

January 13, 2021

1 EXERCISE 3 - ML - Grundverfahren

1.1 1.) Constrained Optimization (6 Points)

You are given the following Optimization problem:

$$\begin{aligned} \min_{\mathbf{x}} \quad & \mathbf{x}^T \mathbf{M} \mathbf{x} + \mathbf{x}^T \mathbf{h} \\ \text{s.t.} \quad & \mathbf{x}^T \mathbf{b} \geq c, \end{aligned}$$

where \mathbf{M} is a positive definit, symmetric Matrix. Note that vectors and matrices are boldsymbol, where Matrices have capital letters. Derive the optimal solution for \mathbf{x} independant of the Lagrangian multiplier(s) (i.e. you have to solve for the dual). Make sure that you mark vectors and matrices as a boldsymbol and small letters and capital letters respectively. Symbols which are not marked as boldsymbols will count as scalar. Take care of vector/matrix multiplication and derivatives. And make use of the properties of \mathbf{M} . Don't forget to look up matrix-vector calculus in the matrix cookbook, if you don't remember the rules.

Hint 1) \mathbf{M} is a positive definit, symmetric Matrix

$$\begin{aligned} \mathbf{M}^T &= \mathbf{M} \\ \mathbf{M}^{-T} &= \mathbf{M}^{-1} \end{aligned}$$

2) $\mathbf{h}^T \mathbf{M} \mathbf{b}$ is scalar

$$\mathbf{h}^T \mathbf{M} \mathbf{b} = \mathbf{b}^T \mathbf{M} \mathbf{h}$$

3) Matrix-vector calculus

$$\begin{aligned} \frac{\partial \mathbf{x}^T \mathbf{M} \mathbf{x}}{\partial \mathbf{x}} &= 2\mathbf{M} \mathbf{x} \\ \frac{\partial \mathbf{x}^T \mathbf{h}}{\partial \mathbf{x}} &= \mathbf{h} \end{aligned}$$

Proof 1) Lagrangian

$$L(\mathbf{x}, \lambda) = \mathbf{x}^T \mathbf{M} \mathbf{x} + \mathbf{x}^T \mathbf{h} - \lambda (\mathbf{x}^T \mathbf{b} - c)$$

2) Find \mathbf{x}^*

$$\mathbf{x}^* = \underset{\mathbf{x}}{\operatorname{argmin}} L(\mathbf{x}, \lambda)$$

$$\frac{\partial L}{\partial \mathbf{x}} = 2\mathbf{M}\mathbf{x} + \mathbf{h} - \lambda\mathbf{b} = 0$$

$$\mathbf{x}^* = \frac{1}{2}\mathbf{M}^{-1}(\lambda\mathbf{b} - \mathbf{h})$$

3) Dual function $g(\lambda)$

$$\begin{aligned} g(\lambda) &= L(\mathbf{x}^*, \lambda) = \mathbf{x}^{*T} \mathbf{M} \mathbf{x}^* + \mathbf{x}^{*T} \mathbf{h} - \lambda \mathbf{x}^{*T} \mathbf{b} + \lambda c \\ &= \frac{1}{4} (\lambda \mathbf{b} - \mathbf{h})^T \mathbf{M}^{-1} (\lambda \mathbf{b} - \mathbf{h}) + \frac{1}{2} (\lambda \mathbf{b} - \mathbf{h})^T \mathbf{M}^{-1} \mathbf{h} - \frac{1}{2} \lambda (\lambda \mathbf{b} - \mathbf{h})^T \mathbf{M}^{-1} \mathbf{b} + \lambda c \\ &= -\frac{\lambda^2}{4} \mathbf{b}^T \mathbf{M}^{-1} \mathbf{b} + \frac{\lambda}{2} \mathbf{b}^T \mathbf{M}^{-1} \mathbf{h} + \lambda c \end{aligned}$$

$$s.t. \lambda \geq 0$$

4) Find λ^*

$$\lambda^* = \underset{\lambda}{\operatorname{argmax}} g(\lambda)$$

$$\frac{\partial g}{\partial \lambda} = -\frac{\lambda}{2} \mathbf{b}^T \mathbf{M}^{-1} \mathbf{b} + \frac{1}{2} \mathbf{b}^T \mathbf{M}^{-1} \mathbf{h} + c = 0$$

$$\lambda^* = \frac{\mathbf{b}^T \mathbf{M}^{-1} \mathbf{h} + 2c}{\mathbf{b}^T \mathbf{M}^{-1} \mathbf{b}}$$

5) optimal solution for \mathbf{x}

- if $\lambda^* < 0$, then $\lambda^* = 0$

$$\begin{aligned} \mathbf{x}^* &= \frac{1}{2} \mathbf{M}^{-1} (\lambda \mathbf{b} - \mathbf{h}) \\ &= -\frac{1}{2} \mathbf{M}^{-1} \mathbf{h} \end{aligned}$$

- if $\lambda^* \geq 0$, then $\lambda^* = \frac{\mathbf{b}^T \mathbf{M}^{-1} \mathbf{h} + 2c}{\mathbf{b}^T \mathbf{M}^{-1} \mathbf{b}}$

$$\begin{aligned} \mathbf{x}^* &= \frac{1}{2} \mathbf{M}^{-1} (\lambda \mathbf{b} - \mathbf{h}) \\ &= -\frac{1}{2} \mathbf{M}^{-1} \mathbf{h} + \lambda \frac{1}{2} \mathbf{M}^{-1} \mathbf{b} \\ &= -\frac{1}{2} \mathbf{M}^{-1} \mathbf{h} + \frac{1}{2} \frac{\mathbf{b}^T \mathbf{M}^{-1} \mathbf{h} + 2c}{\mathbf{b}^T \mathbf{M}^{-1} \mathbf{b}} \mathbf{M}^{-1} \mathbf{b} \end{aligned}$$

1.2 2.) k-Means (7 Points)

Here we will implement one of the most basic approaches to clustering - the k-Means algorithm. Let us start with some basic imports and implementing functionality to visualize our results.

```
[18]: %matplotlib inline

import numpy as np
import matplotlib.pyplot as plt
from typing import Tuple, Optional

def visualize_2d_clustering(data_points: np.ndarray, assignments_one_hot: np.
    ↳ ndarray, centers: np.ndarray, k: int,
                                centers_history: Optional[np.ndarray] = None, title:
    ↳ Optional[str] = None):
    """Visualizes clusters, centers and path of centers"""
    plt.figure(figsize=(6, 6), dpi=100)
    assignments = np.argmax(assignments_one_hot, axis=1)

    for i in range(k):
        # get next color
        c = next(plt.gca()._get_lines.prop_cycler)['color']
        # get cluster
        cur_assignments = assignments == i
        # plot clusters
        plt.scatter(data_points[cur_assignments, 0],
    ↳ data_points[cur_assignments, 1], c=c,
                        label="Cluster {:02d}".format(i))

        #plot history of centers if it is given
        if centers_history is not None:
            plt.scatter(centers_history[:, i, 0], centers_history[:, i, 1],
    ↳ marker="x", c=c)
            plt.plot(centers_history[:, i, 0], centers_history[:, i, 1], c=c)

        plt.scatter(centers[:, 0], centers[:, 1], label="Centers", color="black",
    ↳ marker="X")

    if title is not None:
        plt.title(title)

    plt.legend()
```

Next we going to implement the actual algorithm. As a quick reminder, K-Means works by iterating the following steps:

Start with k randomly picked centers

- 1.) Assign each point to the closest center

- 2.) Adjust centers by taking the average over all points assigned to it

Implementing them will be your task for this exercise

```
[19]: def assignment_step(data_points: np.ndarray, centers: np.ndarray) -> np.ndarray:
    """
    Assignment Step: Computes assignments to nearest cluster
    :param data_points: Data points to cluster (shape: [N x data_dim])
    :param centers: current cluster centers (shape: [k, data_dim])
    :return Assignments (as one hot) (shape: [N, k])
    """

    #####
    # TODO Implement the assignment step of the k-Means algorithm

    # bookkeeping for distances of each data element to centers
    distances = []

    # calculate euclidian distance to all centroids
    for center in centers:
        eucl_dist = np.linalg.norm(data_points - center, axis=1)
        distances.append(eucl_dist)

    # get index of min center
    min_distances = np.argmin(distances, axis=0)

    # one-hot-encode min-distances and return
    assignments = np.eye(data_points.shape[0], centers.shape[0])[min_distances]
    return assignments

    #####

def adjustment_step(data_points: np.ndarray, assignments_one_hot: np.ndarray)
    -> np.ndarray:
    """
    Adjustment Step: Adjust centers given assignment
    :param data_points: Data points to cluster (shape: [N x data_dim])
    :param assignments_one_hot: assignment to adjust to (one-hot
    representation) (shape: [N, k])
    :return Adjusted Centers (shape: [k, data_dim])
    """

    #####
    # TODO Implement the adjustment step of the k-Means algorithm

    # get centroids from one-hot-encoded matrix
    centroids = np.argmax(assignments_one_hot, axis=1)
```

```

    # find all the indices of data points within one cluster. Cluster is
    →identified by allocated column id
    clusters_by_indices = [np.where(centroids == element)[0].tolist() for
    →element in range(0,assignments_one_hot.shape[1])]

    # calculate cluster mean using all data points within cluster
    adj_center = []
    for cluster in clusters_by_indices:
        cluster_center = data_points[cluster].mean(axis = 0)
        # save cluster to mean to list
        adj_center.append(cluster_center)
    return np.array(adj_center)

#####

```

Now to the final algorithm, as said we initialize the centers with random data points and iterate the assignment and adjustment step

```

[20]: def k_means(data_points: np.ndarray, k: int, max_iter: int = 100, vis_interval:
    →int = 3) -> \
        Tuple[np.ndarray, np.ndarray, np.ndarray, np.ndarray]:
    """
    Simple K Means Implementation
    :param data_points: Data points to cluster (shape: [N x data_dim])
    :param k: number of clusters
    :param max_iter: Maximum number of iterations to run if convergence is not
    →reached
    :param vis_interval: After how many iterations to generate the next plot
    :return: - cluster labels (shape: [N])
             - means of clusters (shape: [k, data_dim])
             - SSD over time (shape: [2 * num_iters])
             - History of means over iterations (shape: [num_iters, k,
    →data_dim])
    """
    # Bookkeeping
    i = 0
    means_history = []
    ssd_history = []
    assignments_one_hot = np.zeros(shape=[data_points.shape[0], k])
    old_assignments = np.ones(shape=[data_points.shape[0], k])

    # Initialize with k random data points
    initial_idx = np.random.choice(len(data_points), k, replace=False)
    centers = data_points[initial_idx]
    means_history.append(centers.copy())

```

```

# Iterate while not converged and max number iterations not reached
while np.any(old_assignments != assignments_one_hot) and i < max_iter:
    old_assignments = assignments_one_hot

    # assignment
    assignments_one_hot = assignment_step(data_points, centers)

    # compute SSD
    diffs = np.sum(np.square(data_points[:, None, :] - centers[None, :, :
↪]), axis=-1)
    ssd_history.append(np.sum(assignments_one_hot * diffs))

    # adjustment
    centers = adjustment_step(data_points, assignments_one_hot)

    # compute SSD
    diffs = np.sum(np.square(data_points[:, None, :] - centers[None, :, :
↪]), axis=-1)
    ssd_history.append(np.sum(assignments_one_hot * diffs))

    # Plotting
    if i % vis_interval == 0:
        visualize_2d_clustering(data_points, assignments_one_hot, centers,
↪k, title="Iteration {:02d}".format(i))

    # Bookkeeping
    means_history.append(centers.copy())
    i += 1

    print("Took", i, "iterations to converge")
    return assignments_one_hot, centers, np.array(ssd_history), np.
↪stack(means_history, 0)

```

Finally we run the dataset and visualize the results. Here we provide 4 random datasets, each containing 500 2 samples and you can play around with the number of clusters, k , as well as the seed of the random number generator. Based on this seed the initial centers, and thus the final outcome, will vary.

```

[21]: np.random.seed(42)

data = np.load("samples_3.npy")
k = 8

cluster_labels, centers, ssd_history, centers_history = k_means(data, k)

# plot final clustering with history of centers over iterations

```

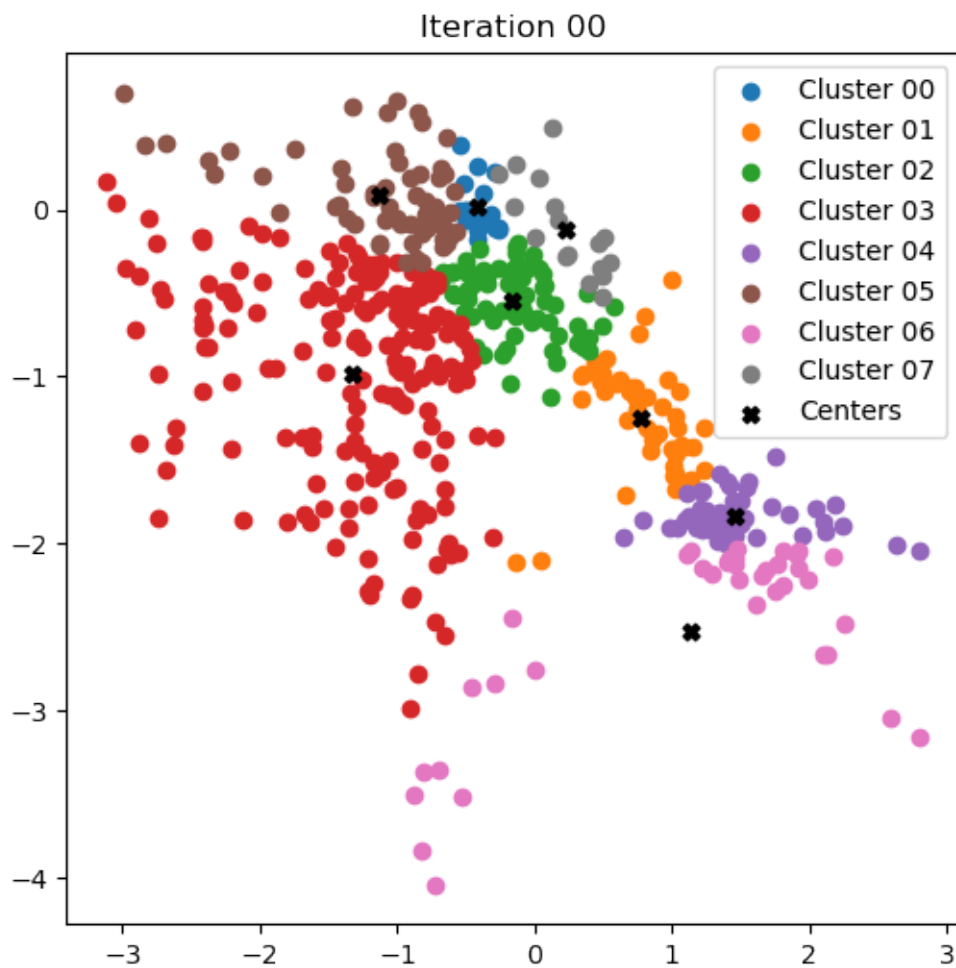
```

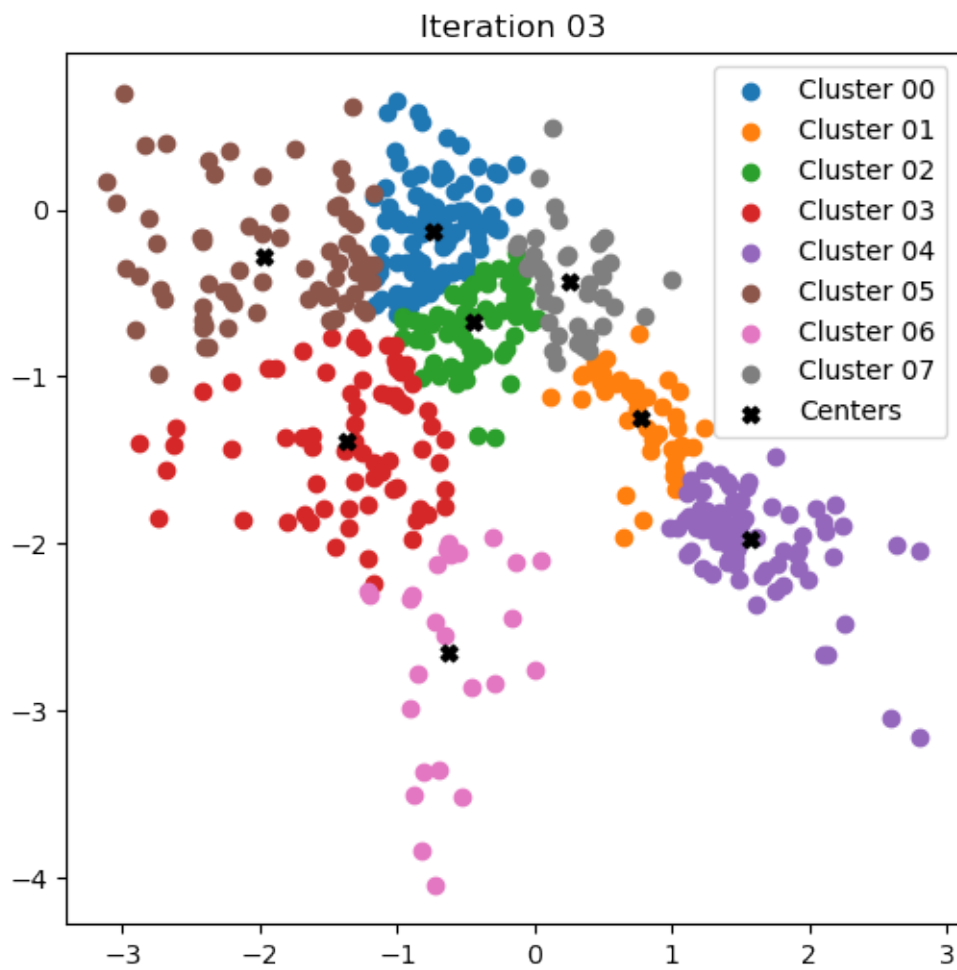
visualize_2d_clustering(data, cluster_labels, centers, k=k,
    ↪centers_history=centers_history, title="Final Clustering")

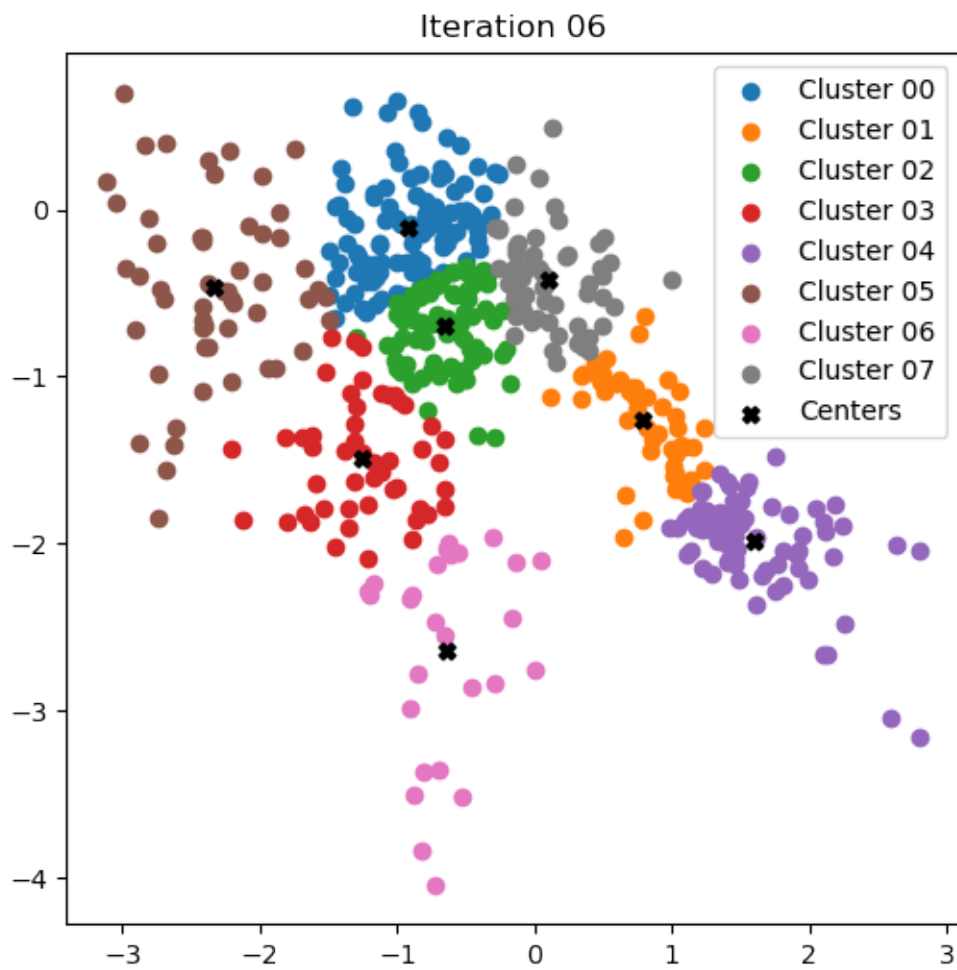
# plot SSD
plt.figure("SSD")
plt.semilogy(np.arange(start=0, stop=len(ssd_history) / 2, step=0.5),
    ↪ssd_history)
plt.xlabel("Iteration")
plt.ylabel("SSD")
plt.show()

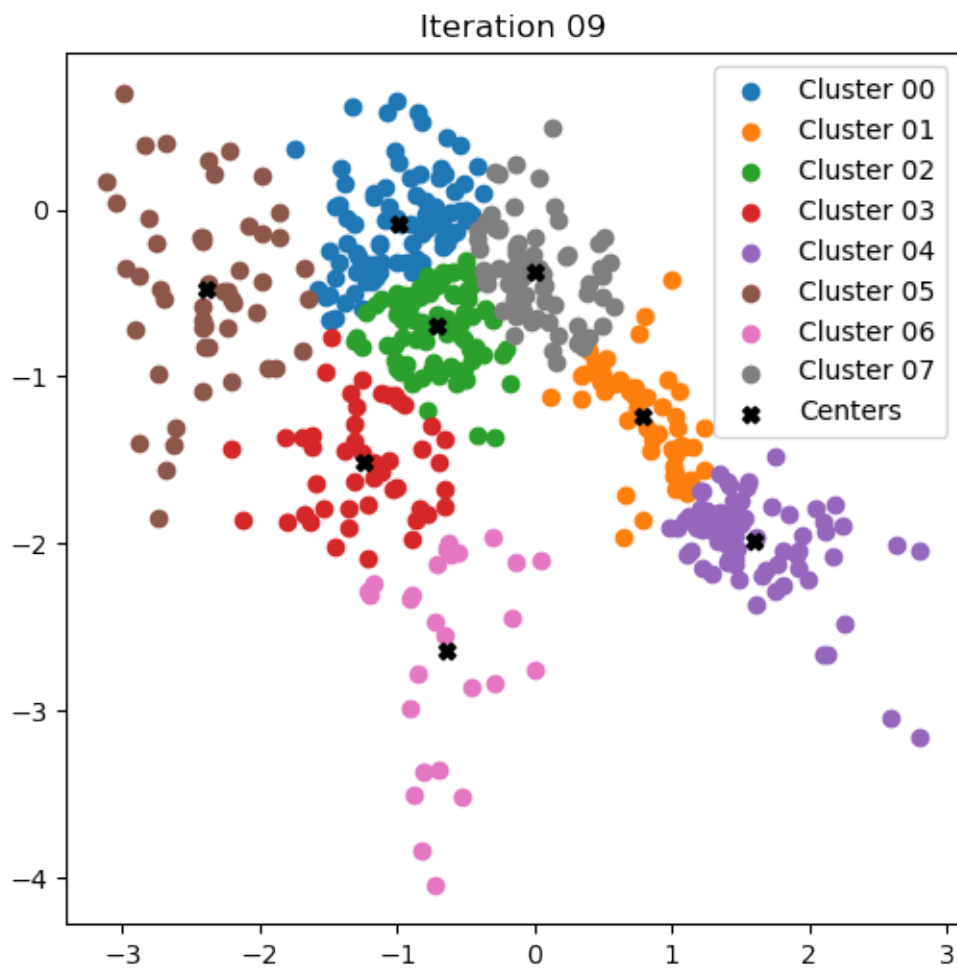
```

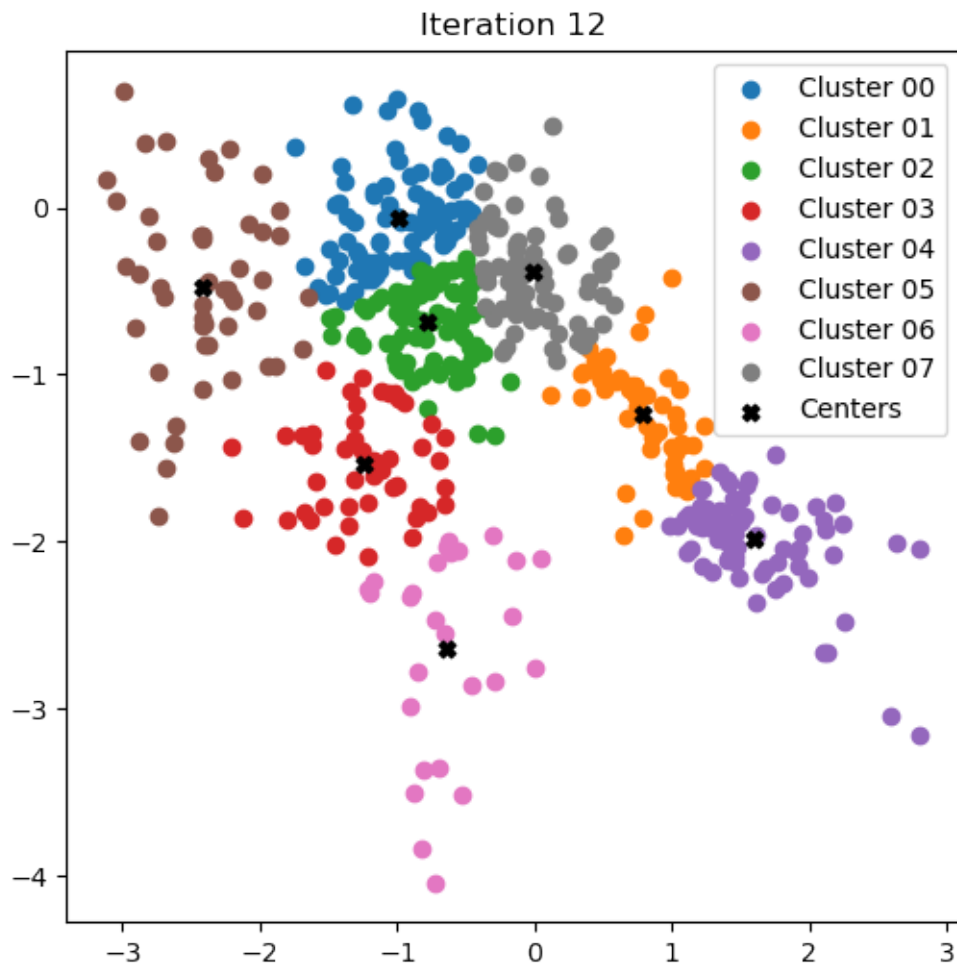
Took 17 iterations to converge

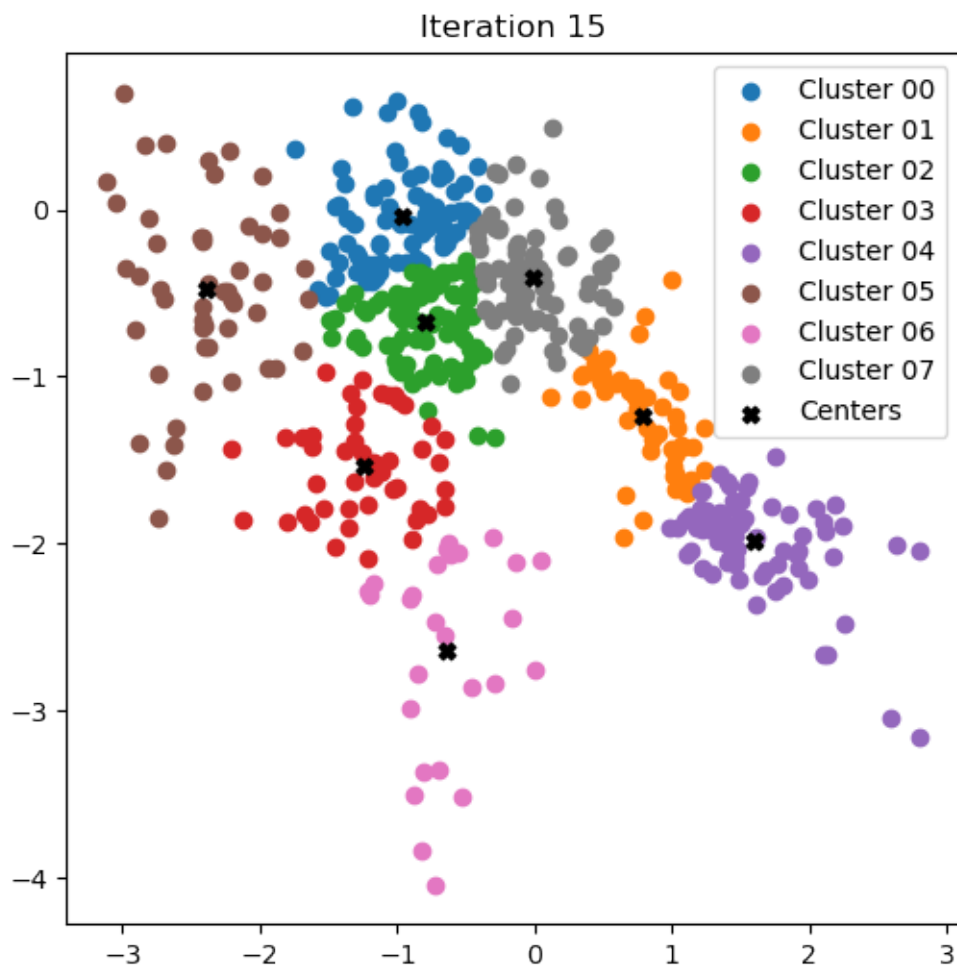


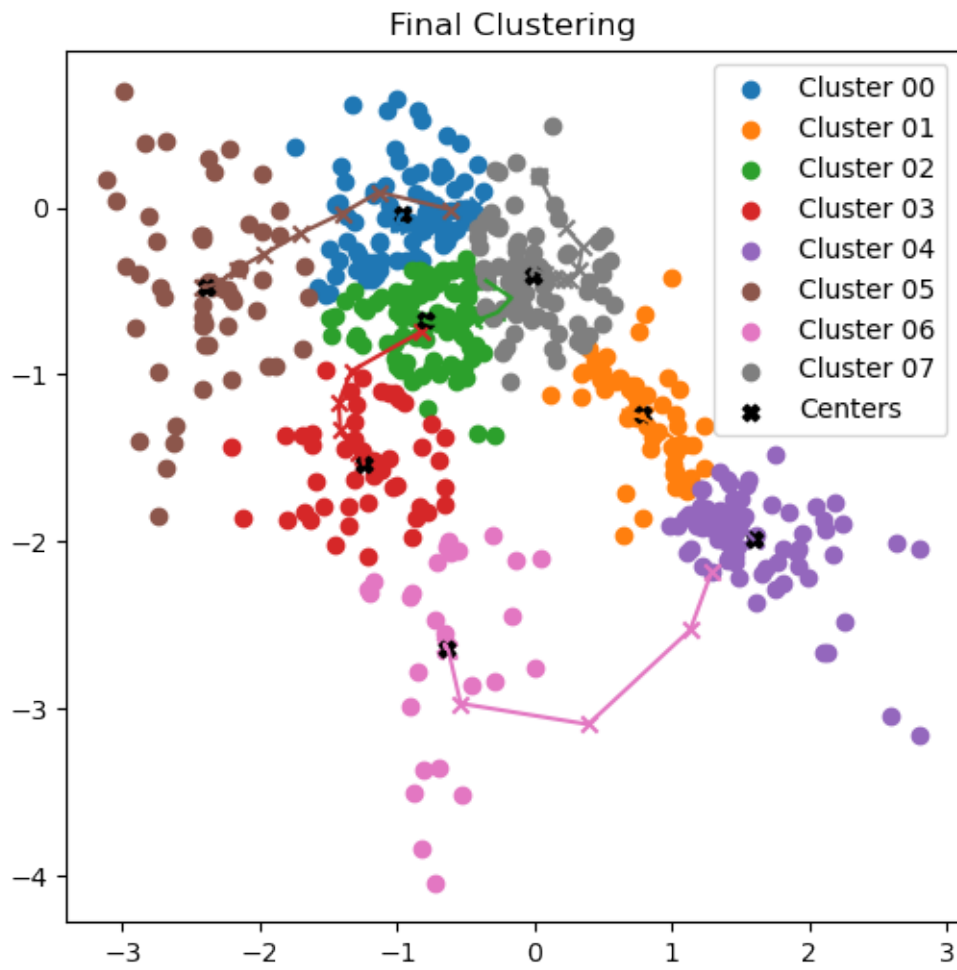


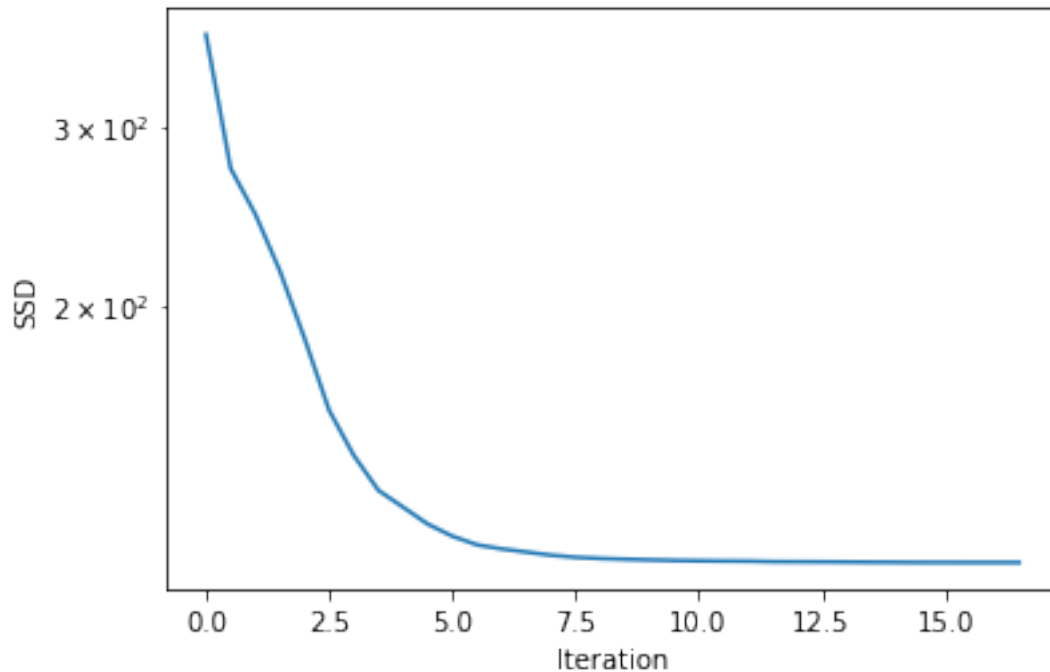












1.3 3.) Expectation Maximization for Gaussian Mixture Models (7 Points)

In the following we implement the Expectation Maximization (EM) Algorithm to fit a Gaussian Mixture Model (GMM) to data. We start with an implementation for the log density of a single Gaussian (take some time to compare this implementation with the one used in the first exercises)...

```
[22]: %matplotlib inline

import numpy as np
import matplotlib.pyplot as plt
from typing import Tuple

def gaussian_log_density(samples: np.ndarray, mean: np.ndarray, covariance: np.
    ↳ ndarray) -> np.ndarray:
    """
    Computes Log Density of samples under a Gaussian Distribution.
    We already saw an implementation of this in the first exercise and noted_
    ↳ there that this was not the "proper"
    way of doing it. Compare the two implementations.
    :param samples: samples to evaluate (shape: [N x dim])
    :param mean: Mean of the distribution (shape: [dim])
    :param covariance: Covariance of the distribution (shape: [dim x dim])
    :return: log N(x|mean, covariance) (shape: [N])
    """
    dim = mean.shape[0]
```

```

chol_covariance = np.linalg.cholesky(covariance)
# Efficient and stable way to compute the log determinant and squared term
→efficiently using the cholesky
logdet = 2 * np.sum(np.log(np.diagonal(chol_covariance) + 1e-25))
# (Actually, you would use scipy.linalg.solve_triangular but I wanted to
→spare you the hustle of setting
# up scipy)
chol_inv = np.linalg.inv(chol_covariance)
exp_term = np.sum(np.square((samples - mean) @ chol_inv.T), axis=-1)
return -0.5 * (dim * np.log(2 * np.pi) + logdet + exp_term)

```

... and some plotting functionality for 2D GMMs:

```

[23]: def visualize_2d_gmm(samples, weights, means, covs, title):
    """Visualizes the model and the samples"""
    plt.figure(figsize=[7,7])
    plt.title(title)
    plt.scatter(samples[:, 0], samples[:, 1], label="Samples", c=next(plt.gca().
→_get_lines.prop_cycler)['color'])

    for i in range(means.shape[0]):
        c = next(plt.gca().get_lines.prop_cycler)['color']

        (largest_eigval, smallest_eigval), eigvec = np.linalg.eig(covs[i])
        phi = -np.arctan2(eigvec[0, 1], eigvec[0, 0])

        plt.scatter(means[i, 0:1], means[i, 1:2], marker="x", c=c)

        a = 2.0 * np.sqrt(largest_eigval)
        b = 2.0 * np.sqrt(smallest_eigval)

        ellipse_x_r = a * np.cos(np.linspace(0, 2 * np.pi, num=200))
        ellipse_y_r = b * np.sin(np.linspace(0, 2 * np.pi, num=200))

        R = np.array([[np.cos(phi), np.sin(phi)], [-np.sin(phi), np.cos(phi)]])
        r_ellipse = np.array([ellipse_x_r, ellipse_y_r]).T @ R
        plt.plot(means[i, 0] + r_ellipse[:, 0], means[i, 1] + r_ellipse[:, 1],
→c=c,
                    label="Component {:02d}, Weight: {:.04f}".format(i,
→weights[i]))
    plt.legend()

```

Now to the actual task: You need to implement 3 functions: - the log likelihood of a GMM for evaluation - the E-Step of the EM algorithm for GMMs - the M-Step of the EM algorithm for GMMs (for this one now for loops are allowed. Using them here will lead to point deduction)

All needed equations are in the slides

```
[24]: # adapted from pset 1
```

```
def mvn_pdf(x: np.ndarray, mu: np.ndarray, sigma: np.ndarray) -> np.ndarray:
    """
    Density of the Multivariate Normal Distribution
    :param x: samples, shape: [N x dimension]
    :param mu: mean, shape: [dimension]
    :param sigma: covariance, shape: [dimension x dimension]
    :return p(x) with  $p(x) = N(\mu, \sigma)$ , shape: [N]
    """
    norm_term = 1 / np.sqrt(np.linalg.det(2 * np.pi * sigma))
    diff = x - np.atleast_2d(mu)
    exp_term = np.sum(np.linalg.solve(sigma, diff.T).T * diff, axis=-1)
    return norm_term * np.exp(-0.5 * exp_term)
```

```
[25]: def gmm_log_likelihood(samples: np.ndarray, weights: np.ndarray, means: np.
    → ndarray, covariances: np.ndarray) -> float:
    """ Computes the Log Likelihood of samples given parameters of a GMM.
    :param samples: samples "x" to compute ess for (shape: [N, dim])
    :param weights: weights (i.e.,  $p(z)$ ) of old model (shape: [num_components])
    :param means: means of old components  $p(x/z)$  (shape: [num_components, dim])
    :param covariances: covariances of old components  $p(x/z)$  (shape:
    → [num_components, dim, dim])
    :return: log likelihood
    """
    #####
    # TODO Implement the log-likelihood for Gaussian Mixtures

    n = samples.shape[0] #number of samples
    num_comp = weights.shape[0] #number of components

    p_x_comp = np.zeros([n,num_comp])

    # Mixture Model:  $p(x)=\sum[p(z)p(x/z)]$ 
    for i in range(num_comp):
        p_x_comp[:,i] = weights[i] * mvn_pdf(samples, means[i],covariances[i])

    # Log of Mixture Model:  $\log[p(x)]$ 
    log_p_x = np.log(np.sum(p_x_comp,axis=-1))

    return np.sum(log_p_x,axis=0)

    #####

def e_step(samples: np.ndarray, weights: np.ndarray, means: np.ndarray,
    → covariances: np.ndarray) -> np.ndarray:
```



```

    """ E-Step of EM for fitting GMMs. Computes estimated sufficient statistics
    ↪(ess),  $p(z|x)$ , using the old model from
    the previous iteration. In the GMM case they are often referred to as
    ↪"responsibilities".
    :param samples: samples "x" to compute ess for (shape: [N, dim])
    :param weights: weights (i.e.,  $p(z)$ ) of old model (shape: [num_components])
    :param means: means of old components  $p(x|z)$  (shape: [num_components, dim])
    :param covariances: covariances of old components  $p(x|z)$  (shape:
    ↪[num_components, dim, dim])
    :return: Responsibilities  $p(z|x)$  (Shape: [N x num_components])
    """
    #####
    # TODO Implement the E-Step for EM for Gaussian Mixtrue Models.

    n = samples.shape[0] #number of samples
    num_comp = weights.shape[0] #number of componets

    p_x_comp = np.zeros([n,num_comp])

    # one component  $p(x,z) = p(z)p(x|z)$ 
    for i in range(num_comp):
        p_x_comp[:,i] = weights[i] * mvn_pdf(samples, means[i],covariances[i])

    # sum of components  $sum[p(x,z)]=sum[p(z)p(x|z)]$ 
    p_x = np.sum(p_x_comp,axis=-1)[:,None]

    #  $q(z) = p(z|x) = p(z)p(x|z) / sum[p(x,z)]=sum[p(z)p(x|z)]$ 
    return p_x_comp / p_x

    #####

def m_step(samples: np.ndarray, responsibilities: np.ndarray) -> Tuple[np.
    ↪ndarray, np.ndarray, np.ndarray]:
    """ M-Step of EM for fitting GMMs. Computes new parameters given samples
    ↪and responsibilities  $p(z|x)$ 
    :param samples: samples "x" to fit model to (shape: [N, dim])
    :param responsibilities:  $p(z|x)$  (Shape: [N x num_components]), as computed
    ↪by E-step
    :return: - new weights  $p(z)$  (shape [num_components])
             - new means of components  $p(x|z)$  (shape: [num_components, dim])
             - new covariances of components  $p(x|z)$  (shape: [num_components,
    ↪dim, dim])
    """
    #####

```

```

    # TODO: Implement the M-Step for EM for Gaussian Mixture models. You are
    ↪not allowed to use any for loops!

    # Hint: Writing it directly without for loops is hard, especially if you
    ↪are not experienced with broadcasting.

    # It's maybe easier to first implement it using for loops and then try
    ↪getting rid of them, one after another.

    n = samples.shape[0] # number of samples

    p_z = np.sum(responsibilities,axis=0) / n #new weights in slide 43

    means_new = np.dot(responsibilities.T,samples) / np.
    ↪sum(responsibilities,axis=0)[:None]
    means_new = means_new # new means in slide 43

    x_mean = samples[:,None,:] - means_new[None,:,:]
    x_covar = x_mean[:, :, :,None] @ x_mean[:, :, None,:]
    covar_new = np.sum(responsibilities[:, :,None,None] * x_covar , axis=0) /
    ↪np.sum(responsibilities,axis=0)[:None,None]
    # new covariances in slide 43

    return p_z,means_new,covar_new

#####

```

We wrap out functions with the actual algorithm, iterating E and M step

```

[26]: def fit_gaussian_mixture(samples: np.ndarray, num_components: int, num_iters:
    ↪int = 30, vis_interval: int = 5):
    """Fits a Gaussian Mixture Model using the Expectation Maximization
    ↪Algorithm
    :param samples: Samples to fit the model to (shape: [N, dim]
    :param num_components: number of components of the GMM
    :param num_iters: number of iterations
    :param vis_interval: After how many iterations to generate the next plot
    :return: - final weights  $p(z)$  (shape [num_components])
             - final means of components  $p(x|z)$  (shape: [num_components, dim])
             - final covariances of components  $p(x|z)$  (shape: [num_components,
    ↪dim, dim]
             - log_likelihoods: log-likelihood of data under model after each
    ↪iteration (shape: [num_iters])
    """
    # Initialize Model: We initialize with means randomly picked from the data,
    ↪unit covariances and uniform
    # component weights. This works here but in general smarter initialization
    ↪techniques might be necessary, e.g.,
    # k-means

```

```

initial_idx = np.random.choice(len(samples), num_components, replace=False)
means = samples[initial_idx]
covs = np.tile(np.eye(data.shape[-1])[None, ...], [num_components, 1, 1])
weights = np.ones(num_components) / num_components

# bookkeeping:
log_likelihoods = np.zeros(num_iters)

# iterate E and M Steps
for i in range(num_iters):
    responsibilities = e_step(samples, weights, means, covs)
    weights, means, covs = m_step(samples, responsibilities)

    # Plotting
    if i % vis_interval == 0:
        visualize_2d_gmm(data, weights, means, covs, title="After Iteration_
→{:02d}".format(i))

    log_likelihoods[i] = gmm_log_likelihood(samples, weights, means, covs)
return weights, means, covs, log_likelihoods

```

Finally we load some data and run the algorithm. Feel free to play around with the parameters a bit.

```

[27]: ## ADAPTABLE PARAMETERS:

np.random.seed(0)
num_components = 5
num_iters = 30
vis_interval = 5

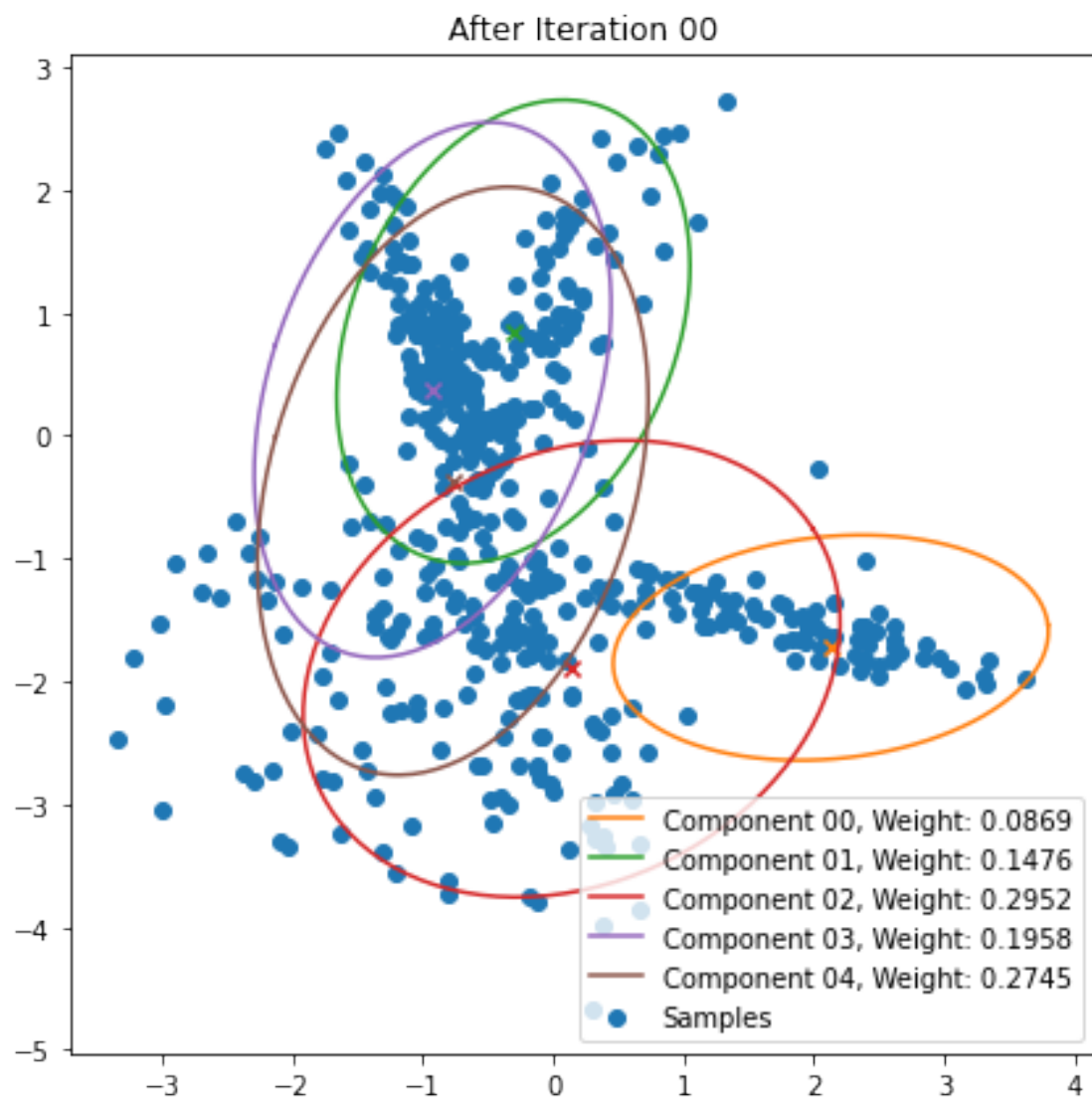
# CHOOSE A DATASET
#data = np.load("samples_1.npy")
data = np.load("samples_2.npy")
#data = np.load("samples_3.npy")
#data = np.load("samples_u.npy")

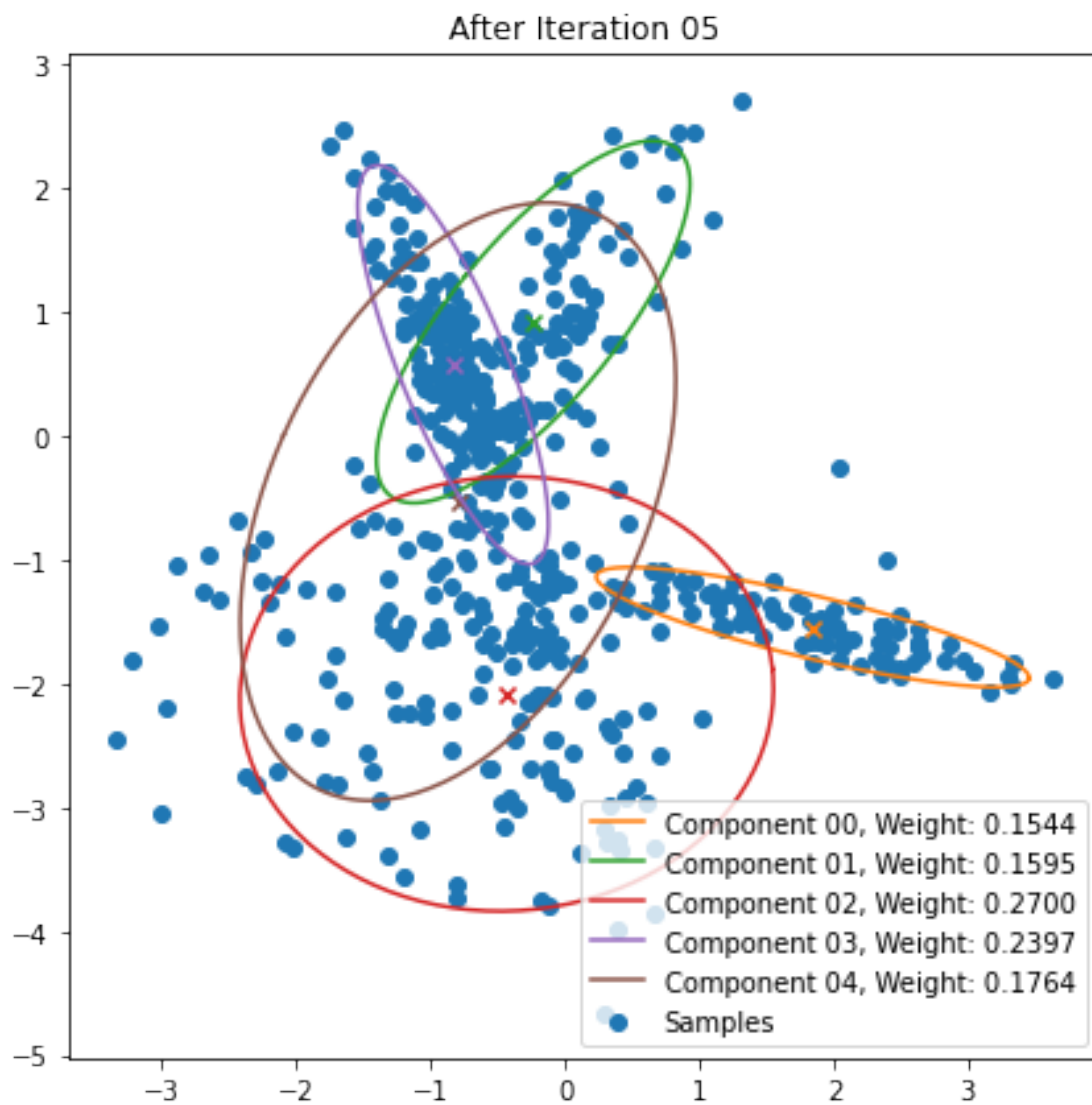
# running and plotting
final_weights, final_means, final_covariances, log_likelihoods = \
    fit_gaussian_mixture(data, num_components, num_iters, vis_interval)
visualize_2d_gmm(data, final_weights, final_means, final_covariances,
→title="Final Model")

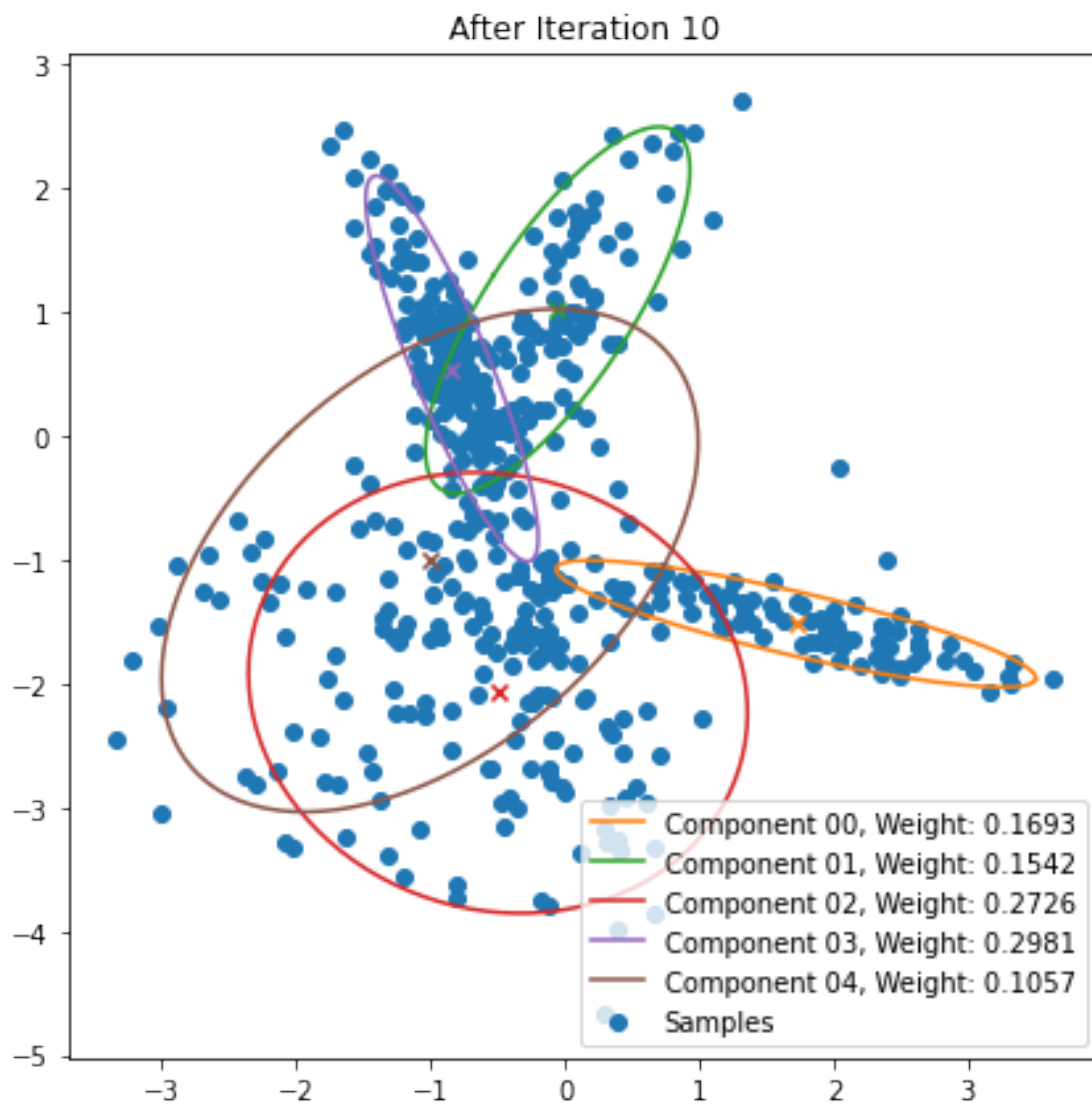
plt.figure()
plt.title("Log-Likelihoods over time")
plt.plot(log_likelihoods)
plt.xlabel("iteration")

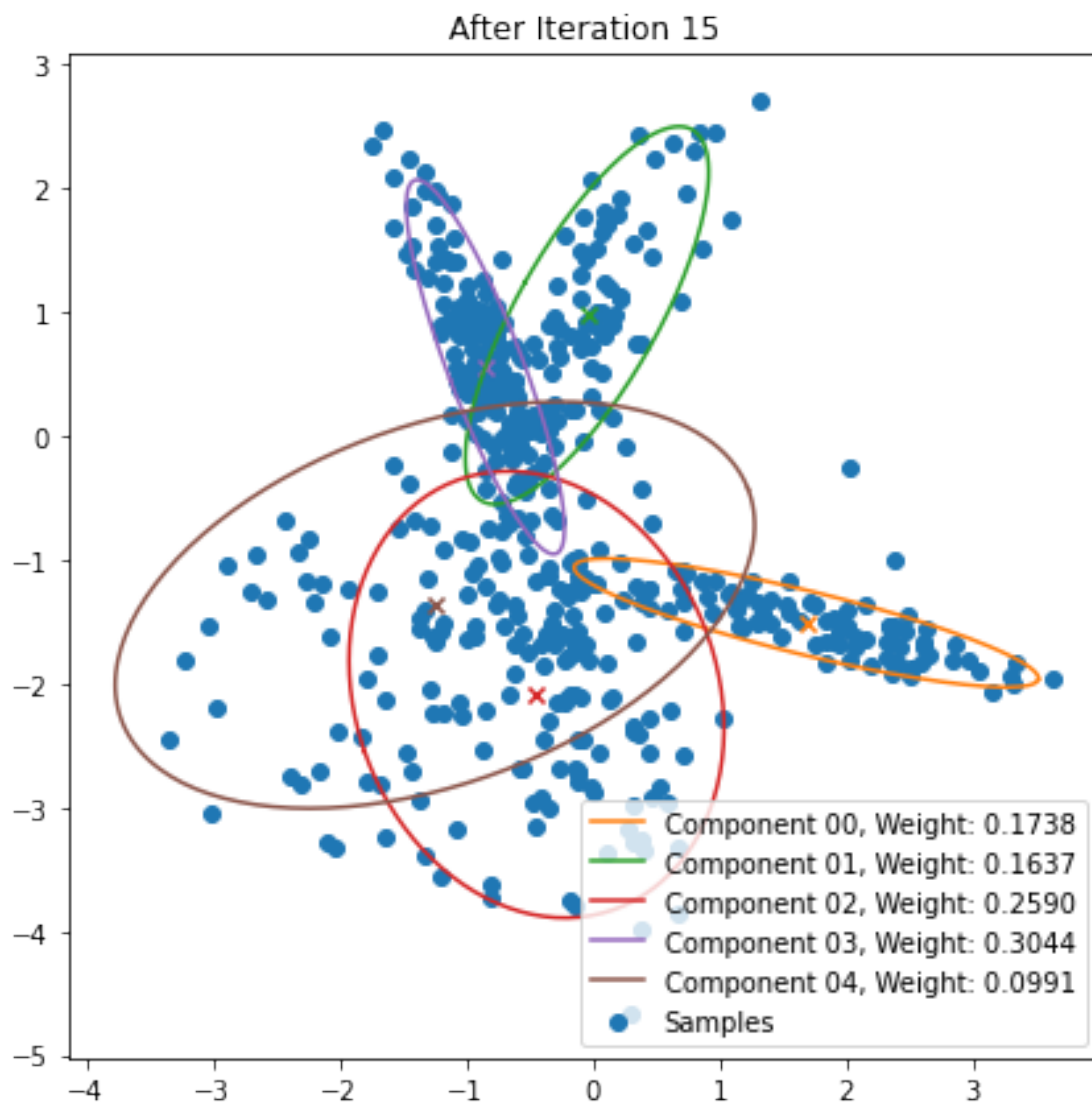
```

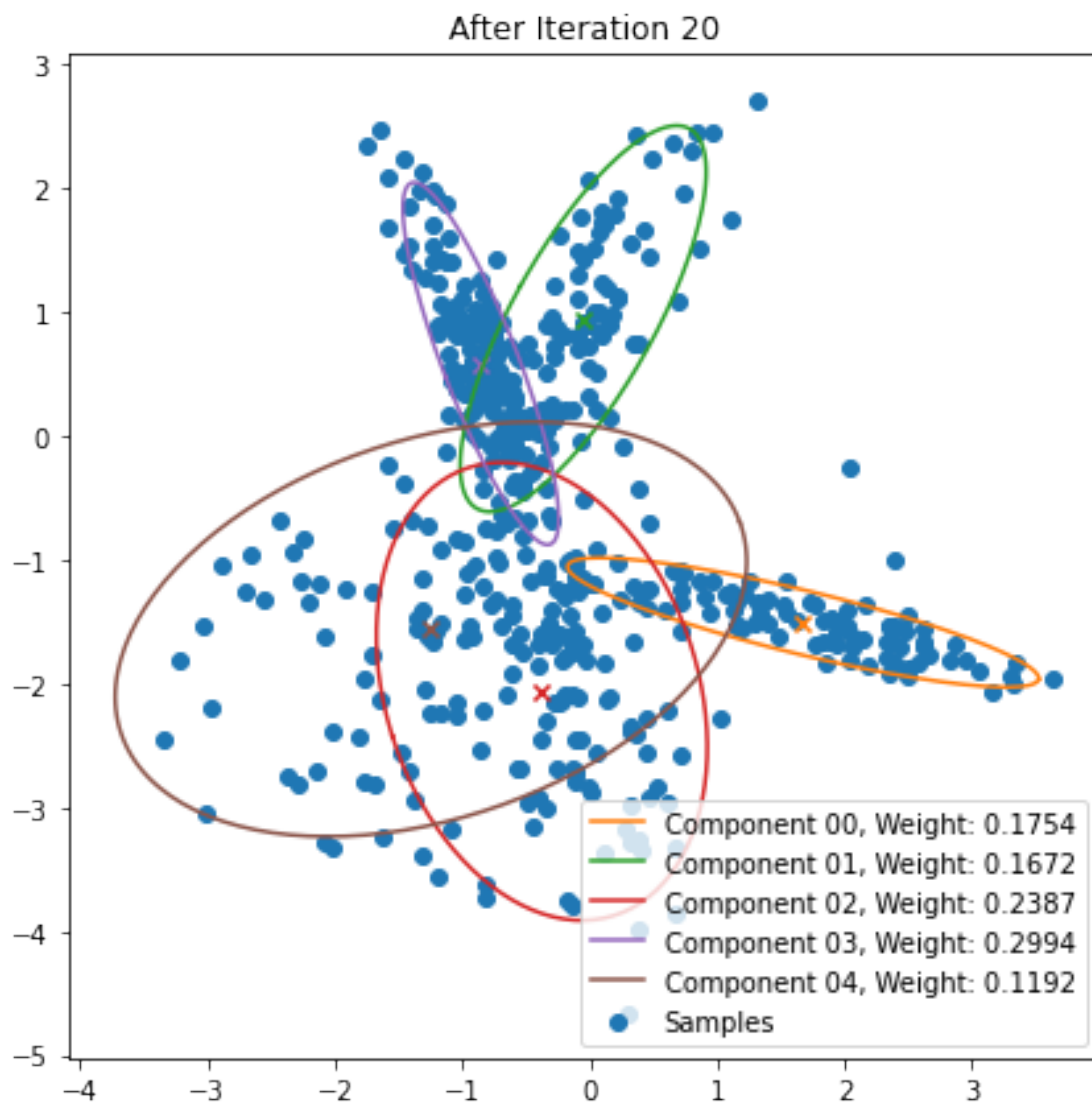
```
plt.ylabel("log-likelihood")  
plt.show()
```

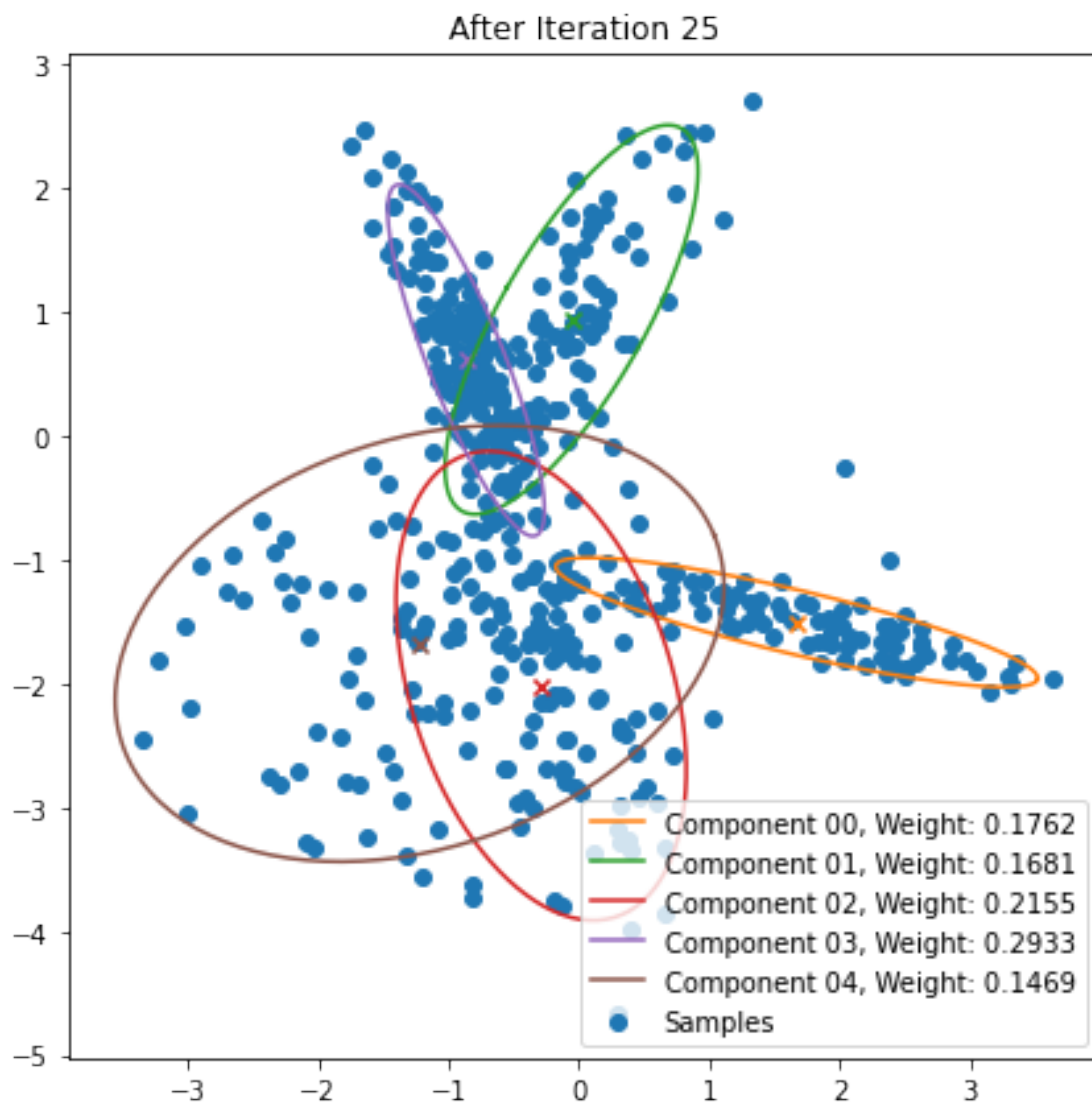


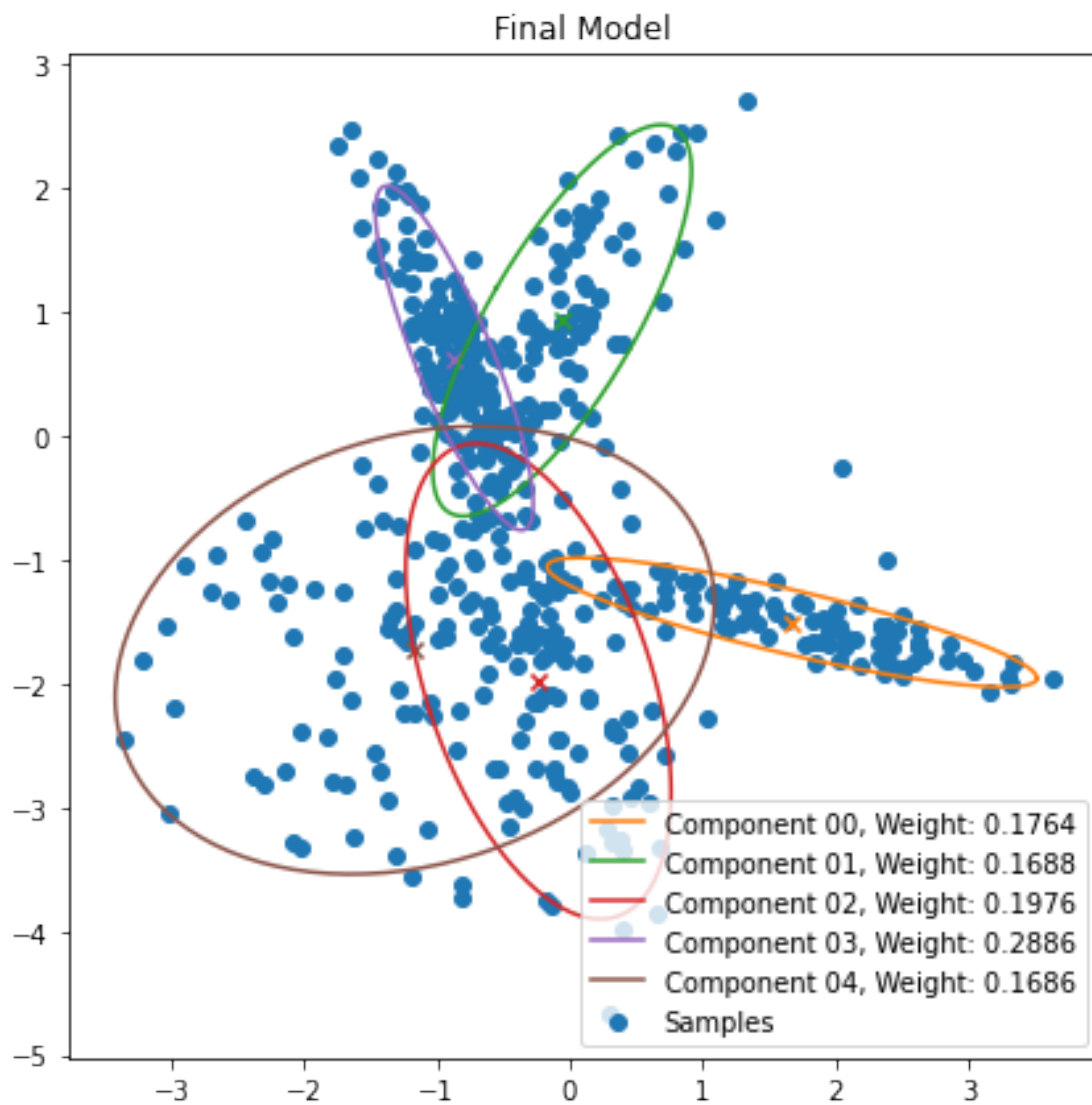


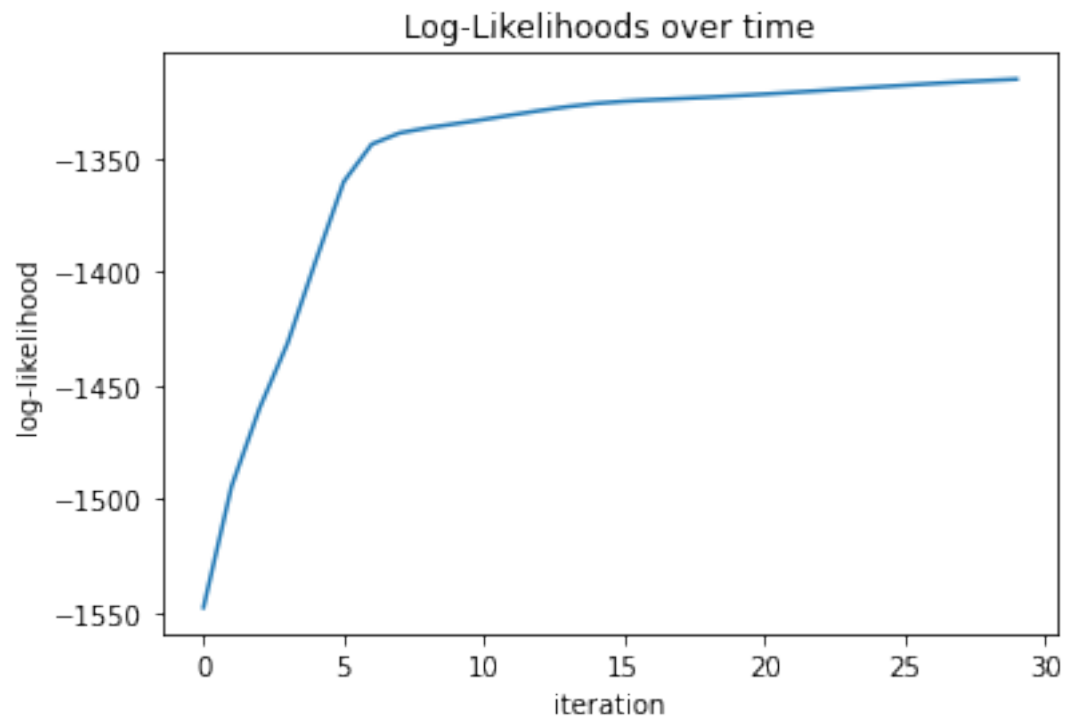












[]: