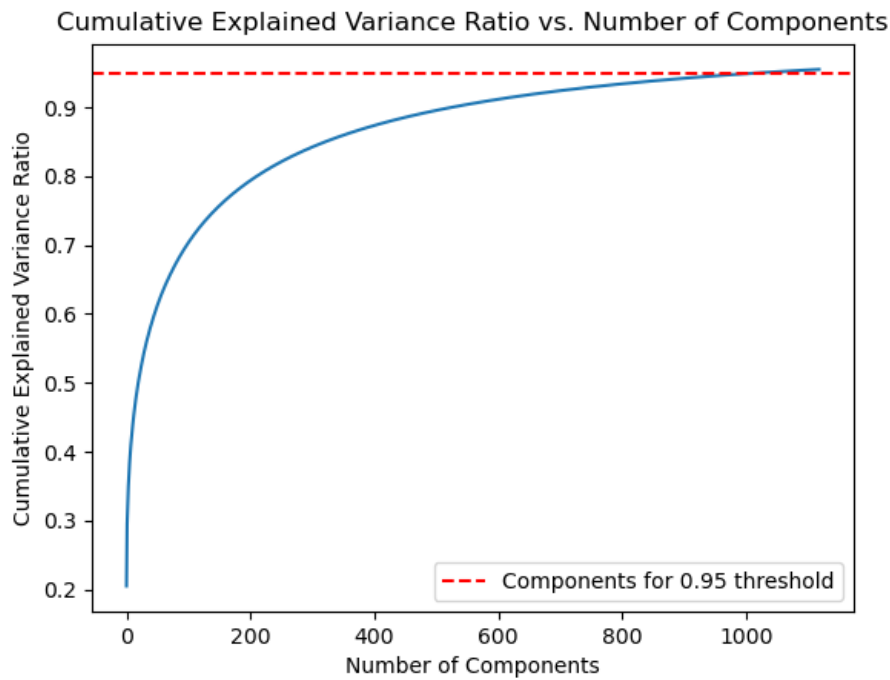
- Loading the images from the provided folders and subfolders
- Extracting the labels from the subfolder names
- Converting the images to grayscale to simplify the analysis

- Resizing the images to a consistent size of 100x100 pixels to standardize all the images and flattening them to prepare it for PCA
- Normalizing the pixel values to be within the range $[0, 1]$ to facilitate computation

2 Principal Component Analysis (PCA)

PCA (Principal Component Analysis) is a technique used to reduce the dimensionality of high-dimensional data while preserving its important features. In our analysis, we applied PCA to preprocess the image data, aiming to decrease its dimensionality while retaining as much variance as possible. To determine the optimal number of principal components, we relied on the explained variance ratio. This ratio indicates the proportion of variance in the data explained by each principal component. By setting a threshold, typically at 95%, we aimed to retain enough variance to accurately represent the data. After analysis, we found that retaining 1119 principal components met our threshold of capturing 95% of the variance. This selection strikes a balance between reducing dimensionality and preserving essential information necessary for subsequent tasks like classification. Including more components may better preserve the informations in the data but increases complexity; selecting too few components may lead to significant information loss.

However, this result can seem a bit strange when we know that, later in this documents, we will plot the decision boundaries of the KNN classifier but only by using the first 2 parameters.



In order to evaluate the effectiveness of PCA in retaining essential image features, we visually compared original medical images with their reconstructed counterparts. This comparison helps assess how well PCA preserves crucial information during dimensionality reduction, guiding the determination of optimal principal components for accurate classification

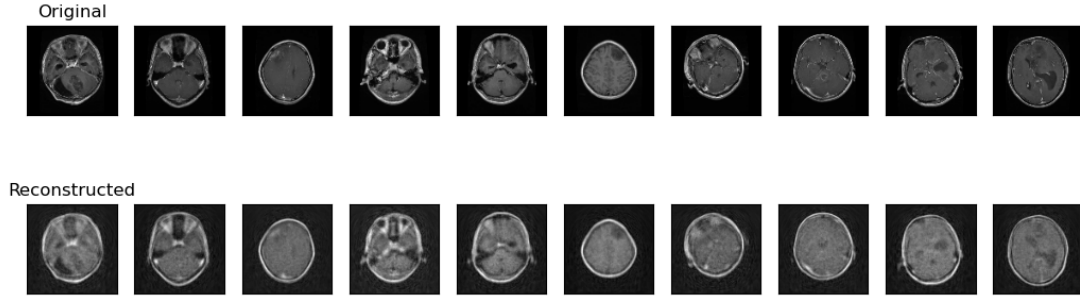


FIGURE 1 – Comparison between original and reconstructed medical images

3 Training and evaluation

With the preprocessed data and the reduced feature space obtained from PCA, we train a KNN classifier. KNN is an effective algorithm for classification tasks, particularly when dealing with labeled datasets. To optimize the performance of the KNN classifier, we explore the impact of different values of the hyperparameter k , which represents the number of neighbors used in the classification. We use a technique called cross-validation to find the optimal value of k that maximizes the classifier's performance. One way to visualize the relationship between the value of k and the classifier's performance is to plot the accuracy against the number of neighbors.

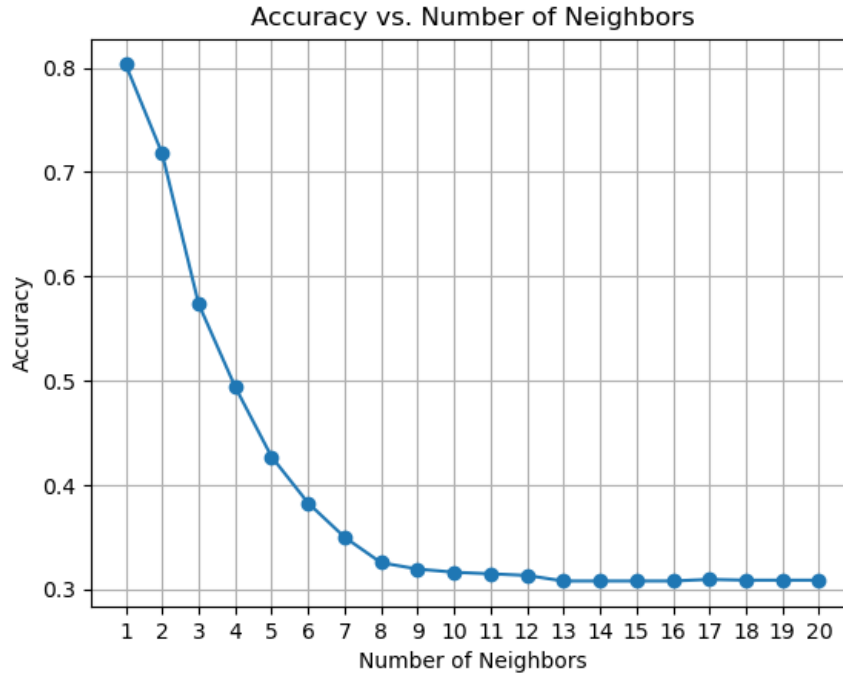


FIGURE 2 – Program performance depending on the number of neighbors

This plot [2] helps us understand how the choice of k affects the accuracy of the classifier. By analyzing this curve, we can determine the optimal value of k that yields the highest accuracy on the test data. We also decided to keep $k = 2$, because the accuracy is very good.

To visualize the performance of the KNN classifier, decision boundaries were additionally plotted for various values of k . Decision boundaries delineate the classification boundaries between different classes within the feature space. By

plotting these boundaries for different k values, the model's complexity can be observed as it varies with the number of neighbors considered for classification. For each k value, data points were colored according to their actual class. Subsequently, the feature space was divided into regions colored based on the class predicted by the KNN classifier. This visualization provides insight into how decision boundaries change with varying k values and their impact on the model's generalization to unseen data.

The uniformity of decision boundaries in the plot can vary with different values of 'k' in the k-nearest neighbors (KNN) algorithm. A smaller 'k' tends to result in more irregular and fragmented decision boundaries, as the classifier relies on a smaller number of nearest neighbors for prediction, making it more sensitive to noise in the data. Conversely, larger values of 'k' lead to smoother and more uniform decision boundaries, as the classifier considers a larger number of neighbors, resulting in a more generalized representation of the feature space. This variation in decision boundary uniformity highlights the relation between bias and variance in the KNN algorithm and emphasizes the importance of selecting an appropriate value of 'k' based on the complexity of the dataset and desired model behavior. However, the result here seems a bit strange to us because with the optimal value found before ($k=2$), the image is not the smoothest one.

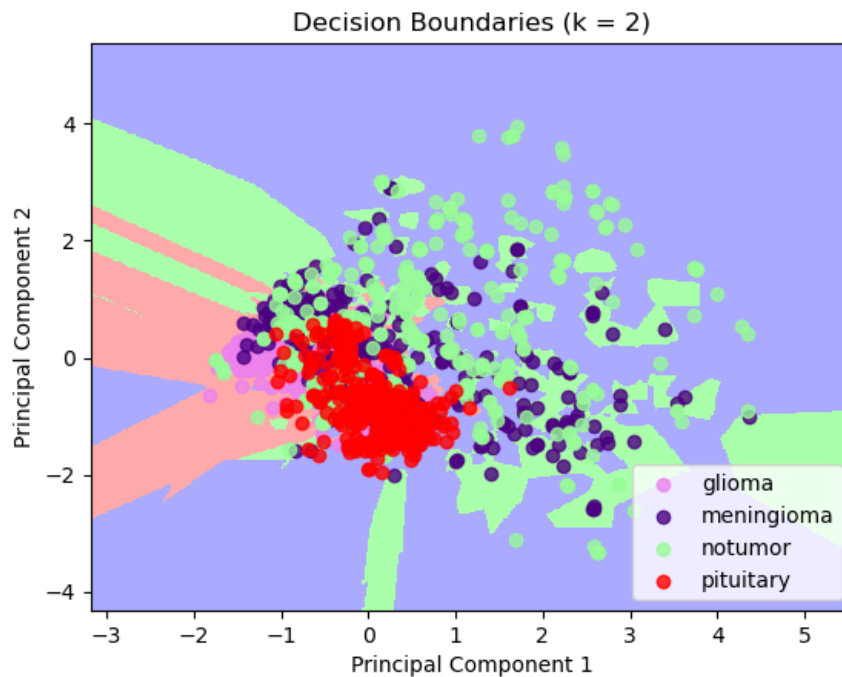


FIGURE 3 – Accuracy : 0.7208237986270023

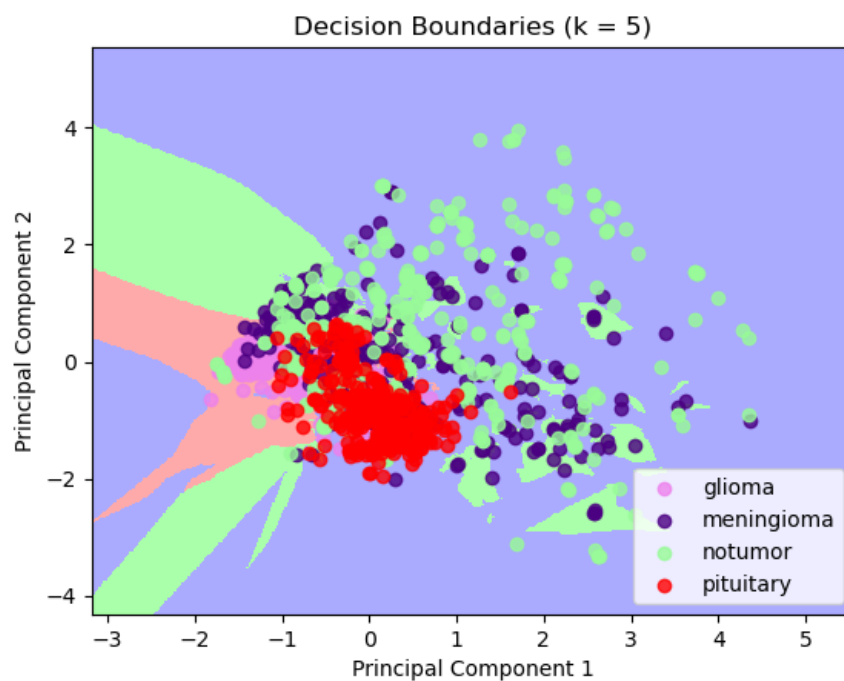


FIGURE 4 – Accuracy : 0.6071700991609459

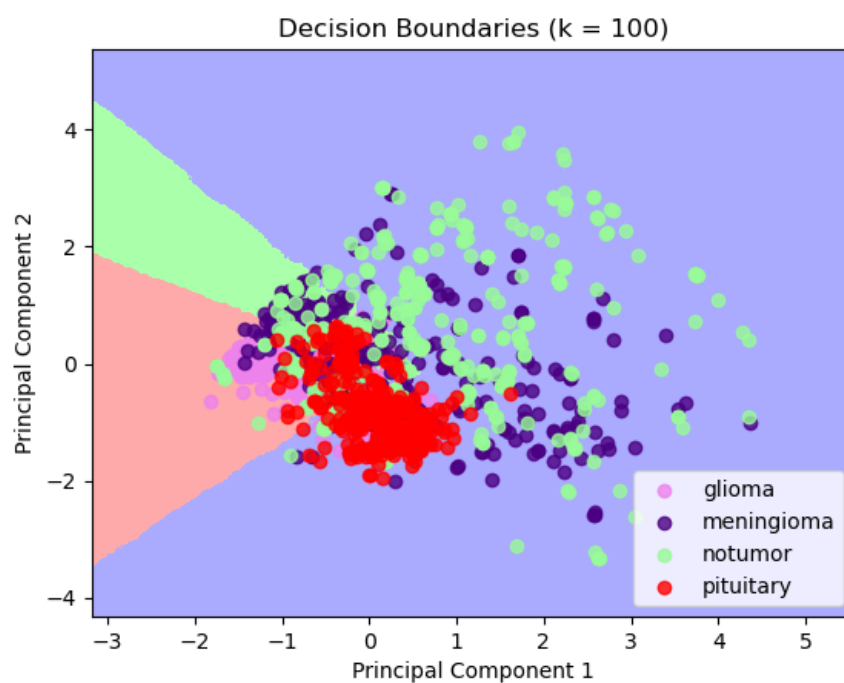


FIGURE 5 – Accuracy : 0.6071700991609459

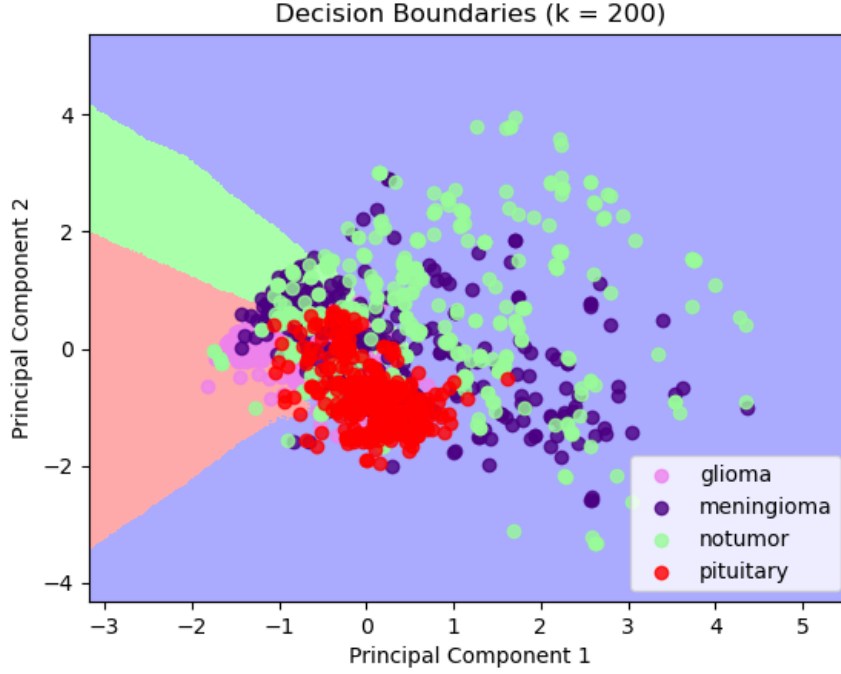


FIGURE 6 – Accuracy : 0.6048817696414951

Analyzing the decision boundary plots for different k values enables the identification of regions where the classifier may be uncertain and prone to classification errors. This aids in understanding the model's limitations and facilitates informed decision-making regarding the selection of the optimal k value to enhance classifier performance.

Results

Task 1 : Dataset Exploration Number of images : 5712 Image dimensions : (512, 512, 3) Label distribution : 'glioma' : 1321, 'meningioma' : 1339, 'notumor' : 1595, 'pituitary' : 1457

K = 2 :

- Accuracy : 0.7208237986270023
- Precision : 0.7905195335685307
- Recall : 0.7147215865751335
- F1-score : 0.7155498319658

Sure, here are the tables rewritten with the provided values :

TABLE 1 – Classification kNN Train Report

	Precision	Recall	F1-score	Support
0	1.00	1.00	1.00	1321
1	1.00	1.00	1.00	1339
2	1.00	1.00	1.00	1595
3	1.00	1.00	1.00	1457
Accuracy : 1.00, Macro avg : 1.00, Weighted avg : 1.00, Total support : 5712				

TABLE 2 – Classification kNN Test Report

	Precision	Recall	F1-score	Support
0	0.68	0.56	0.62	300
1	0.66	0.74	0.70	306
2	0.89	0.97	0.93	405
3	0.76	0.69	0.73	300
Accuracy : 0.76, Macro avg : 0.75, Weighted avg : 0.76, Total support : 1311				

Conclusion

In conclusion, the application of PCA and KNN for medical image analysis showcases a promising way to improve diagnostic accuracy and treatment planning. By effectively preprocessing the images, reducing dimensionality with PCA, and training a KNN classifier, we achieve notable results in classifying medical images into different conditions. However, further optimization and fine-tuning of parameters could enhance the classification performance even more. Exploring alternative algorithms and techniques, such as deep learning approaches, may provide valuable insights for future research in medical image analysis. Moreover, the visualizations of decision boundaries shed light on the model's behavior and highlight the importance of selecting an appropriate value for the hyperparameter k . Understanding the linkss between bias and variance in the KNN algorithm is crucial for achieving optimal performance and generalization to unseen data. In summary, while PCA and KNN offer a solid foundation for medical image analysis, ongoing research and experimentation are essential for pushing the boundaries of diagnostic accuracy and contributing to advancements in healthcare technology.

For the full code implementation and detailed analysis, please refer to the attached source code and documentation.

Link for the code on Google Colab :

<https://colab.research.google.com/drive/1tKH0B46Nh3cwUqRWU6YFNKZbfSw040Ne?usp=sharing>

```
1 import os
2 import numpy as np
3 import cv2
4 import matplotlib
5 # matplotlib.style.use('ggplot')
6 from sklearn.decomposition import PCA
7 from sklearn.model_selection import train_test_split
8 from sklearn.neighbors import KNeighborsClassifier
9 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
10 import matplotlib.pyplot as plt
11 import torch
12 from matplotlib.colors import ListedColormap
13 from sklearn.model_selection import train_test_split, GridSearchCV
14 from sklearn.metrics import classification_report, mean_squared_error
15 from sklearn.metrics import r2_score
16 from sklearn.linear_model import LogisticRegression
17 from sklearn.linear_model import Perceptron
18 from sklearn.tree import DecisionTreeClassifier
19 from sklearn.neighbors import KNeighborsRegressor
20 import pandas
21
22 # Task 1: Dataset Exploration
23 def explore_dataset(dataset_folder):
24     labels = os.listdir(dataset_folder)
25     num_images = 0
26     dimensions = None
27     label_distribution = {}
28
29     for label in labels:
30         label_folder = os.path.join(dataset_folder, label)
31         images = os.listdir(label_folder)
32         num_images += len(images)
33         label_distribution[label] = len(images)
34         if dimensions is None:
35             image_path = os.path.join(label_folder, images[0])
36             image = cv2.imread(image_path)
37             dimensions = image.shape
38
39     print("Number of images:", num_images)
40     print("Image dimensions:", dimensions)
41     print("Label distribution:", label_distribution)
42
43
44 # Task 2: Data Preprocessing
45 def preprocess_images(dataset_folder, image_size=(100, 100)):
46     images = []
47     labels = []
48     label_dict = {'glioma': 0, 'meningioma': 1, 'notumor': 2, 'pituitary': 3}
49     for label in os.listdir(dataset_folder):
50         label_folder = os.path.join(dataset_folder, label)
51         for image_file in os.listdir(label_folder):
52             image_path = os.path.join(label_folder, image_file)
53             image = cv2.imread(image_path)
54             image = cv2.resize(image, image_size)
55             image2 = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)
56             # image = image.flatten()
57             image = cv2.resize(image2, (100, 100))
58             image = image / 255.0
59             images.append(image)
60             labels.append(label_dict[label])
61     # images = np.array(images).view(images.shape[0], -1)
```



```

62     images = np.array(images).flatten().reshape(len(images), 1 * 100 * 100)
63     return images, np.array(labels)
64
65
66 # Task 3: Principal Component Analysis (PCA)
67 def apply_pca(X_train, X_test, n_components):
68     pca = PCA(n_components=n_components, svd_solver='randomized', whiten=True)
69     pca.fit(X_train)
70     X_train_pca = pca.transform(X_train)
71     X_test_pca = pca.transform(X_test)
72     return pca, X_train_pca, X_test_pca
73
74 def determine_components_variance_ratio(pca, threshold=0.95):
75     explained_variance_ratio = pca.explained_variance_ratio_
76     cumulative_variance_ratio = np.cumsum(explained_variance_ratio)
77     num_components = np.argmax(cumulative_variance_ratio >= threshold) + 1
78     plt.plot(cumulative_variance_ratio)
79     plt.xlabel('Number of Components')
80     plt.ylabel('Cumulative Explained Variance Ratio')
81     plt.title('Cumulative Explained Variance Ratio vs. Number of Components')
82     plt.axhline(y=threshold, color='r', linestyle='--', label=f'Components for {threshold} threshold')
83     plt.legend()
84     plt.show()
85     return num_components
86
87 def plot_gallery(images, h, w, n_row=3, n_col=4):
88     plt.figure(figsize=(1.8 * n_col, 2.4 * n_row))
89     plt.subplots_adjust(bottom=0, left=.01, right=.99, top=.90, hspace=.35)
90     for i in range(n_row * n_col):
91         plt.subplot(n_row, n_col, i + 1)
92         plt.imshow(images[i].reshape((h, w)), cmap=plt.cm.gray)
93         plt.xticks(())
94         plt.yticks(())
95
96
97 # Task 4: Training and Evaluation
98 def train_and_evaluate(X_train, X_test, y_train, y_test, n_neighbors):
99     knn = KNeighborsClassifier(n_neighbors=n_neighbors)
100     knn.fit(X_train, y_train)
101     y_pred = knn.predict(X_test)
102     accuracy = accuracy_score(y_test, y_pred)
103     precision = precision_score(y_test, y_pred, average='weighted')
104     recall = recall_score(y_test, y_pred, average='weighted')
105     f1 = f1_score(y_test, y_pred, average='weighted')
106     # Afficher les métriques d'évaluation
107     print("Accuracy:", accuracy)
108     print("Precision:", precision)
109     print("Recall:", recall)
110     print("F1-score:", f1)
111     return knn
112
113
114 # Task 5: Visualization
115 def visualize_decision_boundaries(X_train, y_train, X_test, y_test, n_neighbors, knn):
116     cmap_light = ListedColormap(['#FFAAAA', '#AAFFAA', '#AAAAFF'])
117     cmap_bold = ListedColormap(['#FF0000', '#00FF00', '#0000FF'])
118     clf = KNeighborsClassifier(n_neighbors=n_neighbors)
119     clf.fit(X_train, y_train)
120     y_pred = clf.predict(X_test)
121     accuracy = accuracy_score(y_test, y_pred)
122     print("Accuracy:", accuracy)
123     label_dict = {'glioma': 0, 'meningioma': 1, 'notumor': 2, 'pituitary': 3}
124     inverse_label_dict = {v: k for k, v in label_dict.items()}

```

```

125 x_min, x_max = X_train[:, 0].min() - 1, X_train[:, 0].max() + 1
126 y_min, y_max = X_train[:, 1].min() - 1, X_train[:, 1].max() + 1
127 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.02),
128                      np.arange(y_min, y_max, 0.02))
129 Z = clf.predict(np.c_[xx.ravel(), yy.ravel()])
130 Z = Z.reshape(xx.shape)
131 plt.figure()
132 plt.pcolormesh(xx, yy, Z, cmap=cmap_light)
133 colors = {0: 'violet', 1: 'indigo', 2: 'palegreen', 3: 'red'}
134 ## plt.scatter(X_train_pca[:, 0], X_train_pca[:, 1], c=y_train, cmap=cmap_bold)
135 # plt.scatter(X_test_pca[:, 0], X_test_pca[:, 1], c=y_test, cmap=cmap_bold, marker='x')
136 for label in np.unique(y_test):
137     indices = np.where(y_test == label)
138     plt.scatter(X_test[indices, 0], X_test[indices, 1], c=colors[label], alpha=0.8,
139               label='{}'.format(inverse_label_dict[label]))
140 plt.legend(loc='lower right')
141 plt.xlim(xx.min(), xx.max())
142 plt.ylim(yy.min(), yy.max())
143 plt.title("Decision Boundaries (k = %i)" % n_neighbors)
144 plt.xlabel('Principal Component 1')
145 plt.ylabel('Principal Component 2')
146 plt.show()
147
148 def original_and_reconstructed(images, images_init, pca, knn, labels):
149     print(images.shape)
150     h, w = 100, 100
151     reconstructed_images = pca.inverse_transform(images)
152     n = 10
153     plt.figure(figsize=(20, 4))
154     for i in range(n):
155         ax = plt.subplot(2, n, i + 1)
156         plt.imshow((images_init[i]).reshape((h, w)), cmap='gray')
157         ax.get_xaxis().set_visible(False)
158         ax.get_yaxis().set_visible(False)
159         if i == 0:
160             ax.set_title("Original")
161         ax = plt.subplot(2, n, i + 1 + n)
162         plt.imshow((reconstructed_images[i]).reshape((h, w)), cmap='gray')
163         ax.get_xaxis().set_visible(False)
164         ax.get_yaxis().set_visible(False)
165         if i == 0:
166             ax.set_title("Reconstructed")
167     plt.show()
168
169 def classification_report_function(X_train, y_train, X_test, y_test, n_neighbors):
170     model = KNeighborsClassifier(algorithm="auto")
171     parameters = {"n_neighbors": [n_neighbors],
172                  "weights": ["distance"]}
173     model_optim = GridSearchCV(model, parameters, cv=5, scoring="accuracy")
174     model_optim.fit(X_train, y_train)
175     model_optim.best_estimator_
176     for (i, x, y) in zip(["Train", "Test"], [X_train, X_test], [y_train, y_test]):
177         print("Classification kNN", i, " report:\n", classification_report(y, model_optim.predict(x)))
178
179
180 def plot_accuracy_vs_neighbors(X_train_pca, y_train, X_test_pca, y_test, max_neighbors):
181     accuracies = []
182     neighbors_range = range(1, max_neighbors + 1)
183     for n_neighbors in neighbors_range:
184         knn_classifier = train_and_evaluate(X_train, X_test, y_train, y_test, n_neighbors)
185         y_pred = knn_classifier.predict(X_test)
186         accuracy = accuracy_score(y_test, y_pred)
187         accuracies.append(accuracy)

```

```

188 plt.plot(neighbors_range, accuracies, marker='o')
189 plt.title('Accuracy vs. Number of Neighbors')
190 plt.xlabel('Number of Neighbors')
191 plt.ylabel('Accuracy')
192 plt.xticks(neighbors_range)
193 plt.grid(True)
194 plt.show()
195
196 if __name__ == "__main__":
197     training_folder = r'D:\2A\Machine Learning\Medical Diagnosis\Training'
198     testing_folder = r'D:\2A\Machine Learning\Medical Diagnosis\Testing'
199     n_neighbors = 2
200
201     # Task 1: Dataset Exploration
202     print("Task 1: Dataset Exploration")
203     explore_dataset(training_folder)
204
205     # Task 2: Data Preprocessing
206     print("\nTask 2: Data Preprocessing")
207     X_train, y_train = preprocess_images(training_folder)
208     X_test, y_test = preprocess_images(testing_folder)
209
210     # Task 3: PCA
211     print("\nTask 3: Principal Component Analysis (PCA)")
212     n_components = 1119 # found with determine_components_variance_ratio
213     pca, X_train_pca, X_test_pca = apply_pca(X_train, X_test, n_components)
214     print("Dimension de X_train_pca:", X_train_pca.shape)
215     print("Dimension de X_test_pca:", X_test_pca.shape)
216     print("Proportion de variance expliquée par chaque composante principale:")
217     explained_variance_ratio = pca.explained_variance_ratio_
218     print(explained_variance_ratio)
219     print("Dimension de X_train_pca:", X_train_pca.shape)
220     nb_of_components = determine_components_variance_ratio(pca, threshold=0.95)
221     print(nb_of_components)
222     X_test_pca = pca.transform(X_test)
223
224     # Task 4: Training and Evaluation
225     print("\nTask 4: Training and Evaluation")
226     print(X_test_pca.shape)
227     knn_classifier = train_and_evaluate(X_train_pca, X_test_pca, y_train, y_test, n_neighbors
228 )
229
230     # Task 5: Visualization
231     print("\nTask 5: Visualization")
232     original_and_reconstructed(X_test_pca, X_test, pca, knn_classifier, y_test)
233     visualize_decision_boundaries(X_train_pca[:, :2], y_train, X_test_pca[:, :2], y_test,
234 n_neighbors, knn_classifier)
235
236     # Find optimal K
237     classification_report_function(X_train_pca[:, :2], y_train, X_test_pca[:, :2], y_test,
238 n_neighbors)
239
240     # Plot_accuracy_vs_neighbors
241     max_neighbors = 20
242     plot_accuracy_vs_neighbors(X_train_pca, y_train, X_test_pca, y_test, max_neighbors)

```