

CS 5785 Homework 3

Cornell Tech

Tian Jin^{*†}

Yufan Zhang^{*‡}

October 2023

1 Programming Exercises

1.1 Eigenface for face recognition

Part (a) Download the dataset

Part (b) Load data and plot face

The code in Listing 1 is implemented to load the datasets and plot the sample faces from both the training and test datasets, which are shown in Figure 1.

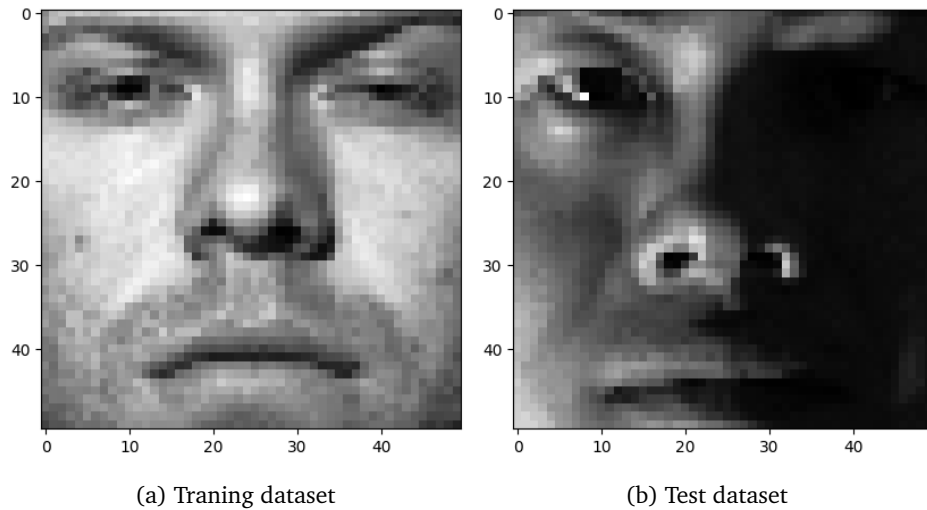


Figure 1: Sample face images from the dataset

^{*}Both authors contributed equally to this project.

[†]tj299@cornell.edu

[‡]yz2894@cornell.edu

```

1  # Loading training data
2  train_labels, train_data = [], []
3  for line in open('./faces/train.txt'):
4      im = imageio.imread(line.strip().split()[0])
5      train_data.append(im.reshape(2500,))
6      train_labels.append(line.strip().split()[1])
7  train_data = np.array(train_data, dtype=float)
8  train_labels = np.array(train_labels, dtype=int)
9
10 print(train_data.shape, train_labels.shape)
11 plt.imshow(train_data[10, :].reshape(50,50), cmap = cm.Greys_r)
12 plt.show()
13
14 # Loading test data
15 test_labels, test_data = [], []
16 for line in open('./faces/test.txt'):
17     im = imageio.imread(line.strip().split()[0])
18     test_data.append(im.reshape(2500,))
19     test_labels.append(line.strip().split()[1])
20 test_data = np.array(test_data, dtype=float)
21 test_labels = np.array(test_labels, dtype=int)
22
23 print(test_data.shape, test_labels.shape)
24 plt.imshow(test_data[10, :].reshape(50,50), cmap = cm.Greys_r)
25 plt.show()

```

Listing 1: Load the datasets and plot sample face

Part (c) Average face

The code in Listing 2 is implemented to get the average face μ from the whole training set, which is shown in Figure 2.

```

1  miu = train_data.sum(axis=0) / train_data.shape[0]
2  plt.imshow(miu.reshape(50,50), cmap = cm.Greys_r)
3  plt.show()

```

Listing 2: Getting the average face

Part (d) Mean subtraction

The code in Listing 3 is implemented to get the faces that were calculated by subtraction of the average face μ . The resulting faces after mean subtraction for both the training and test datasets are shown in Figure 3.

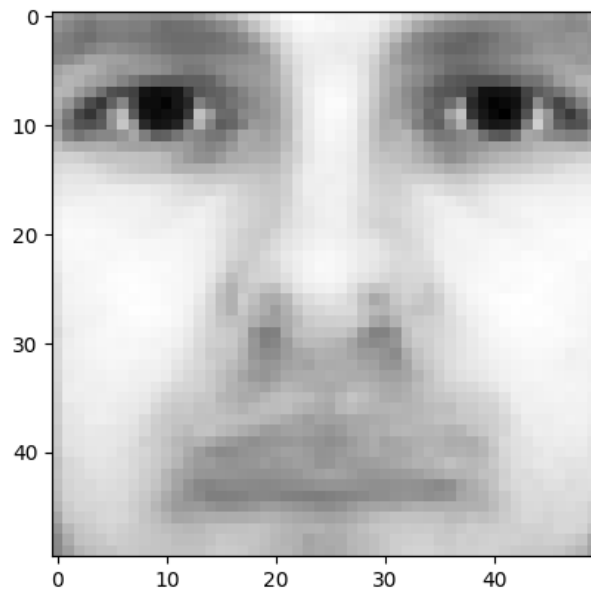


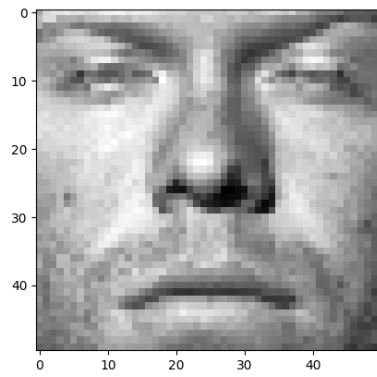
Figure 2: The calculated average face from the whole training set

```
1 new_train = train_data - miu
2 plt.imshow(new_train[10, :].reshape(50,50), cmap = cm.Greys_r)
3 plt.savefig('img/mean_sub_train.png')
4 plt.show()
5
6 new_test = test_data - miu
7 plt.imshow(new_test[10, :].reshape(50,50), cmap = cm.Greys_r)
8 plt.savefig('img/mean_sub_test.png')
9 plt.show()
```

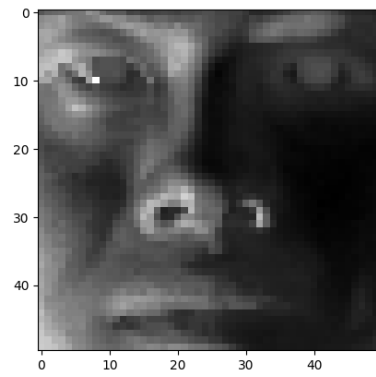
Listing 3: Getting the average face

Part (e) Eigenface

The code in Listing 4 is implemented to get the eigenfaces as instructed in the Homework description, whose results are visualized in Figure 4.



(a) Training dataset



(b) Test dataset

Figure 3: Face images after mean subtraction



(a) Image 9

(b) Image 10

Figure 4: The first ten eigenfaces in greyscale

```

1 def compute_V_T(X):
2     # np.linalg.eig(): the column eigenvectors[:, i] is
3     # the eigenvector corresponding to the eigenvalue eigenvalues[i]
4     eigenvalues, eigenvectors = np.linalg.eig(np.dot(X.T, X))
5     V_T = eigenvectors.T
6     return V_T
7
8 train_V_T = compute_V_T(new_train)
9 for i in range(10):
10     eigenface = np.real(train_V_T[i,:])
11     plt.imshow(eigenface.reshape(50,50), cmap = cm.Greys_r)
12     plt.savefig('img/eigenface_' + str(i) + '.png')
13     plt.show()

```

Listing 4: Getting the eigenfaces

Part (f) Eigenface feature

The code in Listing 5 and 6 are the functions defined to generate r -dimensional feature matrix for training images and test images, respectively.

```

1 train_V_T = compute_V_T(new_train)
2 def getTrainF(r):
3     transpose = train_V_T[:r, :].T
4     f = np.dot(train_data, transpose)
5     return np.real(f)

```

Listing 5: Calculating Eigenface feature for training images

```

1 def getTestF(r):
2     transpose = train_V_T[:r, :].T
3     f = np.dot(test_data, transpose)
4     return np.real(f)

```

Listing 6: Calculating Eigenface feature for test images

Part (g) Face recognition

The code in Listing 7 is implemented for the logistic regression for face recognition, whose classification accuracy on the test set is **0.8**. The classification accuracy on the test set when $r = 1, 2, \dots, 200$ is illustrated in Figure 5.

```

1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import accuracy_score
3
4 # extract training & test features
5 def extractFeatures(r):
6     train_F = getTrainF(r)
7     test_F = getTestF(r)
8     return train_F, test_F
9
10 def logRegOvr(train_F, test_F):
11     # train logistic regression model
12     logRegOvr = LogisticRegression(
13         multi_class="ovr", C=1e5,
14         fit_intercept=True, solver='lbfgs', max_iter=1000
15     )
16     logRegOvr.fit(train_F, train_labels)
17     # predict on test features
18     pred_labels = logRegOvr.predict(test_F)
19     # report accuracy
20     acc = accuracy_score(test_labels, pred_labels)
21     return ACC
22
23 train_F, test_F = extractFeatures(10)
24 acc = logRegOvr(train_F, test_F)
25 print("When r = 10, the classification accuracy on test set is", ACC)
26
27 # When r = 10, the classification accuracy on test set is 0.8
28
29 all_acc = []
30 for r in range(1, 201):
31     train_F, test_F = extractFeatures(r)
32     all_acc.append(logRegOvr(train_F, test_F))
33
34 plt.figure(figsize=(16,8))
35 plt.plot(range(1,201), all_acc, marker = 'o')
36 plt.grid(True)
37 plt.xlabel("Number of eigenfaces (r)")
38 plt.ylabel("Classification accuracy")
39 plt.show()

```

Listing 7: Logistic regression for face recognition

Part (h) Low-Rank data loss

The code in Listing 8 is to compute and plot the Frobenius distances over $r = 1, 2, \dots, 200$, which is illustrated in Figure 6.

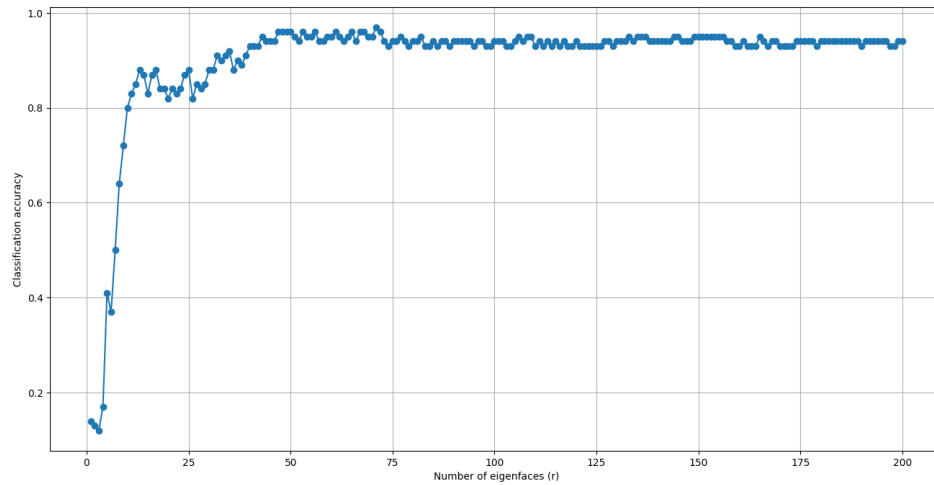


Figure 5: The classification accuracy on the test set when $r = 1, 2, \dots, 200$

```

1 def compute_X_prime(X, V_T, r):
2     # compute F given V_T & r:
3     transpose = V_T[:r, :].T
4     F = np.dot(X, transpose)
5     # compute X_prime:
6     X_prime = np.dot(F, V_T[:r, :])
7     return X_prime
8
9 def compute_Frobenius(X, X_prime):
10     diff = X - X_prime
11     distance = np.sqrt(np.trace(np.dot(diff.T, diff)))
12     return distance
13
14 all_dist = []
15 V_T = compute_V_T(train_data)
16
17 for r in range(1, 201):
18     X_prime = compute_X_prime(train_data, V_T, r)
19     dist = compute_Frobenius(train_data, X_prime)
20     all_dist.append(dist)
21
22 plt.figure(figsize=(16,8))
23 plt.plot(range(1,201), all_dist, marker = 'o')
24 plt.grid(True)
25 plt.xlabel("Number of eigenfaces (r)")
26 plt.ylabel("Average Frobenius distance")
27 plt.show()

```

Listing 8: Compute and plot Frobenius distances over $r = 1, 2, \dots, 200$

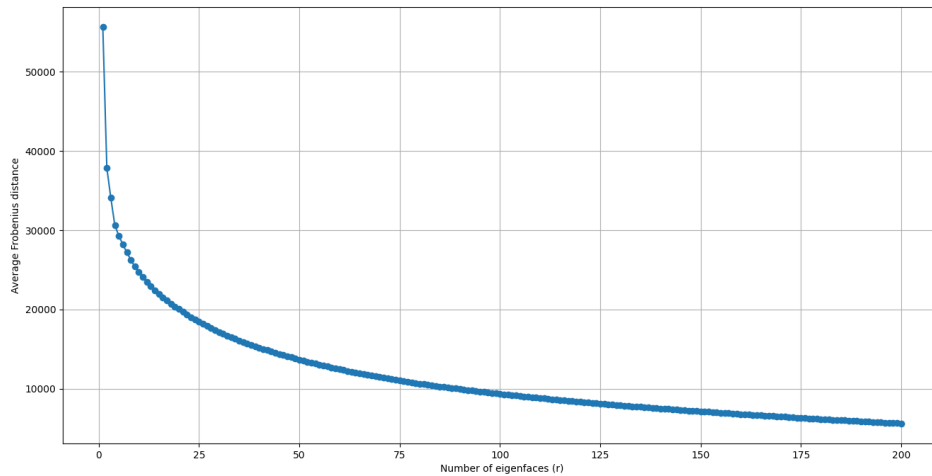


Figure 6: Frobenius distances over $r = 1, 2, \dots, 200$

1.2 Implement EM algorithm

Part (a) Parse and plot all data points on 2-D plane.

Listing 9 shows the script of parsing and plotting all data points on 2-D plane, which is shown in Figure 7.

```

1  with open('OldFaithfulGeyserData.txt') as f:
2      data_array = [line.rstrip('\n') for line in f]
3
4  # remove all text explanations and read data points only
5  idx = data_array.index('    eruptions waiting')
6  data_array = data_array[idx+1:]
7
8  # parse string items to numeric items
9  data_lis = []
10 for str in data_array:
11     entry_lis = list(map(float, str.split()[1:]))
12     data_lis.append(entry_lis)
13
14 # convert python list to numpy array
15 data = np.array(data_lis)
16
17 # plot data points on 2d plane
18 plt.scatter(data[:, 0], data[:, 1], marker='o')
19 plt.xlabel('Eruption time in mins')
20 plt.ylabel('Waiting time to next eruption')
21 plt.grid(True)
22 plt.show()

```

Listing 9: Plot all data points on 2-d plane

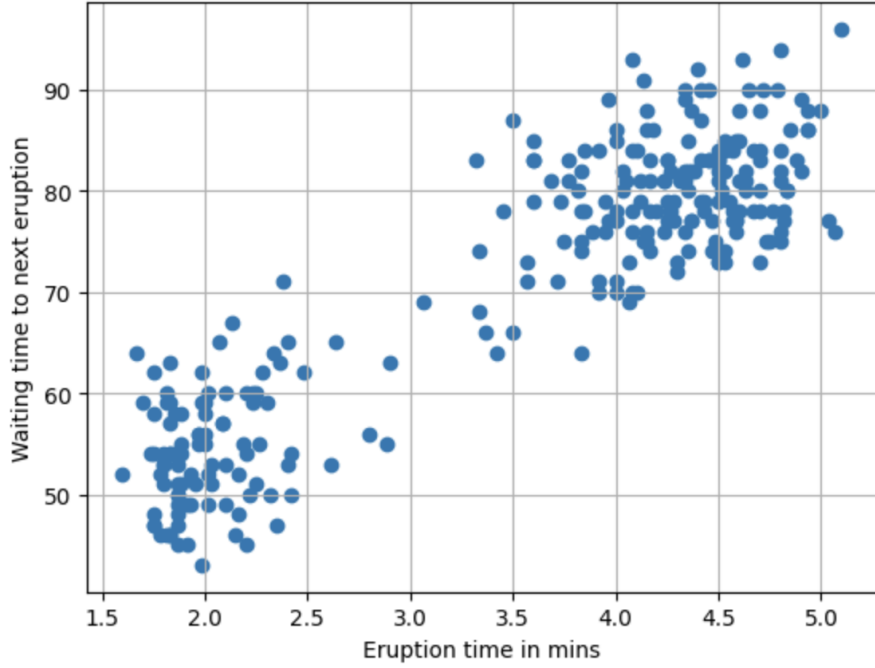


Figure 7: Plotted data points

Part (b) Write down the expression for $P_{\theta_t}(z = k|x)$

$$P_{\theta}(z = k|x) = \frac{P_{\theta}(z = k, x)}{P_{\theta}(x)} = \frac{P(x|z = k)P(z = k)}{\sum_{l=1}^K P_{\theta}(x|z = l)P(z = l)} = \frac{N(x; \mu_k, \Sigma_k) \cdot \phi_k}{\sum_{l=1}^K N(x; \mu_l, \Sigma_l) \cdot \phi_l}$$

Part (c) Write down the formula for μ_k, Σ_k, ϕ

$$\begin{aligned} \text{Goal : } & \max_{\theta} \left(\sum_{k=1}^K \sum_{x \in D} P_{\theta_t}(z_k|x) \log P_{\theta}(x|z_k) + \sum_{k=1}^K \sum_{x \in D} P_{\theta_t}(z_k|x) \log P_{\theta}(z_k) \right) \\ & = \max_{\theta} \left(\sum_{k=1}^K \sum_{x \in D} P_{\theta_t}(z_k|x) (\log P_{\theta}(x|z_k) + \log P_{\theta}(z_k)) \right) \end{aligned}$$

Denoting $P_{\theta_t}(z = k|x)$ as w_k :

$$\text{Goal} \Leftrightarrow \max_{\theta} \left(\sum_{k=1}^K \sum_{x \in D} w_k (\log P_{\theta}(x|z_k) + \log P_{\theta}(z_k)) \right)$$

Taking the derivative of the objective w.r.t μ_k and setting it to zero, we get:

$$\mu_k^{t+1} = \frac{\sum_{x \in D} w_k x}{\sum_{x \in D} w_k}$$

Taking the derivative of the objective w.r.t Σ_k and setting it to zero, we get:

$$\Sigma_k^{t+1} = \frac{\sum_{x \in D} w_k (x - \mu_k^{t+1})(x - \mu_k^{t+1})^T}{\sum_{x \in D} w_k}$$

Taking the derivative of the objective w.r.t ϕ_k and setting it to zero, we get:

$$\frac{\phi_k^{t+1}}{\sum_l \phi_l} = \frac{\sum_{x \in D} w_k}{n}$$

Considering the constraint that $\sum_l \phi_l = 1$, we get:

$$\phi_k^{t+1} = \frac{\sum_{x \in D} w_k}{n}$$

Part (d) Implement and run the EM algorithm

i. Implementing EM algorithm from Scratch

Listing 10 and 11 show the manual implementation of EM algorithm in Python. By running the function `implement_EM`, we can get the learned parameters for the given dataset.

```

1 from scipy.stats import multivariate_normal
2
3 def initialize_parameters(data, num_clusters):
4     num_data = data.shape[0]
5     mu = data[np.random.choice(num_data, num_clusters)]
6     sigma = num_clusters * [np.cov(data, rowvar=False)]
7     phi = num_clusters * [1 / num_clusters]
8     return mu, sigma, phi
9
10 def e_step(data, num_clusters, mu, sigma, phi):
11     # initialize parameters
12     num_data = data.shape[0]
13     posterior = np.zeros((num_data, num_clusters))
14     # compute parameters
15     for k in range(num_clusters):
16         posterior[:, k] = phi[k] * multivariate_normal.pdf(
17             data, mean=mu[k], cov=sigma[k]
18         )
19     return posterior / posterior.sum(axis=1, keepdims=True)
20
21 def m_step(data, num_clusters, posterior):
22     # initialize parameters
23     num_data = data.shape[0]
24     num_features = data.shape[1]
25     mu = np.zeros((num_clusters, num_features))
26     sigma = np.zeros((num_clusters, num_features, num_features))
27     phi = np.zeros(num_clusters)
28     # compute parameters
29     for k in range(num_clusters):
30         phi_k = posterior[:, k].sum()
31         mu[k] = (posterior[:, k].reshape(-1, 1) * data).sum(axis=0) / phi_k
32         numerator = np.dot(
33             (posterior[:, k].reshape(-1, 1) * (data - mu[k])).T, (data - mu[k])
34         )
35         sigma[k] = numerator / phi_k
36         phi[k] = phi_k / num_data
37     return mu, sigma, phi
38
39 def objective(data, num_clusters, mu, sigma, phi):
40     likelihood = np.zeros((data.shape[0], num_clusters))
41     for k in range(num_clusters):
42         likelihood[:, k] = phi[k] * multivariate_normal.pdf(
43             data, mean=mu[k], cov=sigma[k]
44         )
45     return np.sum(np.log(likelihood.sum(axis=1)))

```

Listing 10: The E-Step and M-Step Implementation of the EM algorithm

```

1 def implement_EM(data, num_clusters=2, num_iterations=100, threshold=1e-6):
2     mu, sigma, phi = initialize_parameters(data, num_clusters)
3     log_likelihood_prev = 0
4     for k in range(num_iterations):
5         # e-step
6         posterior = e_step(data, num_clusters, mu, sigma, phi)
7         # m-step
8         mu, sigma, phi = m_step(data, num_clusters, posterior)
9         # compute log-likelihood
10        log_likelihood_current = objective(data, num_clusters, mu, sigma, phi)
11        # check for convergence
12        if abs(log_likelihood_current - log_likelihood_prev) < threshold:
13            print(f"Converged at iteration {k}")
14            break
15        # update log-likelihood
16        log_likelihood_prev = log_likelihood_current
17    return mu, sigma, phi
18
19 mu, sigma, phi = implement_EM(data, num_clusters=2)
20 print(f"mu: {mu}\n")
21 print(f"sigma: {sigma}\n")
22 print(f"phi: {phi}\n")
23
24 # Converged at iteration 9
25 # mu: [[ 4.28966404 79.96814021] [ 2.03639079 54.47853991]]
26 # sigma: [[[ 0.16996581 0.94057589] [ 0.94057589 36.04583496]]
27 # [[ 0.06916953 0.43518701] [ 0.43518701 33.69741426]]]
28 # phi: [0.64412618 0.35587382]

```

Listing 11: Running the implemented EM algorithm

ii. Explanation of the termination criterion

The termination criterion is based on the change in the log-likelihood between consecutive iterations of the EM algorithm. The reasoning is as follows:

1. **Monotonicity of the Log-likelihood:** The log-likelihood should never decrease in the EM algorithm. By monitoring the change in log-likelihood, we can assess whether the algorithm is still making meaningful updates to the parameters.
2. **Early Stopping:** If the improvement in the log-likelihood is smaller than the threshold ($1e-6$ in this case), the algorithm terminates, making the algorithm more efficient by avoiding unnecessary iterations.

iii. Trajectories of the two mean vectors

Before plotting the trajectories of the two mean vectors (μ_1 and μ_2), the modified `implement_EM` in Listing 12 is adopted to store the means during the training process.

```

1  # Modify the implement_EM function to store the posterior probabilities
2  def implement_EM(data, num_clusters=2, num_iterations=100, threshold=1e-6):
3      mu, sigma, phi = initialize_parameters(data, num_clusters)
4      log_likelihood_prev = -np.inf
5      mu_trajectory = {k: [] for k in range(num_clusters)}
6      post_prob = [] # Store posterior probabilities
7
8      for k in range(num_iterations):
9          # Store the means
10         for i in range(num_clusters):
11             mu_trajectory[i].append(mu[i].copy())
12         # e-step
13         posterior = e_step(data, num_clusters, mu, sigma, phi)
14         post_prob.append(posterior.copy()) # Store posterior
15         # m-step
16         mu, sigma, phi = m_step(data, num_clusters, posterior)
17         # compute log-likelihood
18         log_likelihood_current = objective(data, num_clusters, mu, sigma, phi)
19         # check for convergence
20         if abs(log_likelihood_current - log_likelihood_prev) < threshold:
21             break
22         log_likelihood_prev = log_likelihood_current
23
24     return mu, sigma, phi, mu_trajectory, post_prob # Return posterior
25
26 # Run the EM algorithm and get the trajectories and posteriors
27 mu, sigma, phi, mu_trajectory, post_prob = implement_EM(data, num_clusters=2)

```

Listing 12: Modified EM implementation for storing the means

To plot the trajectories of mean vectors, the code in Listing 13 is implemented. Additionally, we plotted the clustered data points in the figure as well to illustrate the final results of the EM algorithm, as shown in Figure 8.

```

1  # Assign colors to data points based on highest posterior probability
2  colors = ["red", "blue"]
3
4  # Create a scatter plot for clustered data points
5  last_posterior = post_prob[-1] # Using posteriors from the last iteration
6  cluster_assignments = np.argmax(last_posterior, axis=1)
7
8  plt.figure(figsize=(10, 6))
9  for i in range(2): # Assuming two clusters
10     cluster_data = data[cluster_assignments == i]
11     plt.scatter(cluster_data[:, 0], cluster_data[:, 1], color=colors[i],
12                alpha=0.2, label=f"Cluster {i+1} data")
13     # Plotting the trajectories of the means
14     mu_x, mu_y = zip(*mu_trajectory[i])
15     plt.plot(mu_x, mu_y, color=colors[i], marker="o",
16             linestyle="--", linewidth=2, markersize=5,
17             label=f"Mean trajectory of cluster {i+1}")
18
19 plt.xlabel("X-coordinate")
20 plt.ylabel("Y-coordinate")
21 plt.title("Clustered Data Points and Trajectories of Mean Vectors")
22 plt.legend()
23 plt.grid()
24 plt.show()

```

Listing 13: Plotting the clustered data points and trajectories of mean vectors

Part (e) Experiment with K-means clustering and comment

Given the same dataset, we implemented a K-means clustering to examine its difference with the GMM model fit using the EM algorithm, as implemented in List 14. The clustering result is visualized in Figure 9. Upon analyzing the clustering results, there are discernible differences between the clusters formed by K-means and the EM algorithm (GMM). As shown in Figure 9, K-means seems to prioritize the data's vertical (Y-axis) spread, resulting in a horizontal boundary between the clusters. In contrast, as Figure 8 shows, the GMM seems to capture more nuanced clusters with elliptical shapes, which is evident from the trajectories of the mean vectors. The K-means algorithm partitions the data in a more definite manner by assigning each to the closest cluster following Euclidean distance, while EM-GMM provides a probabilistic approach, which can result in softer boundaries and more overlap between clusters.

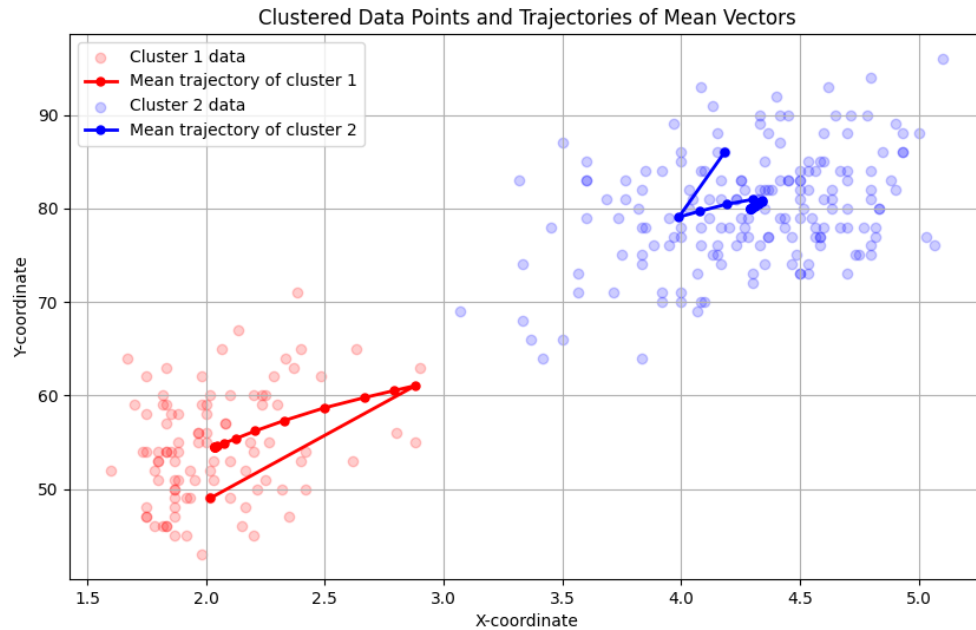


Figure 8: Clustered data points and trajectories of mean vectors

```

1 from sklearn.cluster import KMeans
2
3 # k-means clustering
4 kmeans = KMeans(n_clusters=2, random_state=42)
5 labels = kmeans.fit_predict(data)
6
7 # Plot
8 plt.figure(figsize=(10, 6))
9 plt.scatter(data[:, 0], data[:, 1], c=labels, cmap="viridis", marker="o")
10 plt.scatter(
11     kmeans.cluster_centers[:, 0], kmeans.cluster_centers[:, 1],
12     s=200, c="red", marker="o",
13 )
14 plt.title("K-means Clustering")
15 plt.xlabel("Eruption time in mins")
16 plt.ylabel("Waiting time to next eruption")
17 plt.grid(True)
18 plt.savefig("img/kmeans.png")
19 plt.show()

```

Listing 14: K-means clustering on the same dataset

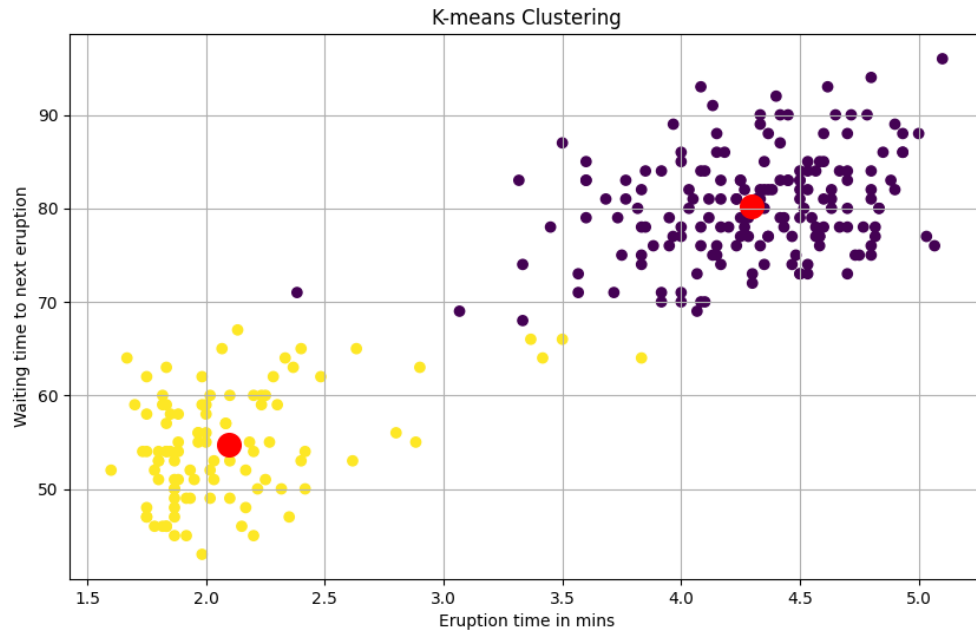


Figure 9: Clustered data points with K-means

1.3 Kernel Density Estimation

Part (a) Plot the density estimate

We are given the following data:

26, 30, 27, 18, 75, 66, 73, 63, 56, 83

Before implementing and plotting the kernel density estimate (KDE) of the given dataset, we need to first use the code in Listing 15 to prepare our data and the range of x values for plotting.

```

1 # Given data
2 data = np.array([26, 30, 27, 18, 75, 66, 73, 63, 56, 83]).reshape(-1, 1)
3
4 # Values of x for the density estimate
5 x_vals = np.linspace(min(data)-20, max(data)+20, 1000).reshape(-1, 1)

```

Listing 15: Prepare the data, bandwidths, and x value for plotting

Listing 16 shows the code to plot the KDE of the given dataset using a Gaussian kernel with Sklearn.


```

1 from sklearn.neighbors import KernelDensity
2
3 # Kernel density estimation for  $\delta = 100$ 
4 kde_100 = KernelDensity(kernel='gaussian', bandwidth=100).fit(data)
5 log_density_100 = kde_100.score_samples(x_vals)
6
7 # Kernel density estimation for  $\delta = 10$ 
8 kde_10 = KernelDensity(kernel='gaussian', bandwidth=10).fit(data)
9 log_density_10 = kde_10.score_samples(x_vals)
10
11 # Plotting the density estimates
12 plt.figure(figsize=(10, 6))
13 plt.plot(x_vals, np.exp(log_density_100), label='$\delta=100$')
14 plt.plot(x_vals, np.exp(log_density_10), label='$\delta=10$')
15 plt.scatter(data, np.zeros_like(data), color='red', label='Data points')
16 plt.title('Kernel Density Estimate with Gaussian Kernel (Using Sklearn)')
17 plt.legend()
18 plt.xlabel('x')
19 plt.ylabel('Density')
20 plt.grid(True)
21 plt.savefig(IMG_PATH + 'sklearn_kde.png')
22 plt.show()

```

Listing 16: Plot the kernel density estimate (KDE) using a Gaussian kernel with Sklearn

Moreover, we manually implement the KDE using NumPy and plot the density estimates. The formula for the KDE with a Gaussian kernel is:

$$\tilde{p}(x) = \frac{1}{n\delta} \sum_{i=1}^n K\left(\frac{x - x_i}{\delta}\right)$$

where

$$K(u) = \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}}$$

Listing 16 shows the code to plot the KDE of the given dataset using a Gaussian kernel, manually calculated with bandwidth parameters $\delta = 100$ and $\delta = 10$.

Figure 10 illustrates the kernel density estimate of the given dataset using a Gaussian kernel, as computed by the `sklearn` library and manual implementation with bandwidth parameters $\delta = 100$ and $\delta = 10$. The results are consistent between these two implementations.

```

1  # Function to calculate Gaussian kernel density estimate
2  def gaussian_kde(x, data, bandwidth):
3      n = len(data)
4      constant = 1 / (np.sqrt(2 * np.pi) * bandwidth)
5      densities = np.sum(
6          [constant * np.exp(-((x - xi) ** 2) / (2 * bandwidth**2)) for xi in data],
7          axis=0,
8      )
9      return densities / n
10
11
12  # Calculating densities manually
13  density_100_manual = gaussian_kde(x_vals, data, bandwidth=100)
14  density_10_manual = gaussian_kde(x_vals, data, bandwidth=10)
15
16  # Plotting the manually calculated densities
17  plt.figure(figsize=(10, 6))
18  plt.plot(x_vals, density_100_manual, label="$\delta=100$ (Manual)")
19  plt.plot(x_vals, density_10_manual, label="$\delta=10$ (Manual)")
20  plt.scatter(data, np.zeros_like(data), color="red", label="Data points")
21  plt.title("Kernel Density Estimate with Gaussian Kernel (Manual Implementation)")
22  plt.xlabel("x")
23  plt.ylabel("Density")
24  plt.legend()
25  plt.grid(True)
26  plt.savefig(IMG_PATH + "manual_kde.png")
27  plt.show()

```

Listing 17: Plot the kernel density estimate (KDE) using a Gaussian kernel with manual implementation

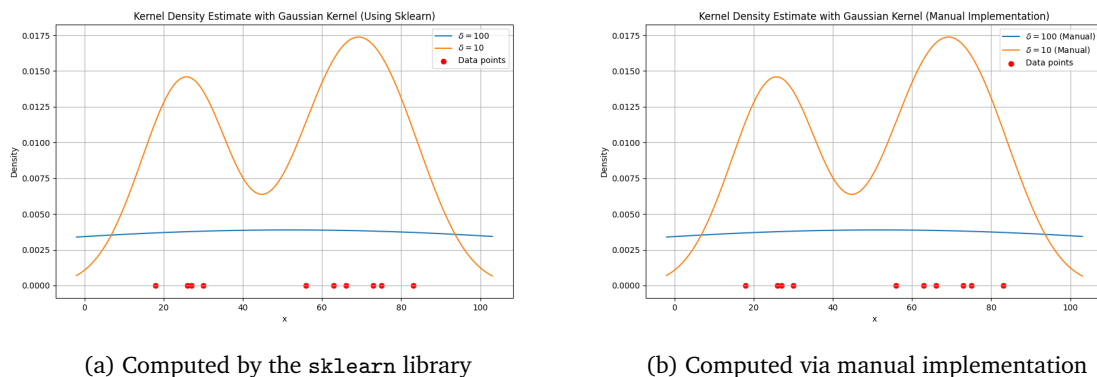


Figure 10: KDE with bandwidth parameters $\delta = 100$ and $\delta = 10$

Part (b) Compare the two bandwidths

The choice of $\delta = 10$ is a preferable choice for this dataset, as justified in the following reasons:

- Bandwidth $\delta = 10$: The curve is more wiggly and appears to fit closely to the data. It captures the individual peaks in the data distribution, which may indicate local variations in the data.
- Bandwidth $\delta = 100$: This curve is smoother and provides a broad overview of the data. It essentially averages out the local variations and might be oversimplifying the data distribution.

Method for choosing δ :

- **Cross-Validation:** We can partition the data into training and validation sets, compute the KDE for the training set using a given bandwidth, and then evaluate how well this KDE predicts the validation set. The bandwidth that performs best (typically measured by log-likelihood) on the validation set is chosen.
- **Visual Inspection:** As demonstrated here, a visual inspection of the KDE plot for different bandwidths can give a good idea of which bandwidth provides a better representation of the data.

Part (c) Probability of new points

The code in Listing 18 is used to compute the density $\tilde{p}(x)$ at the new samples $x = 30$ and $x = 95$ for both $\delta = 100$ and $\delta = 10$.

```
1 kde_100 = KernelDensity(kernel='gaussian', bandwidth=100).fit(data)
2 kde_10 = KernelDensity(kernel='gaussian', bandwidth=10).fit(data)
3
4 # New samples
5 new_samples = np.array([30, 95]).reshape(-1, 1)
6
7 # Calculating the density estimates for the new points
8 density_new_100 = kde_100.score_samples(new_samples)
9 density_new_10 = kde_10.score_samples(new_samples)
10
11 # Returning the density values of the new points
12 print(np.exp(density_new_100), np.exp(density_new_10))
13
14 # [0.00380071 0.00355839] [0.01358759 0.00292187]
```

Listing 18: Calculate the probability of two new samples $x = 30$ and $x = 95$

The densities $\tilde{p}(x)$ at the new samples $x = 30$ and $x = 95$ were calculated using `sklearn` implementation. The results are as follows:

- For $\delta = 100$:

$$\tilde{p}(30) \approx 0.003801, \tilde{p}(95) \approx 0.003558$$

- For $\delta = 10$:

$$\tilde{p}(30) \approx 0.013588, \tilde{p}(95) \approx 0.002922$$

Furthermore, we use the code in Listing 19 to visualize these new points on the previous plots, as shown in Figure 11. The KDE are displayed with the new samples $x = 30$ and $x = 95$ highlighted as green dots for bandwidth parameters $\delta = 100$ and $\delta = 10$.

```

1  # Plotting the density estimates with new points
2  plt.figure(figsize=(10, 6))
3  plt.plot(x_vals, np.exp(log_density_100), label="$\delta=100$")
4  plt.plot(x_vals, np.exp(log_density_10), label="$\delta=10$")
5  plt.scatter(data, np.zeros_like(data), color="red", label="Data points")
6  plt.scatter(
7      new_samples, np.zeros_like(new_samples), color="blue",
8      label="New points"
9  )
10
11 # Highlighting new samples
12 plt.scatter(
13     new_samples, np.exp(density_new_100), color="green",
14     label="New points (Density)"
15 )
16 plt.scatter(
17     new_samples, np.exp(density_new_10), color="green",
18     label="New points (Density)"
19 )
20
21 plt.title("Kernel Density Estimate with New Points")
22 plt.xlabel("x")
23 plt.ylabel("Density")
24 plt.legend()
25 plt.grid(True)
26 plt.savefig(IMG_PATH + "new_points.png")
27 plt.show()

```

Listing 19: Visualize the new samples in the previous KDE plot

Part (d) Rule to accept or reject points

A possible rule to determine whether a point is an outlier based on KDE is by setting a threshold density value, τ . A point x is considered part of the data distribution if its density $\tilde{p}(x)$ is above the threshold:

$$\tilde{p}(x) > \tau$$

Otherwise, it is considered an outlier.

However, this strategy might not always work effectively, especially when the data distribution is multi-

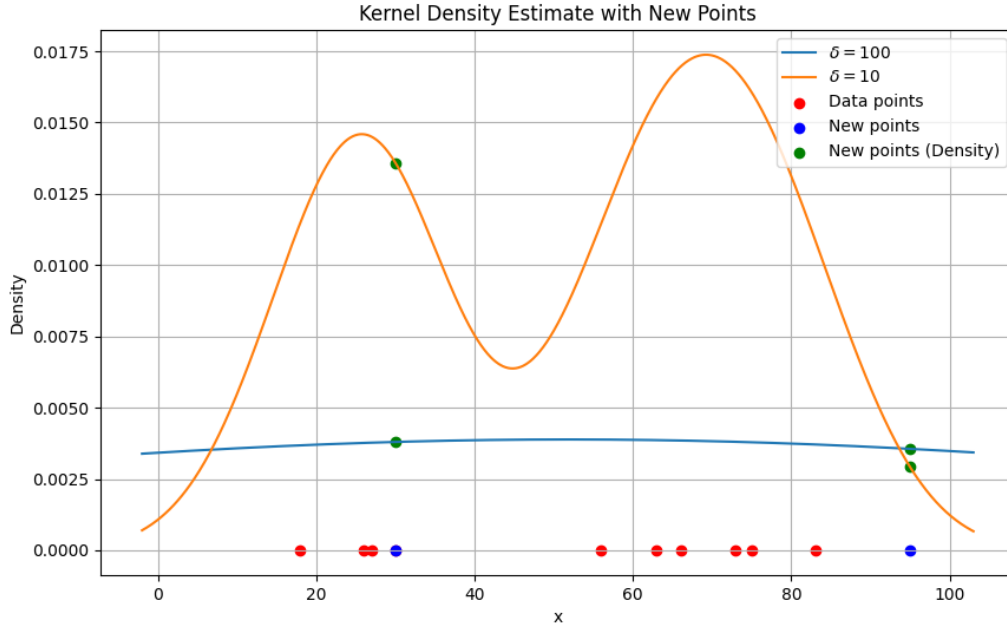


Figure 11: KDE with new samples highlighted

modal or has varying densities in different regions. In such cases, a global threshold might not suitably differentiate between typical points and outliers. Additionally, the choice of kernel and bandwidth also influences the KDE, making the rule sensitive to these hyperparameters.

Part (e) Compare Gaussian and Tophat kernels

The code in Listing 20 change the kernel to a Tophat kernel and plot the KDE for $\delta = 10$. The Tophat kernel is defined as:

$$K(x, z; \delta) = \begin{cases} 1 & \text{if } \|x - z\| \leq \frac{\delta}{2} \\ 0 & \text{otherwise} \end{cases}$$

Figure 12 visualizes the kernel density estimate (KDE) using the Tophat kernel with a bandwidth parameter $\delta = 10$.

One shortcoming of using the Tophat kernel for outlier detection compared to the Gaussian kernel is its abrupt transition at the boundaries of the kernel window. This can lead to a less smooth and less generalized estimate of the data distribution, making it harder to differentiate between typical points and outliers effectively.

Let's consider a one-dimensional dataset mostly distributed between values 10 and 20 with a single outlier at value 35.

- **Tophat Kernel:** For a Tophat kernel with a bandwidth parameter, $\delta = 20$. the outlier might not be effectively identified as the KDE might be influenced heavily due to the uniform weighting within the kernel window.
- **Gaussian Kernel:** For a Gaussian kernel, even if the kernel window includes the outlier, it would

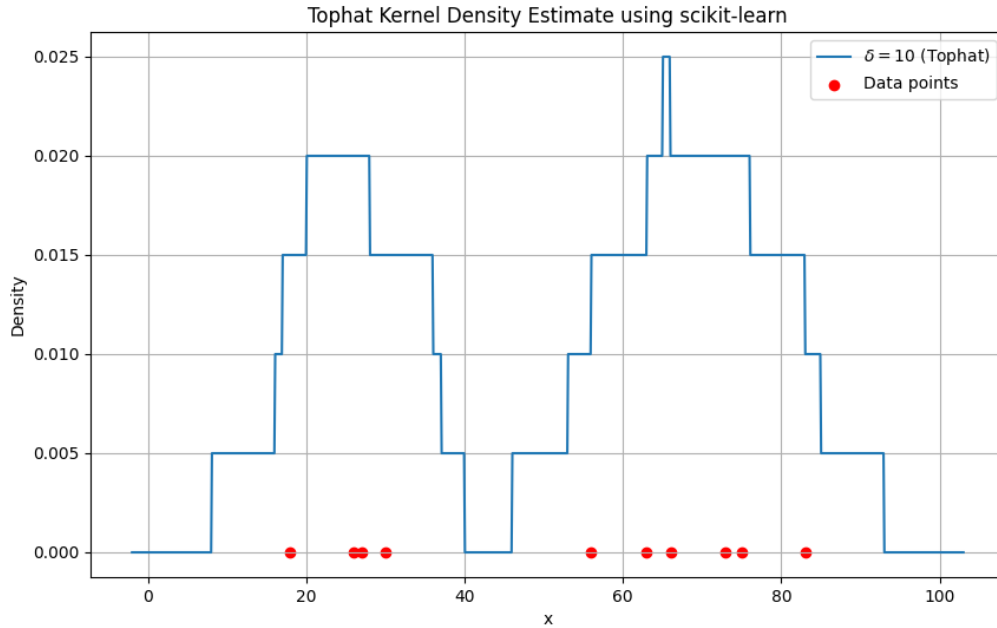


Figure 12: KDE using Tophat kernel with $\delta = 10$

assign a lower weight to it due to the decreasing weight assigned to points as they move away from the center of the kernel. In other words, the Gaussian kernel, with its smooth weight transition, can more effectively identify the outlier by assigning lower weights to distant points.

```

1  # Kernel density estimation with Tophat kernel for  $\delta = 10$ 
2  kde_tophat_10 = KernelDensity(kernel="tophat", bandwidth=10).fit(
3      data.reshape(-1, 1)
4  )
5  log_density_tophat_10 = kde_tophat_10.score_samples(x_vals)
6
7  # Plotting the Tophat kernel density estimates
8  plt.figure(figsize=(10, 6))
9  plt.plot(x_vals, np.exp(log_density_tophat_10), label="$\delta=10$ (Tophat)")
10 plt.scatter(data, np.zeros_like(data), color="red", label="Data points")
11 plt.title("Tophat Kernel Density Estimate using scikit-learn")
12 plt.xlabel("x")
13 plt.ylabel("Density")
14 plt.legend()
15 plt.grid(True)
16 plt.savefig(IMG_PATH + "tophat_kde.png")
17 plt.show()

```

Listing 20: Plot the Tophat density estimate of the dataset using $\delta = 10$

2 Written Exercises

2.1 K-means clustering

Part (a) $J(c_{t-1}, f_t) \leq J(c_{t-1}, f_{t-1})$

Proof. Consider a data point $x^{(i)}$. Let $c_{t-1}^{(a)}$ be the closest centroid to $x^{(i)}$ at time $t-1$, and let $c_{t-1}^{(b)}$ be the centroid assigned to $x^{(i)}$ at time $t-1$. Thus, we have:

$$\|x^{(i)} - c_{t-1}^{(a)}\|_2 \leq \|x^{(i)} - c_{t-1}^{(b)}\|_2 \quad (1)$$

Squaring both sides to work with squared distances, and summing over all data points, we get:

$$\sum_{i=1}^n \|x^{(i)} - c_{t-1}^{(a)}\|_2^2 \leq \sum_{i=1}^n \|x^{(i)} - c_{t-1}^{(b)}\|_2^2 \quad (2)$$

After Step 1 of the K-means iteration, $c_{t-1}^{(a)}$ becomes the centroid for the new assignment of $x^{(i)}$, i.e., $c_{t-1}^{f_t(x^{(i)})}$. Using the notation $c_{t-1}^{f_t(x^{(i)})}$ and $c_{t-1}^{f_{t-1}(x^{(i)})}$ to represent the centroids, we rewrite the equation as:

$$\sum_{i=1}^n \|x^{(i)} - c_{t-1}^{f_t(x^{(i)})}\|_2^2 \leq \sum_{i=1}^n \|x^{(i)} - c_{t-1}^{f_{t-1}(x^{(i)})}\|_2^2 \quad (3)$$

This simplifies to:

$$J(c_{t-1}, f_t) \leq J(c_{t-1}, f_{t-1}) \quad (4)$$

Concluding, the objective function does not increase and therefore, the K-means algorithm is monotonically decreasing. \square

Part (b): $J(c_t, f_t) \leq J(c_{t-1}, f_t)$

Proof. Consider a single cluster, denoted by the index k , and let the points in this cluster be given by $\{x^{(i)} : f_t(x^{(i)}) = k\}$. The objective is to show that the updated centroid, $c_t^{(k)}$, which is the mean of the points in cluster k , minimizes the sum of squared Euclidean distances to the points in the cluster. Let us denote this sum as $M_k(c)$, which is formulated as follows:

$$M_k(c) = \sum_{i: f_t(x^{(i)})=k} \|x^{(i)} - c\|_2^2$$

To find the value of c that minimizes $M_k(c)$, we take the derivative with respect to c and set it to zero:

$$\frac{dM_k(c)}{dc} = -2 \sum_{i: f_t(x^{(i)})=k} (x^{(i)} - c)$$

Solving the equation $\frac{dM_k(c)}{dc} = 0$, we get:

$$\begin{aligned}
\sum_{i: f_t(x^{(i)})=k} (x^{(i)} - c) &= 0 \\
\sum_{i: f_t(x^{(i)})=k} x^{(i)} &= cS^{(k)} \\
c &= \frac{1}{S^{(k)}} \sum_{i: f_t(x^{(i)})=k} x^{(i)} \\
c &= c_t^{(k)}
\end{aligned}$$

This confirms that the value of c which minimizes $M_k(c)$ is indeed the updated centroid $c_t^{(k)}$. Consequently, updating the centroids in this manner also ensures that the objective function $J(c_t, f_t)$ is minimized, since it is the sum of Euclidean distances, establishing that:

$$J(c_t, f_t) \leq J(c_{t-1}, f_t)$$

□

2.2 Local Optima in K-means

Consider a dataset with the following 1D points:

$$[0, 2, 6, 11, 12, 13]$$

We want to partition these points into two clusters ($K = 2$).

Initialization 1: Centroids: $C_1 : 2, C_2 : 12$

$$\begin{aligned}
\text{Iteration 1: } & \begin{cases} \text{Points closer to } C_1 : [0, 2, 6] \\ \text{Points closer to } C_2 : [11, 12, 13] \end{cases} \\
& \text{New centroids: } C_1 : \frac{0+2+6}{3} = 2.67, C_2 : \frac{11+12+13}{3} = 12
\end{aligned}$$

$$\begin{aligned}
\text{Iteration 2: } & \begin{cases} \text{Points closer to } C_1 : [0, 2, 6] \\ \text{Points closer to } C_2 : [11, 12, 13] \end{cases} \\
& \text{New centroids: } C_1 : \frac{0+2+6}{3} = 2.67, C_2 : \frac{11+12+13}{3} = 12
\end{aligned}$$

Stop the iteration here because the centroids do not change anymore.

Initialization 2: Centroids: $C_1 : 0, C_2 : 6$

$$\begin{aligned} \text{Iteration 1: } & \begin{cases} \text{Points closer to } C_1 : [0, 2] \\ \text{Points closer to } C_2 : [6, 11, 12, 13] \end{cases} \\ & \text{New centroids: } C_1 : \frac{0+2}{2} = 1, C_2 : \frac{6+11+12+13}{4} = 10.5 \end{aligned}$$

$$\begin{aligned} \text{Iteration 2: } & \begin{cases} \text{Points closer to } C_1 : [0, 2] \\ \text{Points closer to } C_2 : [6, 11, 12, 13] \end{cases} \\ & \text{New centroids: } C_1 : \frac{0+2}{2} = 1, C_2 : \frac{6+11+12+13}{4} = 10.5 \end{aligned}$$

Stop the iteration here because the centroids do not change anymore.

In **Initialization 1**, the final clusters are $[0, 2, 6]$ and $[11, 12, 13]$. In **Initialization 2**, the final clusters are $[0, 2]$ and $[6, 11, 12, 13]$, showcasing the sensitivity of K-means to initial centroid choices.