

hw3__mixtures__of__gaussians__EM

November 24, 2023

1 Mixtures of Gaussians and Expectation Maximization algorithm

The following manuscript from PRML derives the EM algorithm for Mixtures of Gaussians. You can refer to PRML 9.2 for more details.

Recall that the Gaussian mixture distribution can be written as a linear superposition of Gaussians,

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

Let us introduce a K -dimensional binary random variable \mathbf{z} having a 1 -of- K representation in which a particular element z_k is equal to 1 and all other elements are equal to 0 . The values of z_k therefore satisfy $z_k \in \{0, 1\}$ and $\sum_k z_k = 1$, and we see that there are K possible states for the vector \mathbf{z} according to which element is nonzero. We shall define the joint distribution $p(\mathbf{x}, \mathbf{z})$ in terms of a marginal distribution $p(\mathbf{z})$ and a conditional distribution $p(\mathbf{x} \mid \mathbf{z})$. The marginal distribution over \mathbf{z} is specified in terms of the mixing coefficients π_k , such that

$$p(z_k = 1) = \pi_k$$

where the parameters $\{\pi_k\}$ must satisfy

$$0 \leq \pi_k \leq 1$$

together with

$$\sum_{k=1}^K \pi_k = 1$$

in order to be valid probabilities. Because \mathbf{z} uses a 1 -of- K representation, we can also write this distribution in the form

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}.$$

Similarly, the conditional distribution of \mathbf{x} given a particular value for \mathbf{z} is a Gaussian

$$p(\mathbf{x} \mid z_k = 1) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

which can also be written in the form

$$p(\mathbf{x} \mid \mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}.$$

The joint distribution is given by $p(\mathbf{z})p(\mathbf{x} | \mathbf{z})$, and the marginal distribution of \mathbf{x} is then obtained by summing the joint distribution over all possible states of \mathbf{z} to give

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x} | \mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

Thus the marginal distribution of \mathbf{x} is a Gaussian mixture of the form (9.7). If we have several observations $\mathbf{x}_1, \dots, \mathbf{x}_N$, then, because we have represented the marginal distribution in the form $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$, it follows that for every observed data point \mathbf{x}_n there is a corresponding latent variable \mathbf{z}_n .

We shall use $\gamma(z_k)$ to denote $p(z_k = 1 | \mathbf{x})$, whose value can be found using Bayes' theorem

$$\begin{aligned} \gamma(z_k) \equiv p(z_k = 1 | \mathbf{x}) &= \frac{p(z_k = 1) p(\mathbf{x} | z_k = 1)}{\sum_{j=1}^K p(z_j = 1) p(\mathbf{x} | z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \end{aligned}$$

The log of the likelihood function is given by

$$\log p(\mathbf{X}) = \sum_{n=1}^N \log \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}.$$

Let us begin by writing down the conditions that must be satisfied at a maximum of the likelihood function. Setting the derivatives of $\ln p(\mathbf{X} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ with respect to the means $\boldsymbol{\mu}_k$ of the Gaussian components to zero, we obtain

$$0 = - \sum_{n=1}^N \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}}_{\gamma(z_{nk})} \boldsymbol{\Sigma}_k (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

where we have made use of the form (2.43) for the Gaussian distribution. Note that the posterior probabilities, or responsibilities, given by (9.13) appear naturally on the right-hand side. Multiplying by $\boldsymbol{\Sigma}_k^{-1}$ (which we assume to be nonsingular) and rearranging we obtain

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

where we have defined

$$N_k = \sum_{n=1}^N \gamma(z_{nk})$$

If we set the derivative of $\ln p(\mathbf{X} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ with respect to $\boldsymbol{\Sigma}_k$ to zero, and follow a similar line of reasoning, making use of the result for the maximum likelihood solution for the covariance matrix of a single Gaussian, we obtain

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

which has the same form as the corresponding result for a single Gaussian fitted to the data set, but again with each data point weighted by the corresponding posterior probability and with the denominator given by the effective number of points associated with the corresponding component.

Finally, we maximize $\ln p(\mathbf{X} \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ with respect to the mixing coefficients π_k . Here we must take account of the constraint (9.9), which requires the mixing coefficients to sum to one. This can be achieved using a Lagrange multiplier and maximizing the following quantity

$$\ln p(\mathbf{X} \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

which gives

$$0 = \sum_{n=1}^N \frac{\mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} + \lambda$$

where again we see the appearance of the responsibilities. If we now multiply both sides by π_k and sum over k making use of the constraint (9.9), we find $\lambda = -N$. Using this to eliminate λ and rearranging we obtain

$$\pi_k = \frac{N_k}{N}$$

so that the mixing coefficient for the k^{th} component is given by the average responsibility which that component takes for explaining the data points.

1.1 EM for Gaussian Mixtures

Given a Gaussian mixture model, the goal is to maximize the likelihood function with respect to the parameters (comprising the means and covariances of the components and the mixing coefficients).

1. Initialize the means $\boldsymbol{\mu}_k$, covariances $\boldsymbol{\Sigma}_k$ and mixing coefficients π_k , and evaluate the initial value of the log likelihood.

2. E step. Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}.$$

3. M step. Re-estimate the parameters using the current responsibilities

$$\begin{aligned} \boldsymbol{\mu}_k^{\text{new}} &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \\ \boldsymbol{\Sigma}_k^{\text{new}} &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^T \\ \pi_k^{\text{new}} &= \frac{N_k}{N} \end{aligned}$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}).$$

4. Evaluate the log likelihood

$$\ln p(\mathbf{X} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

and check for convergence of either the parameters or the log likelihood. If the convergence criterion is not satisfied return to step 2 .

```
[ ]: !pip install numpy
      !pip install scipy
      !pip install matplotlib
```

```
Looking in indexes: https://pypi.tuna.tsinghua.edu.cn/simple,
https://pypi.ngc.nvidia.com
Requirement already satisfied: numpy in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (1.21.6)
Looking in indexes: https://pypi.tuna.tsinghua.edu.cn/simple,
https://pypi.ngc.nvidia.com
Requirement already satisfied: scipy in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (1.7.3)
Requirement already satisfied: numpy<1.23.0,>=1.16.5 in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from scipy) (1.21.6)
Looking in indexes: https://pypi.tuna.tsinghua.edu.cn/simple,
https://pypi.ngc.nvidia.com
Requirement already satisfied: matplotlib in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (3.5.3)
Requirement already satisfied: cyclor>=0.10 in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from matplotlib)
(0.11.0)
Requirement already satisfied: fonttools>=4.22.0 in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from matplotlib)
(4.38.0)
Requirement already satisfied: kiwisolver>=1.0.1 in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from matplotlib)
(1.4.5)
Requirement already satisfied: numpy>=1.17 in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from matplotlib)
(1.21.6)
Requirement already satisfied: packaging>=20.0 in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from matplotlib)
(23.2)
Requirement already satisfied: pillow>=6.2.0 in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from matplotlib)
(9.5.0)
Requirement already satisfied: pyparsing>=2.2.1 in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from matplotlib)
(3.1.1)
Requirement already satisfied: python-dateutil>=2.7 in
```

```

/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from matplotlib)
(2.8.2)
Requirement already satisfied: typing-extensions in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from
kiwisolver>=1.0.1->matplotlib) (4.7.1)
Requirement already satisfied: six>=1.5 in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from python-
dateutil>=2.7->matplotlib) (1.16.0)

```

```

[ ]: import numpy as np
import matplotlib.pyplot as plt
from matplotlib.patches import Ellipse
from scipy.stats import multivariate_normal

# The seed is fixed for reproducibility.
np.random.seed(42)

```

```

[ ]: def log_likelihood(X, pi, mu, sigma):
    N, d = X.shape
    K = len(pi)
    ll = 0

    ↳ #####
    # TODO: ┐
    # ↳
    # Calculate the log-likelihood (while this is not essential for vanilla EM) ┐
    # ↳
    # Hint: try use multivariate_normal. ┐
    # ↳

    ↳ #####
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    # print("k = ", k, ", n = ", n, ", d = ", d)
    # k = 5 , n = 5000, d = 2
    # pi : 5*1
    # mu : 5*2
    # sigma : 5*2*2

    # l = log P(X/mu, sigma, pi) = \sum_{n=1}^N ln{ \sum_{k=1}^K pi_k * ┐
    ↳ N(x_n/mu_k, sigma_k) }

    ll = 0
    for k in range(K):
        ll += pi[k] * multivariate_normal.pdf(X, mu[k], sigma[k])
    ll = np.log(ll).sum()

```

```
# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
return ll
```

```
[ ]: # Draw color points according to their cluster using pyplot, you should not
      ↪ edit this function and can skip it safely.
def draw(X, pi, mu, sigma, iter):
    n, d = X.shape
    k = len(mu)
    gamma = np.zeros((n, k))
    for i in range(n):
        for j in range(k):
            gamma[i, j] = pi[j] * multivariate_normal(mean=mu[j], cov=sigma[j]).
            ↪ pdf(X[i])
        gamma[i] /= np.sum(gamma[i])
    y = np.argmax(gamma, axis=1)
    plt.scatter(X[:, 0], X[:, 1], c=y, s=1)
    # plot the mean of each cluster with striking points
    plt.scatter(mu[:, 0], mu[:, 1], c='black', s=50)
    plt.axis('equal')
    plt.title('iter: {}'.format(iter))
    plt.show()
```

1.1.1 Implementation of EM algorithm

```
[ ]: def EM(X, K, max_iter=10, plot=True):

    N, d = X.shape
    pi = np.ones(K) / K
    # k-means++ initialization
    mu = np.zeros((K, d))
    mu[0] = X[np.random.choice(N)]
    for j in range(1, K):
        dist = np.zeros(N)
        for i in range(N):
            dist[i] = np.min(np.sum((X[i] - mu[:j]) ** 2, axis=1))
        mu[j] = X[np.random.choice(N, p=dist / np.sum(dist))]
    sigma = np.array([np.eye(d) for _ in range(K)])
    ll = log_likelihood(X, pi, mu, sigma)

    draw(X, pi, mu, sigma, '-1')

    for iter in range(max_iter):
        # E-step
```

```

    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    # gamma(z_{nk}) = pi_k * N(x_n/mu_k, sigma_k) / \sum_{j=1}^K pi_j *
    N(x_n/mu_j, sigma_j)

    gamma = np.zeros((N, K))
    for k in range(K):
        gamma[:, k] = pi[k] * multivariate_normal.pdf(X, mu[k], sigma[k])
    gamma /= gamma.sum(axis=1, keepdims=True)

    # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    # M-step

    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    # N_k = \sum_{n=1}^N gamma(z_{nk})
    N_k = gamma.sum(axis=0)

    # pi_k = N_k / N
    pi = N_k / N

    # mu_k = \sum_{n=1}^N gamma(z_{nk}) * x_n / N_k
    mu = gamma.T.dot(X) / N_k.reshape(-1, 1)

```

```

# sigma_k = \sum_{n=1}^N gamma(z_{nk}) * (x_n - mu_k) * (x_n - mu_k)^T /
↪ N_k
sigma = np.zeros((K, d, d))
for k in range(K):
    sigma[k] = (X - mu[k]).T.dot(np.diag(gamma[:, k])).dot(X - mu[k]) / ↪
↪ N_k[k]

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

ll_new = log_likelihood(X, pi, mu, sigma)
if np.abs(ll_new - ll) < 1e-5:
    break
ll = ll_new

# plot the current cluster
if iter % 1 == 0:
    print('Iteration: {}, log-likelihood: {}'.format(iter, ll))
    if plot:
        draw(X, pi, mu, sigma, iter)

return pi, mu, sigma

```

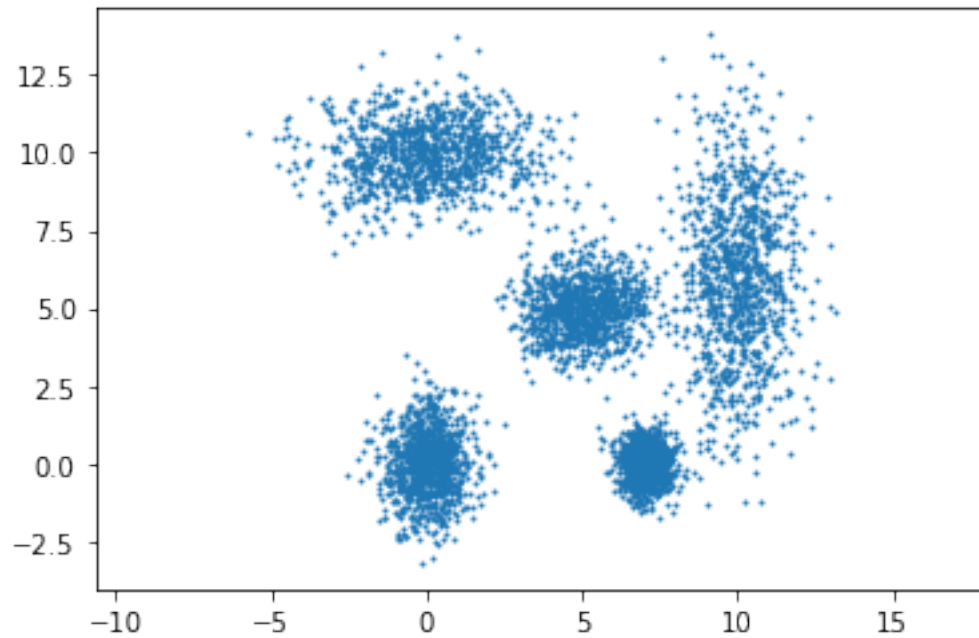
1.1.2 1. Unsupervised Clustering [40 points]

Let's validate our implementation on synthetic data of 2D Gaussian mixture with 5 components.

```

[ ]: # The data is geneated from 5 gaussians with different means, differnet ↪
↪ covariance matrices
K = 5
X = np.loadtxt('synthetic.csv', delimiter=',')
plt.scatter(X[:, 0], X[:, 1], s=1)
plt.axis('equal')
plt.show()

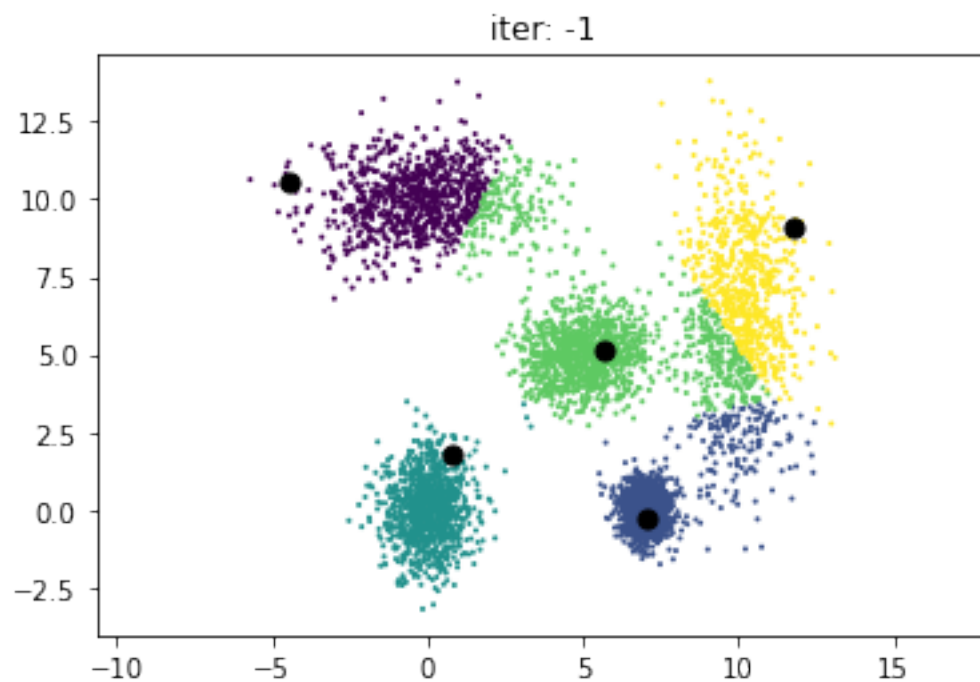
```

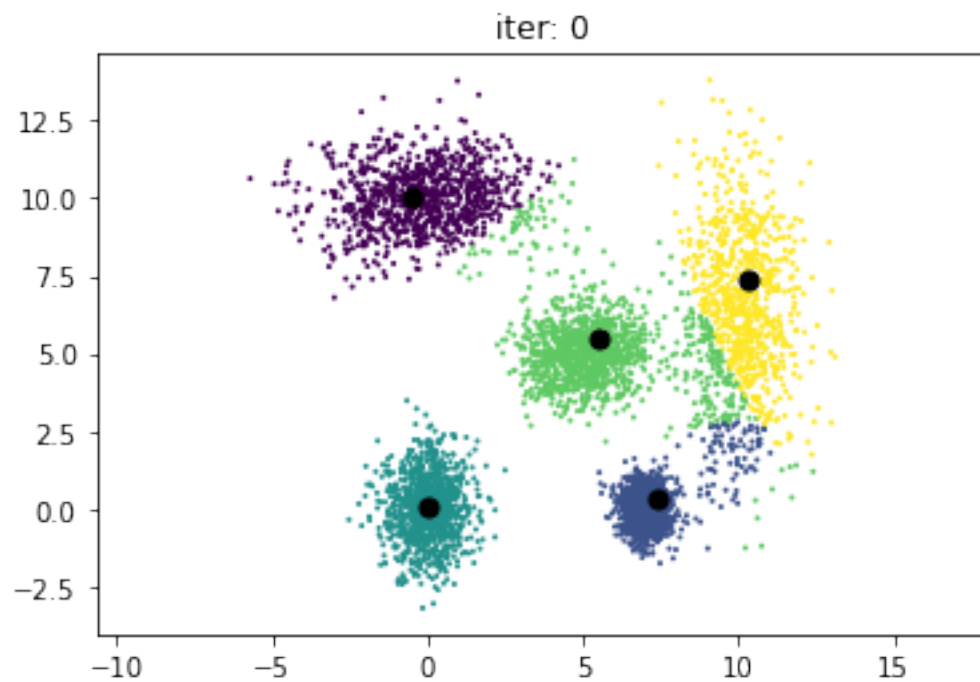
```
[ ]: # run EM algorithm
pi, mu, sigma = EM(X, K)

# print(pi) # pi : 5*1
# print(mu) # mu : 5*2
# print(sigma) # sigma : 5*2*2

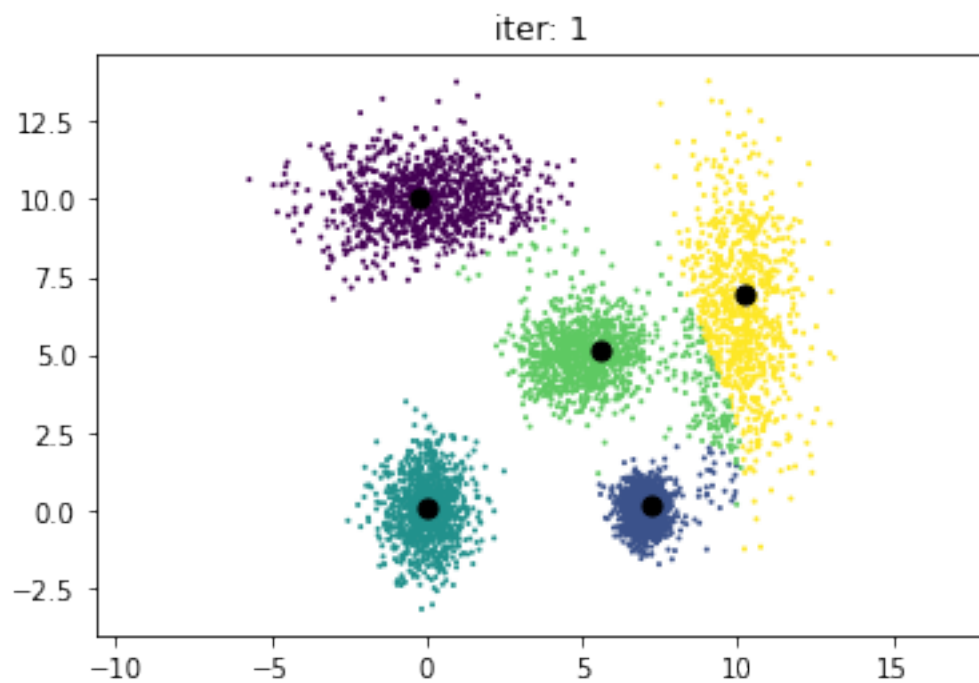
## A quick validation of your implementation:
## log-likelihood should increase after each iteration, and is around ~-20000.
```



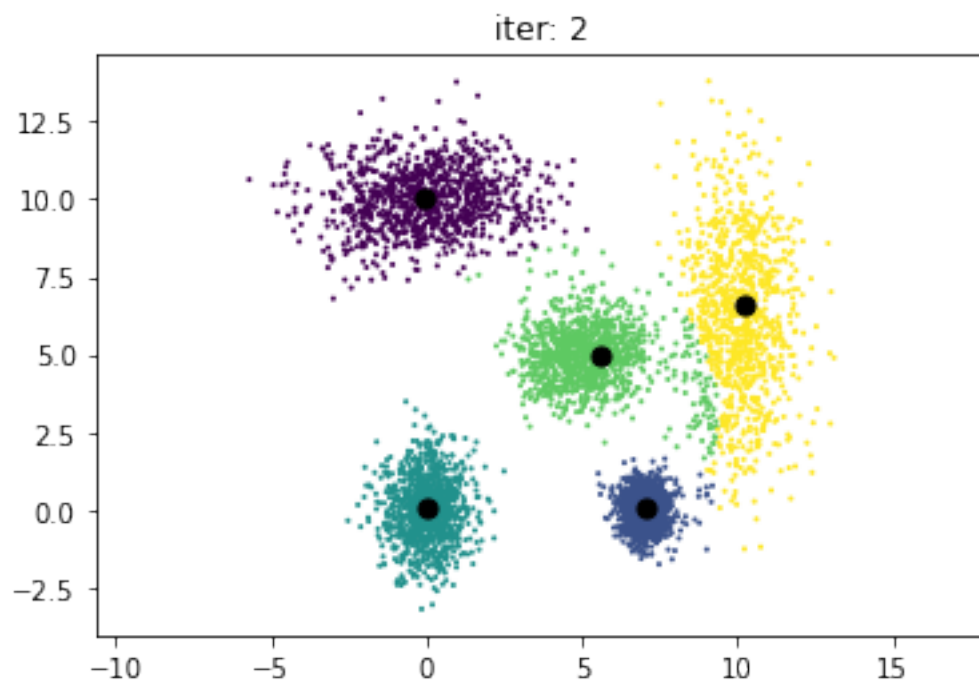
Iteration: 0, log-likelihood: -23110.476800635322



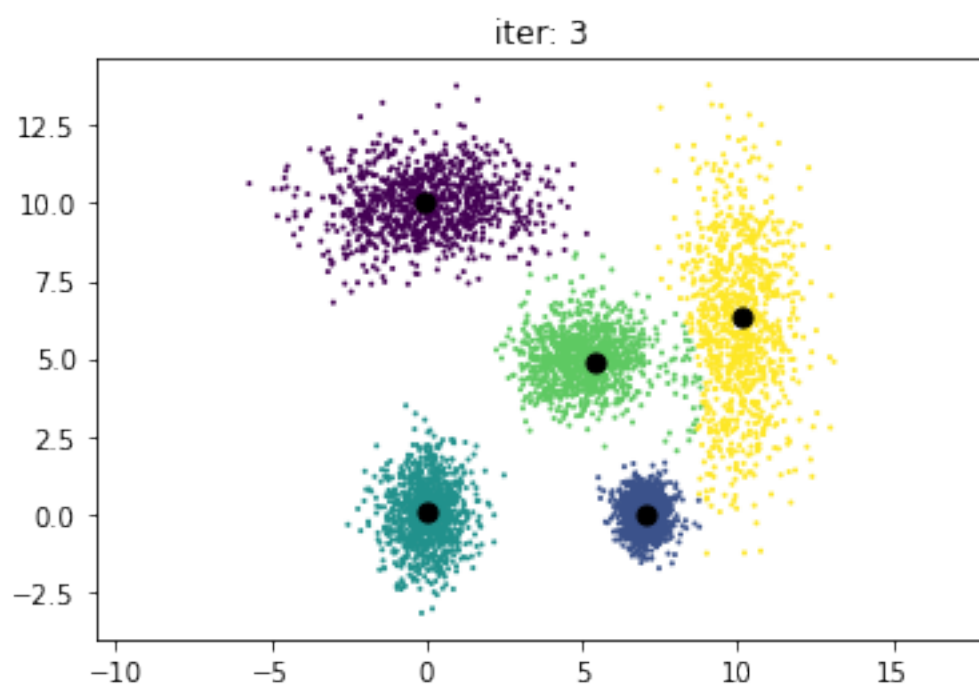
Iteration: 1, log-likelihood: -22611.194696984392



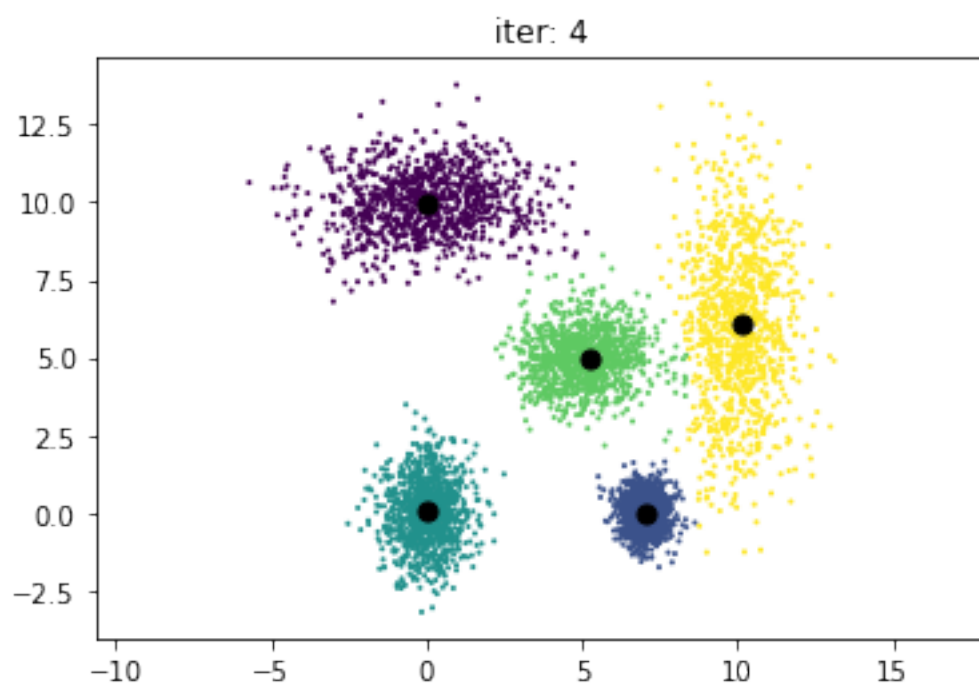
Iteration: 2, log-likelihood: -22178.018903589415



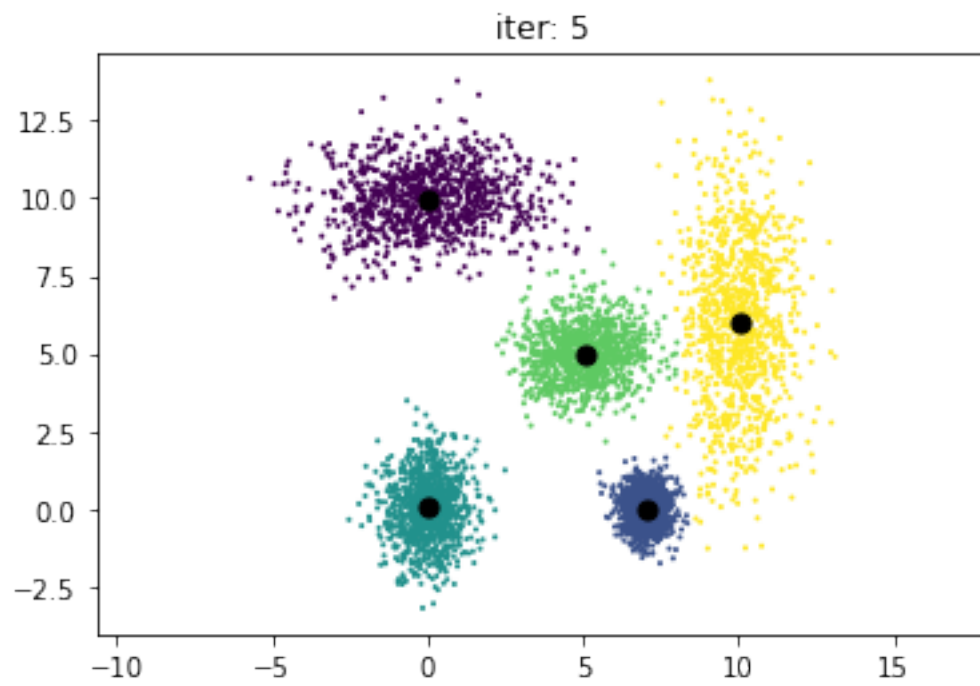
Iteration: 3, log-likelihood: -21952.192471972892



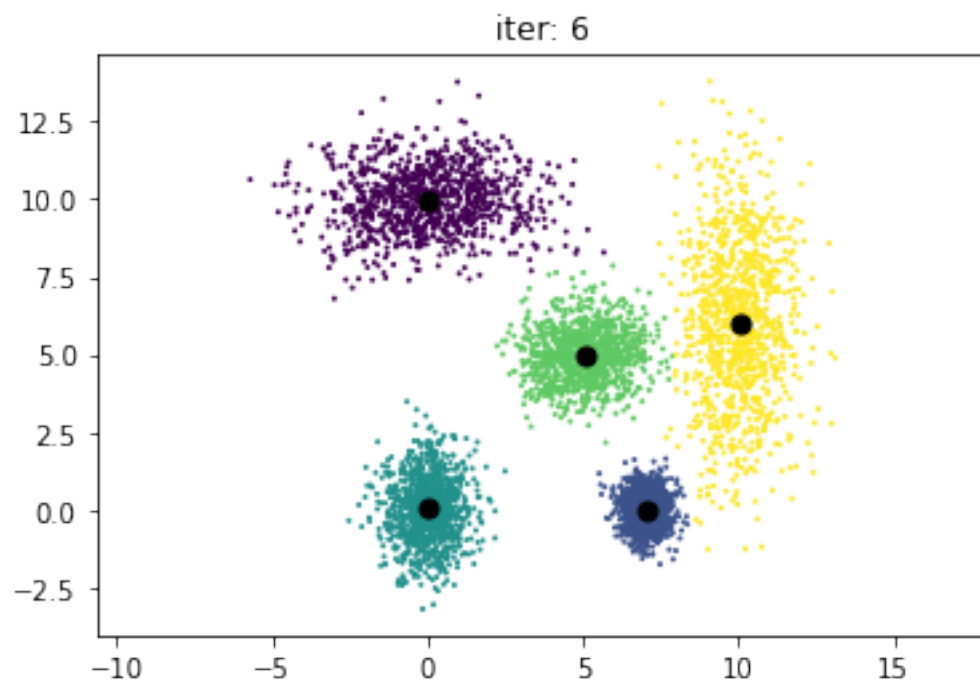
Iteration: 4, log-likelihood: -21804.825654169887



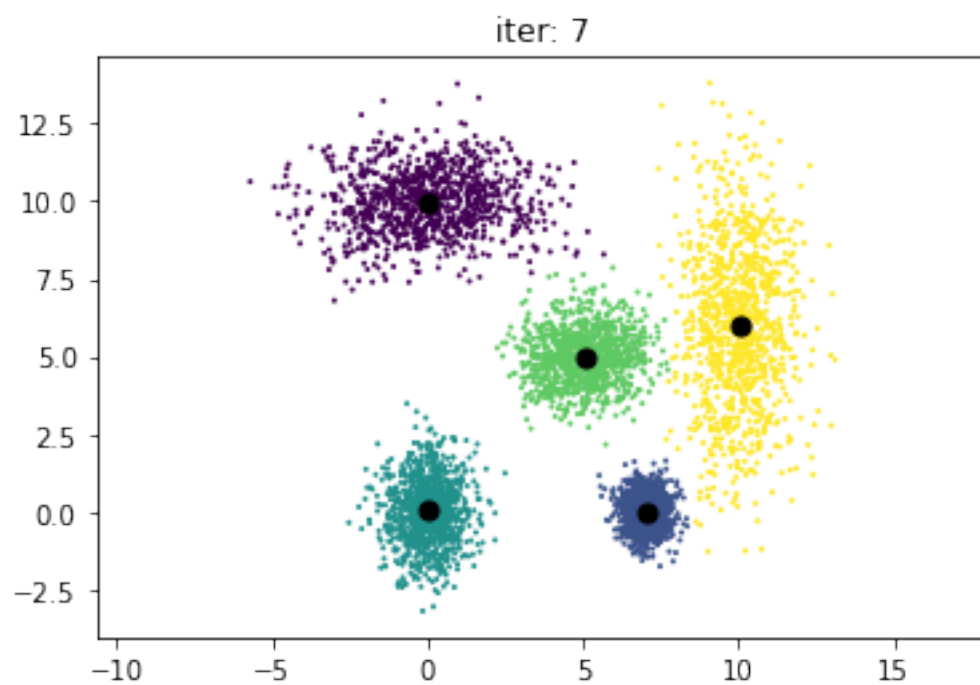
Iteration: 5, log-likelihood: -21736.666532626896



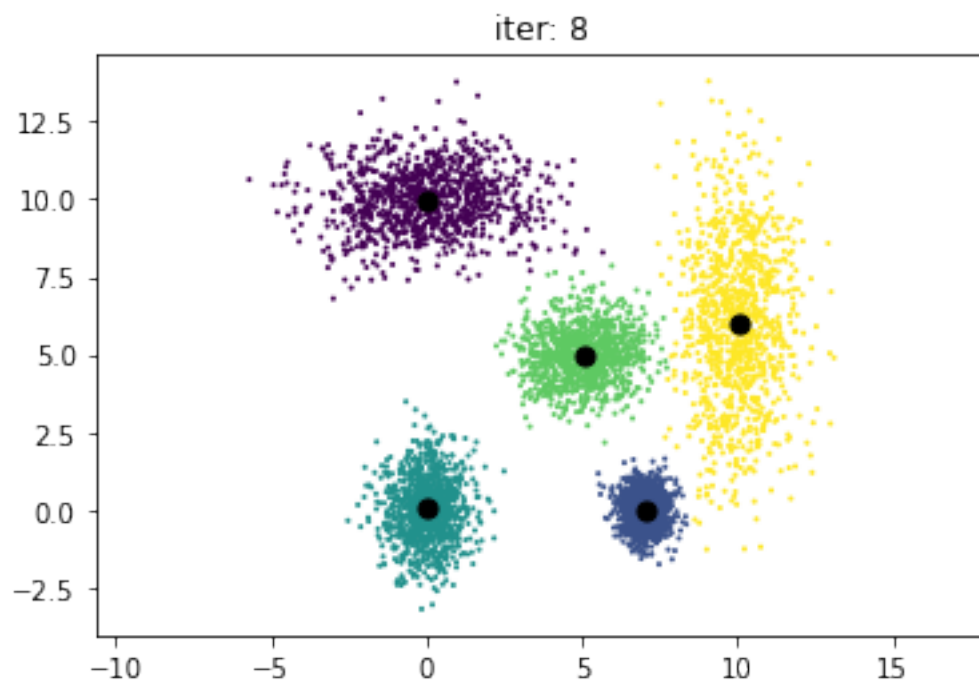
Iteration: 6, log-likelihood: -21726.15813383701



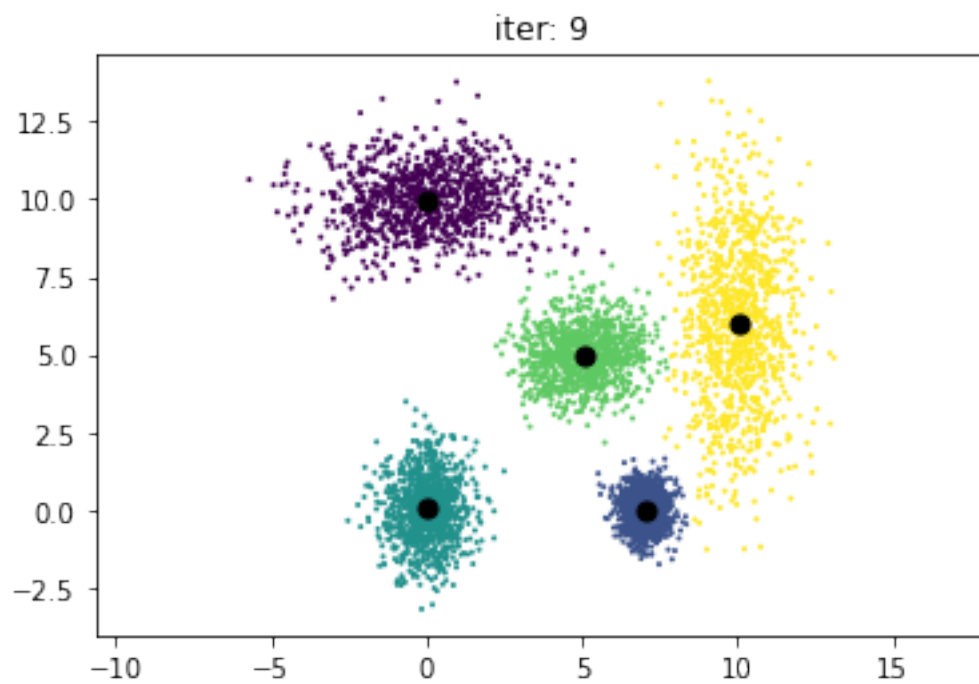
Iteration: 7, log-likelihood: -21725.32599697586



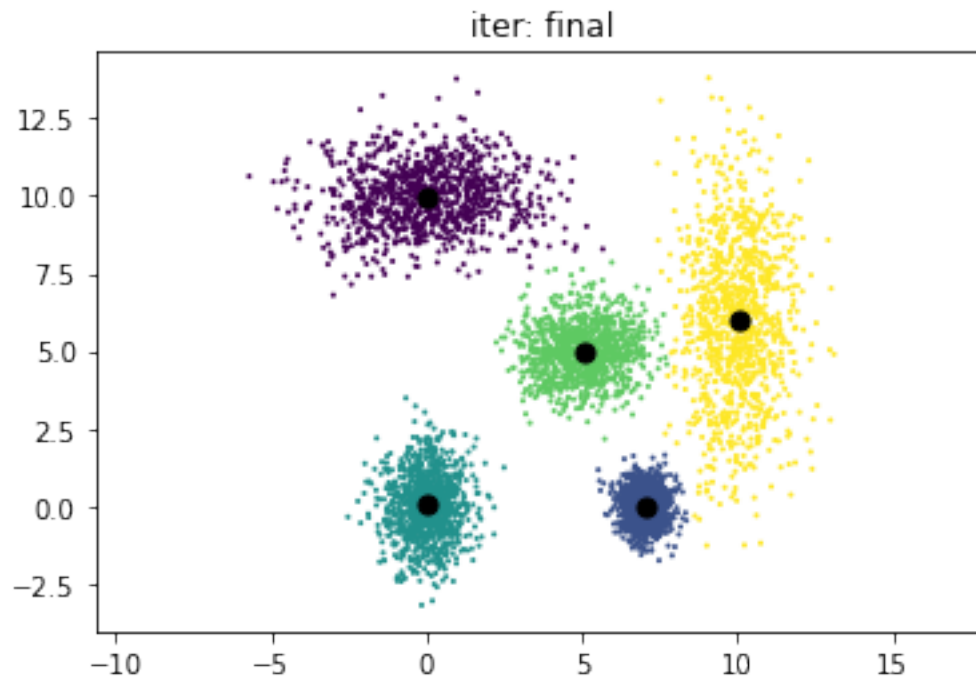
Iteration: 8, log-likelihood: -21725.263882574327



Iteration: 9, log-likelihood: -21725.259092899465

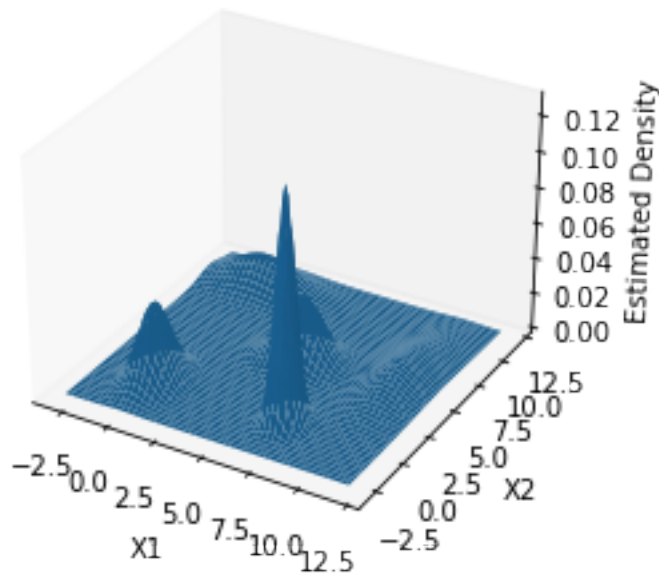


```
[ ]: # The following draws the means and clusters of the final result.
draw(X, pi, mu, sigma, 'final')
```



```
[ ]: # The following plots the estimated density on a 3D grid.
x1 = np.linspace(-3, 12, 100)
x2 = np.linspace(-3, 12, 100)
X1, X2 = np.meshgrid(x1, x2)
Z = np.zeros(X1.shape)
for j in range(K):
    Z += pi[j] * multivariate_normal(mean=mu[j], cov=sigma[j]).pdf(np.
    ↳dstack((X1, X2)))

fig = plt.figure()
ax = fig.add_subplot(projection='3d')
ax.grid(False)
ax.plot_surface(X1, X2, Z)
ax.set_xlabel('X1')
ax.set_ylabel('X2')
ax.set_zlabel('Estimated Density')
plt.show()
```

1.1.3 2. Image Compression [20 points]

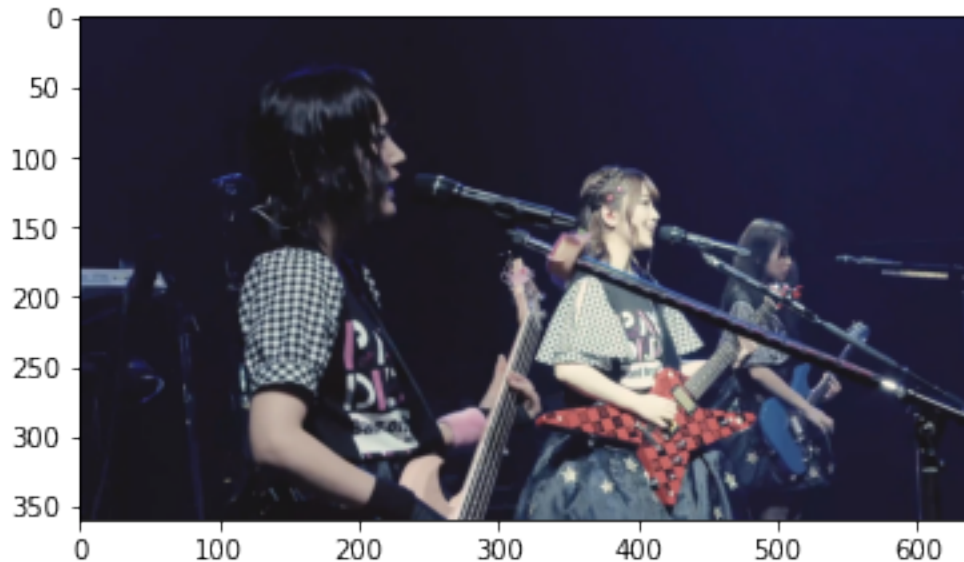
Now let's explore the application of Gaussian mixture model to image compression. We will use the GMM from sklearn package to compress the image as it provides much more stable and faster implementation! The GMM will be trained on the pixels of the image and the cluster assignment for each pixel will be used to replace the original pixel value. The number of clusters is a hyper-parameter that can be tuned to control the compression rate. The higher the number of clusters, the higher the compression rate.

```
[ ]: !pip install scikit-learn
      from sklearn.mixture import GaussianMixture
```

```
Looking in indexes: https://pypi.tuna.tsinghua.edu.cn/simple,
https://pypi.ngc.nvidia.com
Requirement already satisfied: scikit-learn in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (1.0.2)
Requirement already satisfied: numpy>=1.14.6 in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from scikit-learn)
(1.21.6)
Requirement already satisfied: scipy>=1.1.0 in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from scikit-learn)
(1.7.3)
Requirement already satisfied: joblib>=0.11 in
/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from scikit-learn)
(1.3.2)
Requirement already satisfied: threadpoolctl>=2.0.0 in
```

/home/zjk/anaconda3/envs/cs182/lib/python3.7/site-packages (from scikit-learn)
(3.1.0)

```
[ ]: image = plt.imread('compression.png')  
plt.imshow(image)  
plt.show()
```



```
[ ]: # image compression using EM in sklearn  
def GMM_compression(image, k):  
    X = image.reshape(-1, image.shape[2])  
    image_compressed = np.zeros(X.shape)  
  
    # TODO:  
    # Refer to https://scikit-learn.org/stable/modules/generated/sklearn.  
    # mixture.GaussianMixture.html  
    # Create a GaussianMixture object with k components and fit the image data  
    # (X).  
    # Then, predict the cluster label of each pixel and replace each pixel by  
    # its corresponding cluster mean.  
    # Finally, reshape the compressed image to the original image shape.  
  
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
```

```

gm = GaussianMixture(n_components=k)
gm.fit(X)
labels = gm.predict(X)
means = gm.means_
for i in range(len(labels)):
    image_compressed[i] = means[labels[i]]
image_compressed = image_compressed.reshape(image.shape)

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

original_size = image.shape[0] * image.shape[1] * image.shape[2] * 32
# the compressed image stores the label of each pixel, which can be
→ represented by np.log2(k) bits
# and the mean of each cluster, which can be represented by 4 * 32 bits
compressed_size = image_compressed.shape[0] * image_compressed.shape[1] *
→ np.log2(k) + 4 * 32 * k
compression_rate = original_size / compressed_size
plt.title(f'Compression rate: {compression_rate:.2f}, k={k}')
plt.imshow(image_compressed)
plt.show()

```

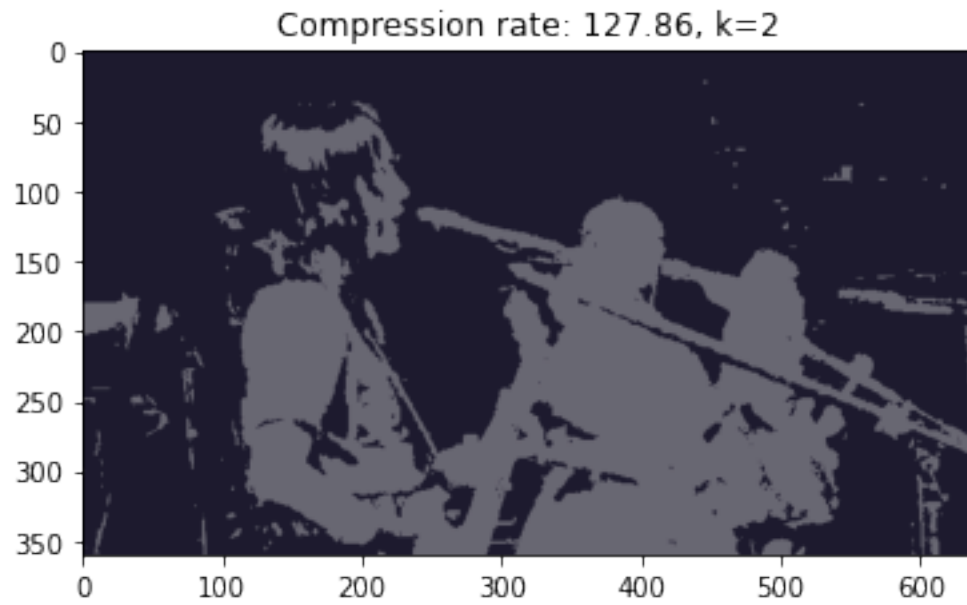
Now let's see how the image looks like after compression under different number of clusters. Here Compression rate is defined as the theoretical number of bits required to represent the image divided by the number of bits required to represent the compressed image.

```

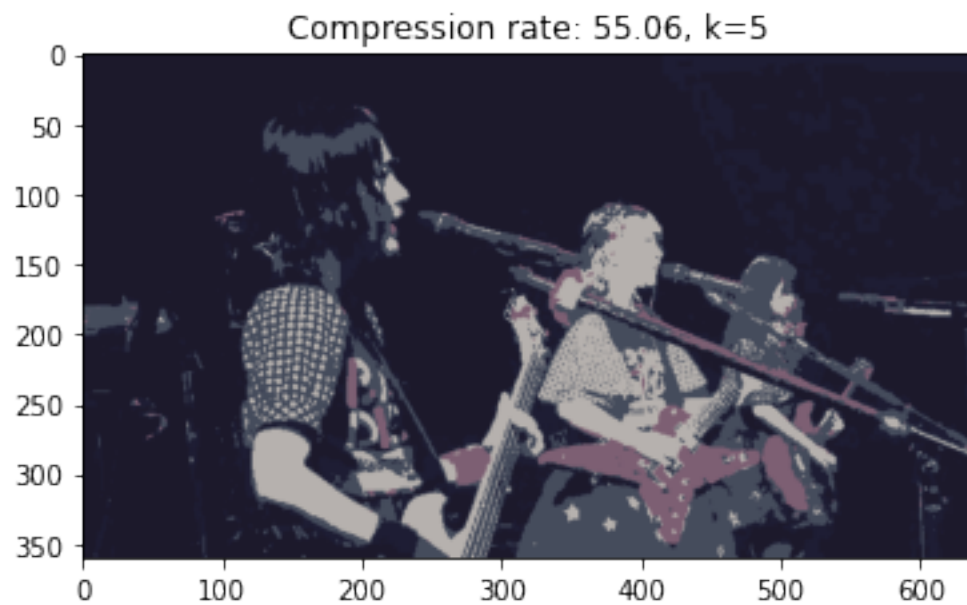
[ ]: for k in [2, 5, 10, 20, 30]:
    GMM_compression(image, k)

```

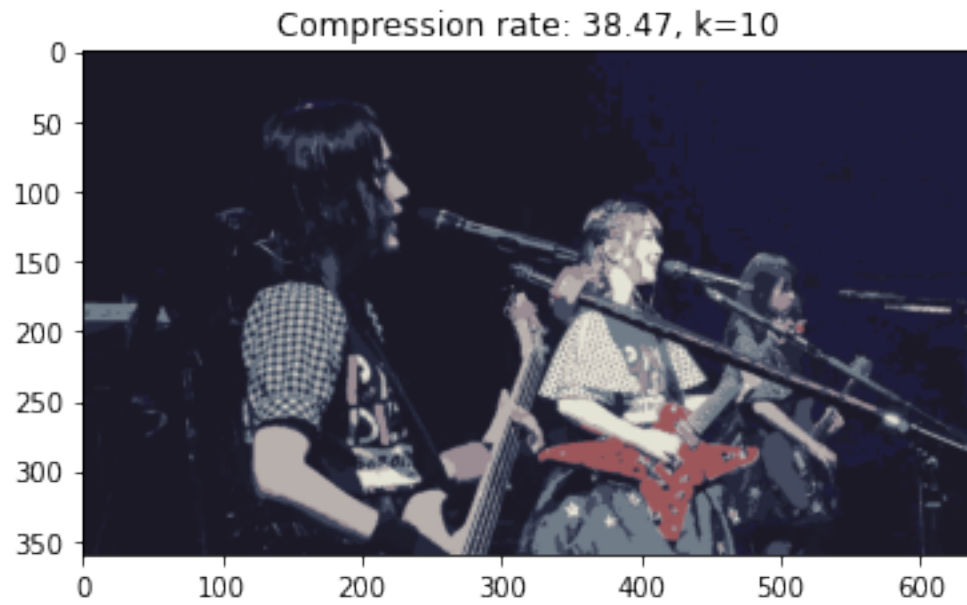
Clipping input data to the valid range for imshow with RGB data ([0..1] for floats or [0..255] for integers).



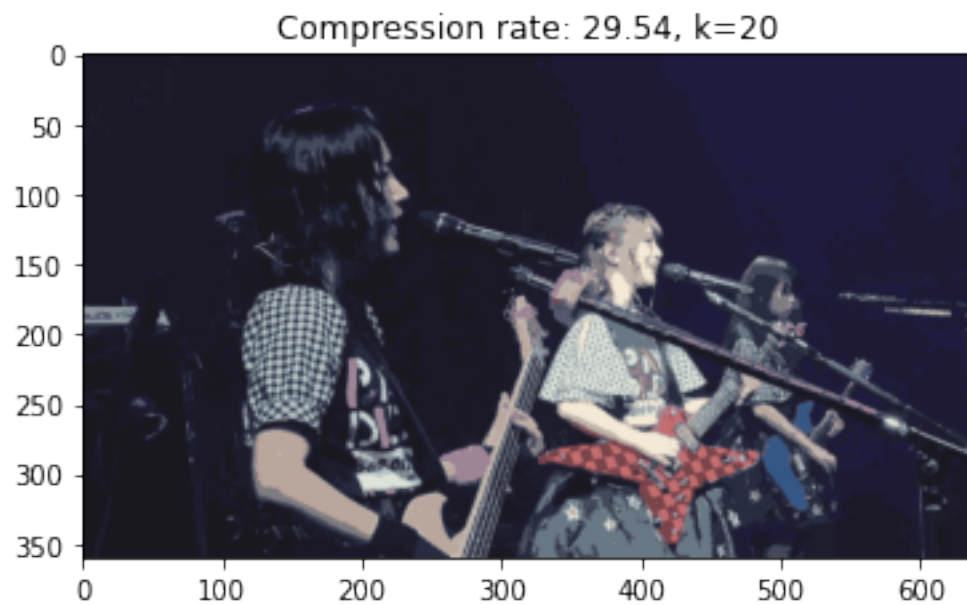
Clipping input data to the valid range for imshow with RGB data ([0..1] for floats or [0..255] for integers).



Clipping input data to the valid range for imshow with RGB data ([0..1] for floats or [0..255] for integers).



Clipping input data to the valid range for imshow with RGB data ([0..1] for floats or [0..255] for integers).



Clipping input data to the valid range for imshow with RGB data ([0..1] for floats or [0..255] for integers).

