```python
import numpy as np
from matplotlib.image import imread
import matplotlib.pyplot as plt


def main():
    A = imread('peppers-large.tiff')
    plt.imshow(A)
    plt.show()

    im_small = imread('peppers-small.tiff')
    plt.imshow(im_small)
    plt.show()

    k = 16
    centroid = kmeans(im_small, k)

    # assign each example in the large image to the closest cluster usi
    dim = A.shape[0]
    A = np.reshape(A, (-1, 3))
    diffs = []
    for c in centroid:
        diff = np.linalg.norm(A - c, axis=1)
        diffs.append(diff)

    # Join the array "diff" along a new axis
    c_i = np.argmin(diffs, axis=0)

    # Compress the large image A
    compress_A = np.zeros((A.shape[0], A.shape[1]), dtype=int)
    for j in range(k):
        ind_j = np.where(c_i == j)
        compress_A[ind_j] = centroid[j]

    compress_A = compress_A.reshape(dim, dim, 3)
    plt.imshow((compress_A))
    plt.show()
```

```python
def kmeans(A, k):
    # initialize centroid by randomly picking k training examples,
    # and set the cluster centroids to be equal to the values of these k examples
    A = np.reshape(A, (-1, 3))
    m = A.shape[0]
    ind = np.random.choice(np.arange(m), size=k, replace=False)
    centroid = A[ind]

    iter = 0
    centroid = np.array(centroid)
    c_i = c_i_old = None
    while c_i_old is None or not np.array_equal(c_i, c_i_old):
        iter += 1
        c_i_old = c_i
        # Assigning each training example x_i to the closest cluster centroid miu_j
        diffs = []
        for c in centroid:
            diff = np.linalg.norm(A - c, axis=1)
            diffs.append(diff)

        c_i = np.argmin(diffs, axis=0)
        #print("c_i_old: ", c_i_old)
        #print("c_i: ", c_i)

        # Moving each cluster centroid miu_j to the mean of the points assigned to it
        miu_js = []
        for j in range(k):
            ind_j = np.where(c_i == j)
            miu_j = A[ind_j].mean(axis=0)
            miu_js.append(miu_j)

        centroid = np.array(miu_js)
        print("iteration: ", iter)
    return centroid


if __name__ == "__main__":
    main()
```
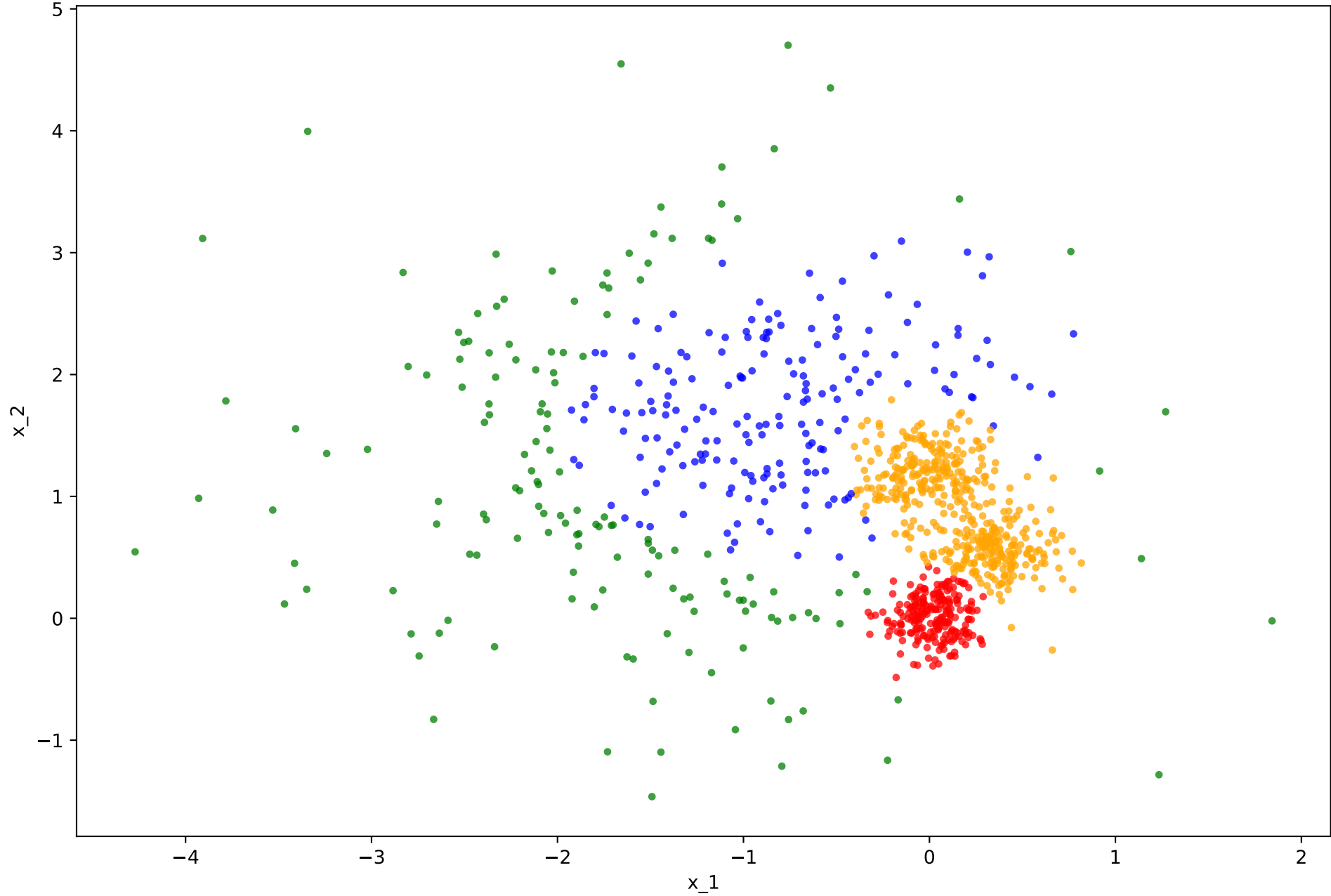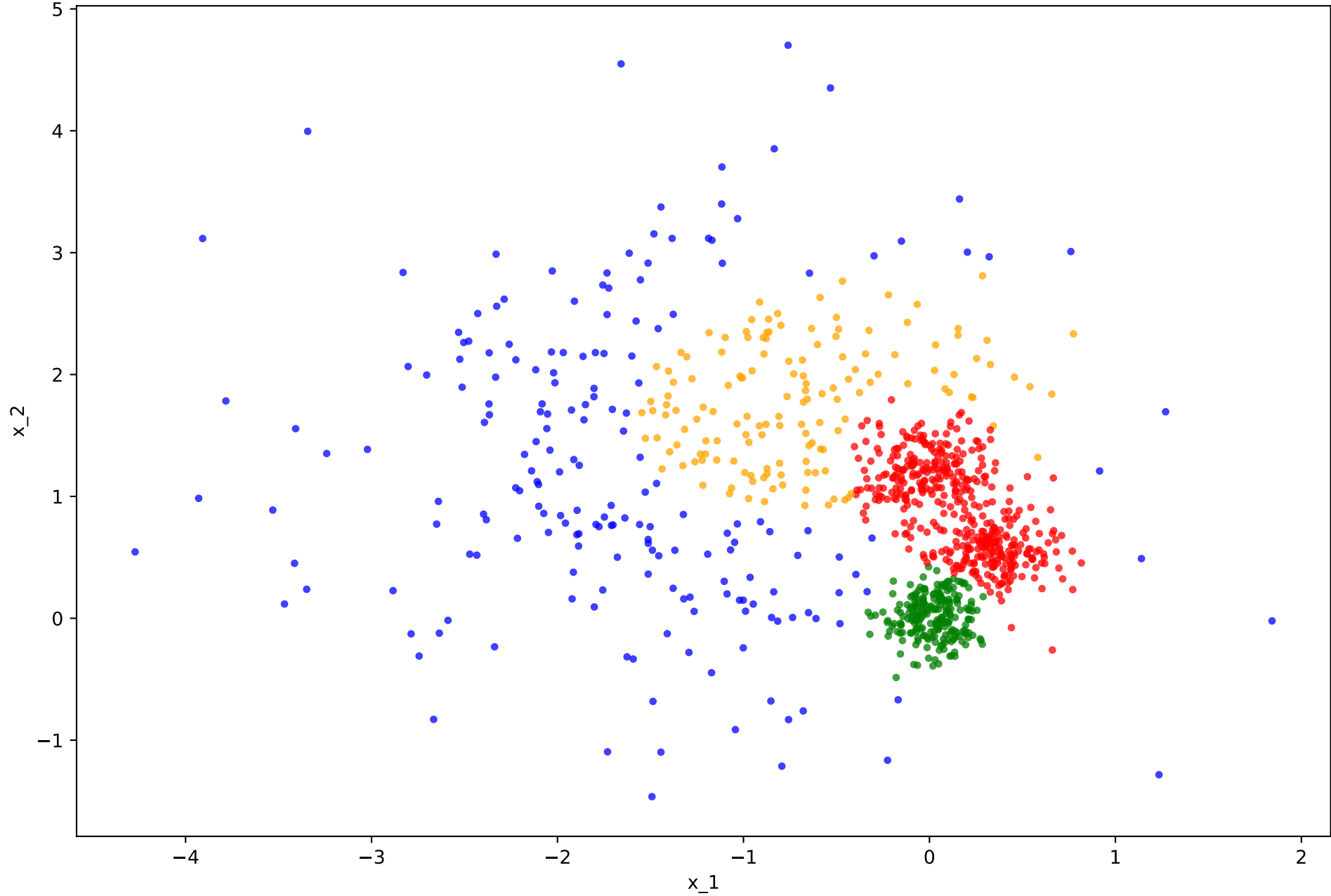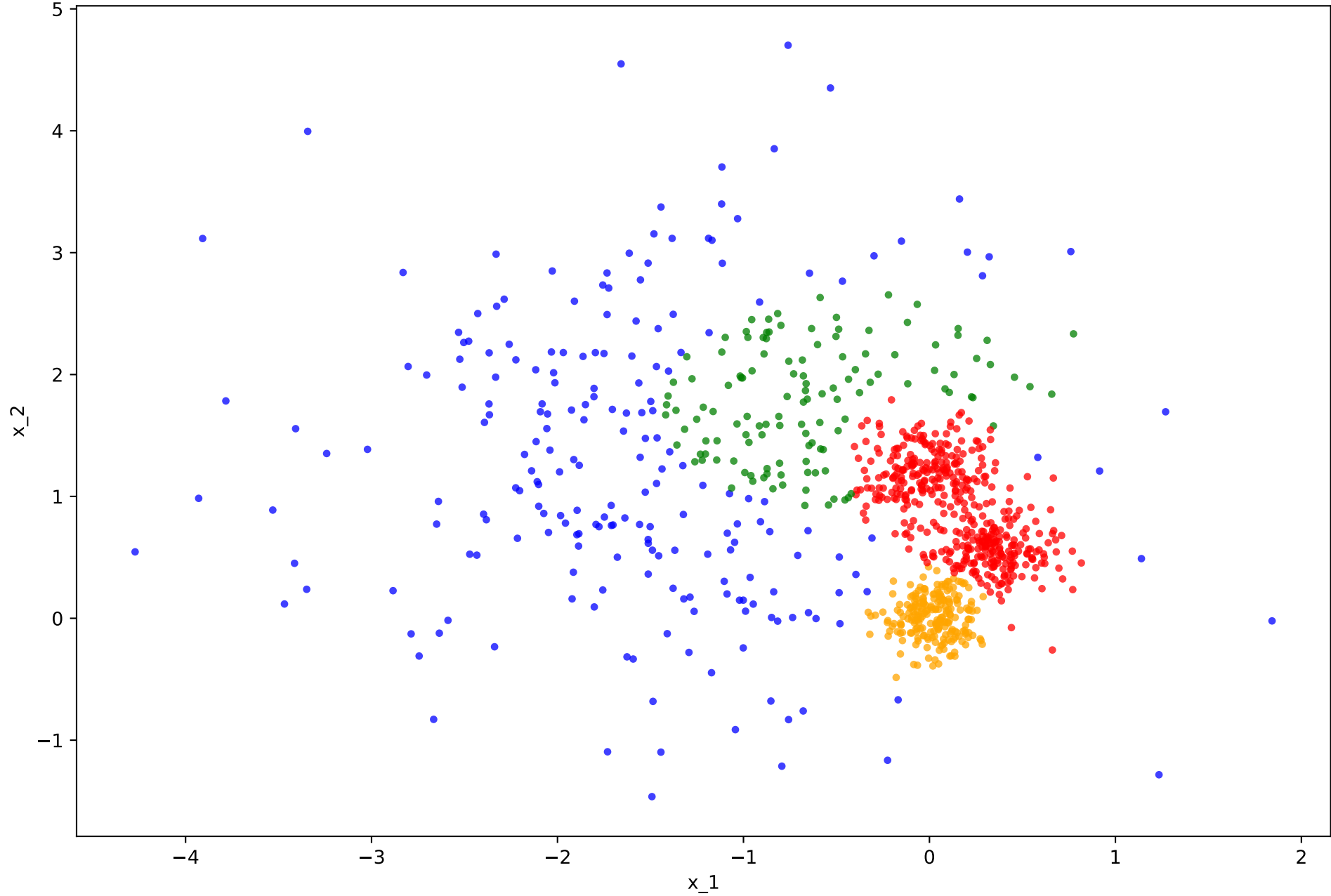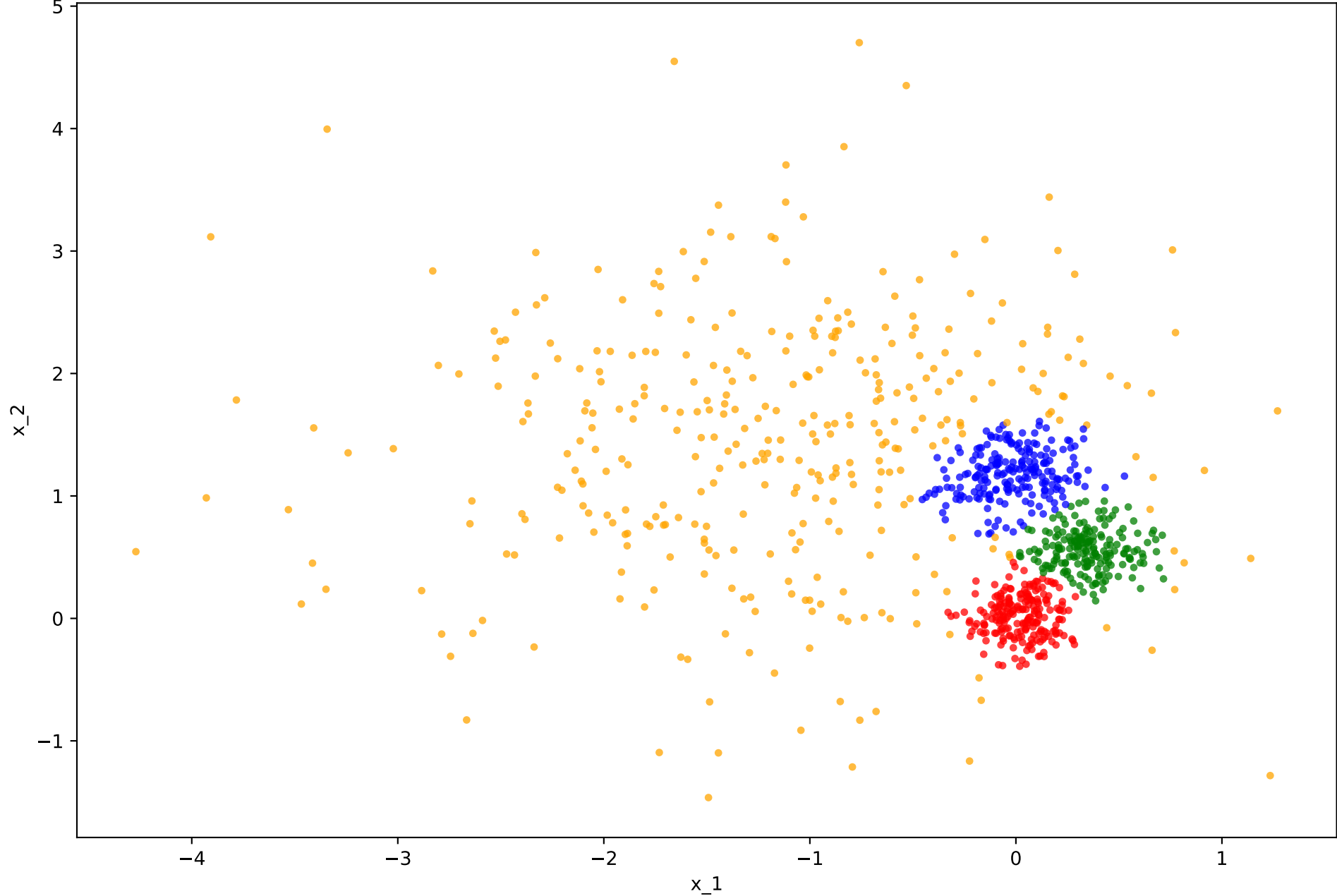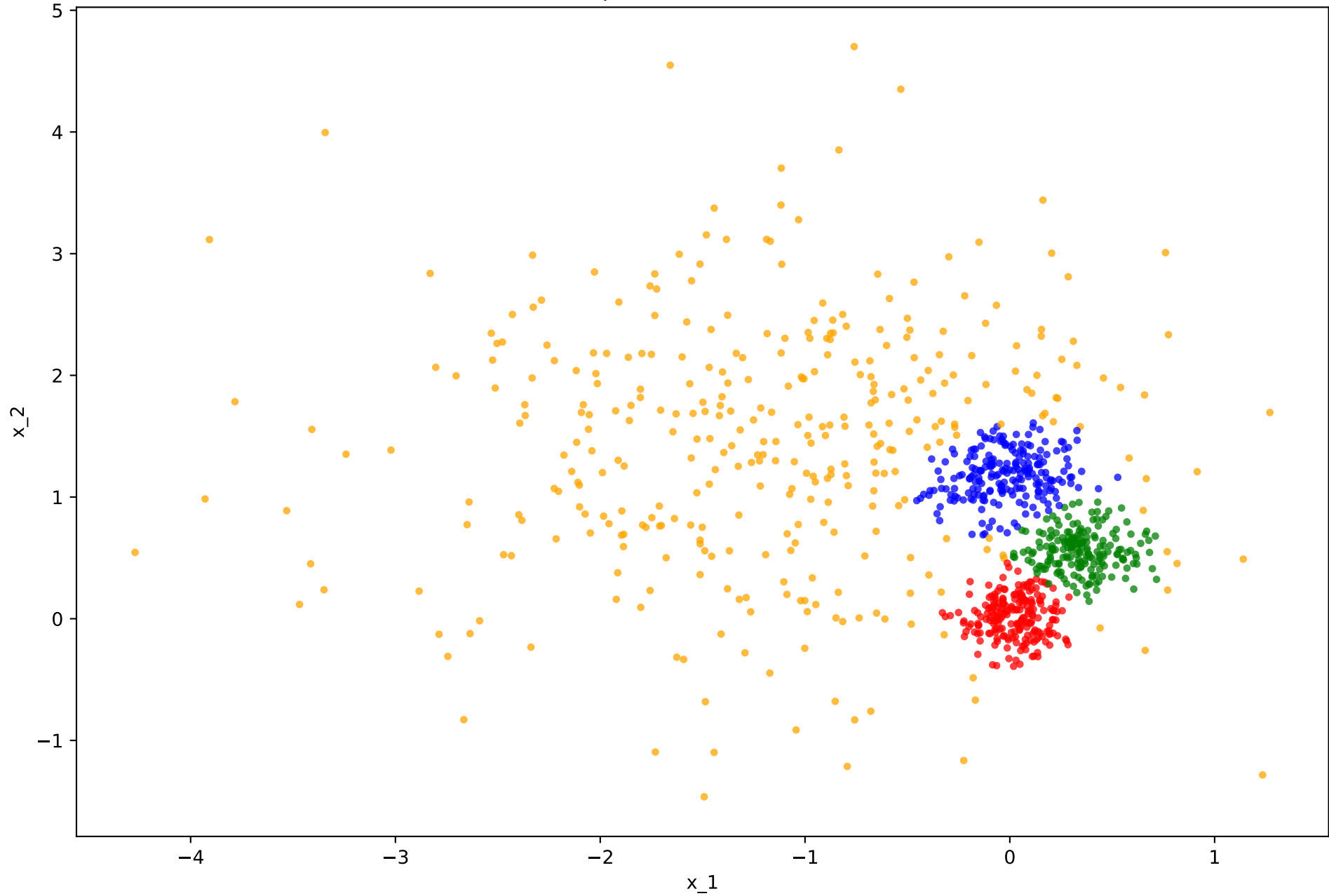
Unsupervised GMM Predictions

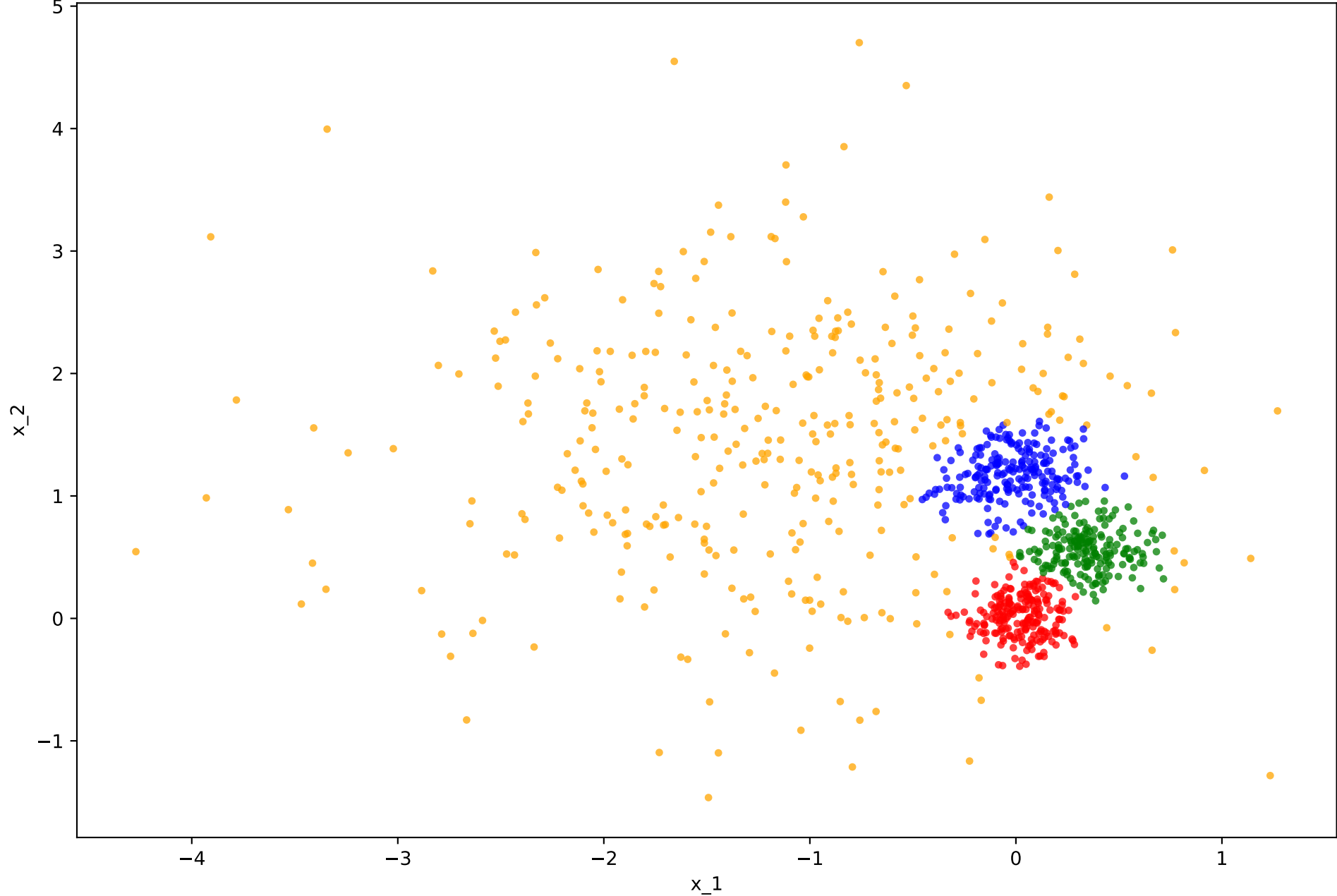Unsupervised GMM Predictions

Semi-supervised GMM Predictions

Semi-supervised GMM Predictions

CS 229 HW 3
Mengxuan Jin
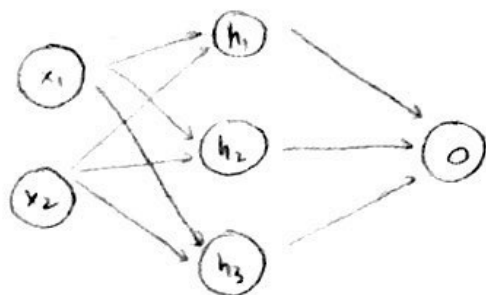
## Problem 1

m: sample size.    n = 2 feature



$$W_{1,1}^{[1]}, \quad W_{2,1}^{[1]} \qquad\qquad W_{0,1}^{[1]}$$
$$W_{1,2}^{[1]}, \quad W_{2,2}^{[1]} \qquad\qquad W_{0,2}^{[1]}$$
$$W_{1,3}^{[1]}, \quad W_{2,3}^{[1]} \qquad\qquad W_{0,3}^{[1]}$$
$$w_1^{[2]} \; w_2^{[2]} \; w_3^{[2]} \qquad\qquad W_0^{[2]}$$

(a)
$$z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \end{bmatrix} = \begin{bmatrix} W_{1,1}^{[1]} & W_{2,1}^{[1]} \\ W_{1,2}^{[1]} & W_{2,2}^{[1]} \\ W_{1,3}^{[1]} & W_{2,3}^{[1]} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} W_{0,1}^{[1]} \\ W_{0,2}^{[1]} \\ W_{0,3}^{[1]} \end{bmatrix}, \quad \sigma(z^{[1]}) = \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix}$$

$$z^{[2]} = \begin{bmatrix} w_1^{[2]} & w_2^{[2]} & w_3^{[2]} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} + W_0^{[2]}, \quad \sigma(z^{[2]}) = O_{\text{utput}}$$

then
$$z_2^{[1](i)} = W_{1,2}^{[1]} x_1^{(i)} + W_{2,2}^{[1]} x_2^{(i)} + W_{0,2}^{[1]}$$
$$h_2^{(i)} = \sigma(z_2^{[1](i)})$$
$$z^{[2](i)} = w_1^{[2]} h_1^{(i)} + w_2^{[2]} h_2^{(i)} + w_3^{[2]} h_3^{(i)} + W_0^{[2]}$$
$$\text{Output}^{(i)} = \sigma(z^{[2](i)})$$

$$\frac{\partial \ell}{\partial W_{1,2}^{[1]}} = \frac{\partial \ell}{\partial O^{(i)}} \cdot \frac{\partial O^{(i)}}{\partial z^{[2](i)}} \cdot \frac{\partial z^{[2](i)}}{\partial h_2^{(i)}} \cdot \frac{\partial h_2^{(i)}}{\partial z_2^{[1](i)}} \cdot \frac{\partial z_2^{[1]}}{\partial W_{1,2}^{[1]}}$$

$$= \left( \frac{1}{m} \sum_{i=1}^{m} 2 O^{(i)} \right) \cdot O^{(i)} (1 - O^{(i)}) \cdot w_2^{[2]} \cdot h_2^{(i)} (1 - h_2^{(i)}) \cdot x_1^{(i)}$$

$$W_{1,2}^{[1]} = W_{1,2}^{[1]} - \alpha \frac{\partial \ell}{\partial W_{1,2}^{[1]}}, \quad \text{where } \frac{\partial \ell}{\partial W_{1,2}^{[1]}} \text{ is given above}$$

(b)



$$f(z) = \begin{cases} 1, & z \geq 0 \\ 0, & z < 0 \end{cases} \qquad z = w^T x + b$$

$$z^{[1]} = \begin{bmatrix} W_{0,1}^{[1]} & W_{1,1}^{[1]} & W_{2,1}^{[1]} \\ W_{0,2}^{[1]} & W_{1,2}^{[1]} & W_{2,1}^{[1]} \\ W_{0,3}^{[1]} & W_{1,3}^{[1]} & W_{2,3}^{[1]} \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

we can construct this matrix s.t.

$f(z^{[1]}) = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$ for $x^{(i)}$ within triangle.

$(0.5, 0.5)$ $(3.5, 0.5)$ $(0.5, 3.5)$ are critical points

$$\begin{bmatrix} -1 & 2 & 0 \\ -1 & 0 & 2 \\ 4 & -1 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

$4 = 3.5 + 0.5$ ; $-1, -1 = $ slope

$2 = \frac{1}{0.5}$

then for the second matrix we can write as

$$z^{[2]} = \begin{bmatrix} w_0^{[2]} & w_1^{[2]} & w_2^{[2]} & w_3^{[2]} \end{bmatrix} \begin{bmatrix} 1 \\ \vdots \\ 1 \\ \vdots \end{bmatrix}$$, want $f(z^{[2]})$ be negative for each entry, but nonnega

if any of the entry in last 3 is $-1$.

choose $\begin{bmatrix} w_0^{[2]} & w_1^{[2]} & w_2^{[2]} & w_3^{[2]} \end{bmatrix} = \begin{bmatrix} 2.5 & -1 & -1 & -1 \end{bmatrix}$ satisfy the conditions desired.

(c) It's not possible. Because data set B not linearly separable. if using $f(x) = X$ on $h_1, h_2, h_3$, it's just linear regression on $x^{(i)}$ and so con't seperate the two classes. Then no matter what activation fn used in o, we can't achieve 100% accuracy

## Problem 2

(a) $\quad D_{KL}(P \| Q) = \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)} = \sum_{x \in X} P(x) \log \left(\frac{Q(x)}{P(x)}\right)^{-1} = -\sum_{x \in X} P(x) \log \left(\frac{Q(x)}{P(x)}\right)$

By def of Expected value: $\quad = -E\left(\log \frac{Q(x)}{P(x)}\right) = E\left(-\log \frac{Q(x)}{P(x)}\right)$

- log function is strictly convex:

$$\geq \log E\left(\frac{Q(x)}{P(x)}\right) = -\log \sum_{x \in X} P(x) \frac{Q(x)}{P(x)} = -\log \sum_{x \in X} Q(x) = -\log 1 = 0$$

$\therefore D_{KL}(P \| Q) \geq 0$

~~$\frac{Q(x)}{P(x)} = \log E\left(\frac{Q(x)}{P(x)}\right)$~~ $\quad \frac{Q(x)}{P(x)}$ is constant by ~~Jensen's inequality~~

- if $D_{KL}(P \| Q) = 0$, then the equality holds above, which means

$$\frac{Q(x)}{P(x)} = E\left(\frac{Q(x)}{P(x)}\right) = \sum_x P(x) \frac{Q(x)}{P(x)} = \sum_x Q(x) = 1 \implies Q(x) = P(x)$$

- If $P = Q$, then $D_{KL}(P \| Q) = \sum_{x \in X} P(x) \cdot \log 1 = 0$

(b) $\quad D_{KL}(P(X|Y) \| Q(X|Y)) = \sum_y P(y) \left( \sum_x P(x|y) \log \frac{P(x|y)}{Q(x|y)} \right)$

$\quad D_{KL}(P(X) \| Q(X)) = \sum_x P(x) \log \frac{P(x)}{Q(x)} \quad \underbrace{}_{Eq①}$

$\quad D_{KL}(P(Y|X) \| Q(Y|X)) = \sum_x P(x) \left( \sum_y P(y|x) \log \frac{P(y|x)}{Q(y|x)} \right) \quad \underbrace{}_{Eq②}$

$D_{KL}(P(X,Y) \| Q(X,Y)) = \sum_{x,y} P(x,y) \log \frac{P(x,y)}{Q(x,y)} = \sum_{x,y} P(x,y) \log \frac{P(x)P(y|x)}{Q(x)Q(y|x)}$

$= \sum_{x,y} P(x,y) \log \frac{P(x)}{Q(x)} + \sum_{x,y} P(x,y) \log \frac{P(y|x)}{Q(y|x)} = \sum_{x,y} P(x)P(y|x) \log \frac{P(x)}{Q(x)} + \sum_{x,y} P(x)P(y|x) \log \frac{P(y|x)}{Q(y|x)}$

$= \sum_x P(x) \log \frac{P(x)}{Q(x)} + \sum_x P(x) \sum_y P(y|x) \log \frac{P(y|x)}{Q(y|x)}$

$= Eq① + Eq②$

$\therefore D_{KL}(P(X,Y) \| Q(X,Y)) = D_{KL}(P(X) \| Q(X)) + D_{KL}(P(Y|X) \| Q(Y|X))$

(c) $\quad D_{KL}(\hat{P} \| P_\theta) = \sum_x \hat{P}(x) \log \frac{\hat{P}(x)}{P_\theta(x)} = \sum_x \frac{1}{m} \sum_{i=1}^{m} 1\{x^{(i)}=x\} \log \frac{\frac{1}{m}\sum_{i=1}^{m} 1\{x^{(i)}=x\}}{P_\theta(x)}$

$= \frac{1}{m} \sum_{i=1}^{m} \sum_x 1\{x^{(i)}=x\} \log \frac{\frac{1}{m}\sum_{i=1}^{m}1}{P_\theta(x^{(i)})} = \frac{1}{m} \sum_{i=1}^{m} \log \frac{1}{P_\theta(x^{(i)})} = -\frac{1}{m} \sum_{i=1}^{m} \log P_\theta(x^{(i)})$

so $\arg\min\limits_\theta D_{KL}(\hat{P} \| P_\theta) = \arg\min\limits_\theta -\frac{1}{m} \sum_{i=1}^{m} \log P_\theta(x^{(i)})$

$= \arg\max\limits_\theta \frac{1}{m} \sum_{i=1}^{m} \log P_\theta(x^{(i)})$

$= \arg\max\limits_\theta \sum_{i=1}^{m} \log P_\theta(x^{(i)})$

Problem 3

(a) $E_{y \sim p(y;\theta)} \left[ \nabla_{\theta'} \log p(y;\theta') |_{\theta'=\theta} \right]$

$= \int_{-\infty}^{\infty} p(y;\theta) \nabla_\theta \log p(y;\theta) \, dy$

$= \int_{-\infty}^{\infty} p(y;\theta) \frac{1}{p(y;\theta)} \nabla_\theta p(y;\theta) \, dy$

$= \nabla_\theta \int_{-\infty}^{\infty} p(y;\theta) \, dy = \nabla_\theta 1 = 0$

1) By definition of Covariance and covariance matrix

then $\text{Cov}_{y\sim p(y;\theta)}$ is a covariance matrix

$\cancel{\text{cov}(x,x) = \text{Var}(x) = E(x^2) - E(x)^2}$

$\text{Cov}_{y\sim p(y;\theta)}\left[\nabla_{\theta'}\log P(y;\theta')\big|_{\theta'=\theta}\right] = E\left(\nabla_\theta\log P(y;\theta')\,\nabla_\theta\log P(y;\theta')^T\big|_{\theta'=\theta}\right) -$

$\underbrace{E\left(\nabla_\theta\log P(y;\theta')\big|_{\theta'=\theta}\right)}_{=0 \text{ by part b)}}\underbrace{E\left(\nabla_\theta\log P(y;\theta'\big|_{\theta'=\theta}\right)^T}_{0}$

$\therefore \text{Cov}_{y\sim p(y;\theta)}\left[\nabla_\theta\log P(y;\theta')\big|_{\theta'=\theta}\right] = E\left(\nabla_\theta\log P(y;\theta')\,\nabla_\theta\log P(y;\theta')^T\big|_{\theta'=\theta}\right)$

c) $\nabla_\theta^2\log P(y;\theta')\big|_{\theta'=\theta} = H_{\log P(y;\theta)} = J\left(\dfrac{\nabla P(y;\theta)}{P(y;\theta)}\right)$, where $J$ is Jacobian operator

$= -P(y;\theta)^{-2}\cdot\nabla_\theta P(y;\theta)\,\nabla_\theta P(y;\theta)^T - \dfrac{\nabla_\theta^2 P(y;\theta)}{P(y;\theta)}$

$= \dfrac{\nabla_\theta^2 P(y;\theta)}{P(y;\theta)} - \dfrac{\nabla_\theta P(y;\theta)\,\nabla_\theta P(y;\theta)^T}{P(y;\theta)\,P(y;\theta)}$

$= \dfrac{H_{P(y;\theta)}}{P(y;\theta)} - \dfrac{\nabla P(y;\theta)}{P(y;\theta)}\left(\dfrac{\nabla_\theta P(y;\theta)}{P(y;\theta)}\right)^T$

then $E\left(\nabla_\theta^2\log P(y;\theta)\right) = E\left[\dfrac{H_{P(y;\theta)}}{P(y;\theta)} - \dfrac{\nabla_\theta P(y;\theta)}{P(y;\theta)}\left(\dfrac{\nabla_\theta P(y;\theta)}{P(y;\theta)}\right)^T\right]$

$= E\left[\dfrac{H_{P(y;\theta)}}{P(y;\theta)}\right] - E\left[\dfrac{\nabla_\theta P(y;\theta)}{P(y;\theta)}\left(\dfrac{\nabla_\theta P(y;\theta)}{P(y;\theta)}\right)^T\right]$

$= \int P(y;\theta)\dfrac{H_{P(y;\theta)}}{P(y;\theta)}\,dy - \underbrace{E\left[\nabla_\theta\log P(y;\theta)\,\nabla_\theta\log P(y;\theta)^T\right]}_{I(\theta)}$

$= \int H_{P(y;\theta)}\,dy - I(\theta)$

$= H_{\int P(y;\theta)dy} - I(\theta) = 0 - I(\theta)$

$\therefore E\left(\nabla_\theta^2\log P(y;\theta)\right) = -I(\theta) \Rightarrow E\left(-\nabla_\theta^2\log P(y;\theta)\right) = I(\theta)$

(d) $D_{KL}(P_\theta \| P_{\theta+d}) = \int_{-\infty}^{\infty} P_\theta(x) \log\frac{P_\theta(x)}{P_{\theta+d}(x)} dx = -\int_{-\infty}^{\infty} P_\theta(x) \log\frac{P_{\theta+d}(x)}{P_\theta(x)} dx$

$\approx -\int_{-\infty}^{\infty} P_\theta(x) \left( \log\frac{P_\theta(x)}{P_\theta(x)}^{\cancel{=0}} + d^T \nabla_\theta \log P_\theta(x) + \frac{1}{2} d^T \nabla_\theta^2 \log P_\theta(x) d \right) dx$

$= -\int_{-\infty}^{\infty} P_\theta(x) \left( d^T \nabla_\theta \log P_\theta(x) + \frac{1}{2} d^T \nabla_\theta^2 \log P_\theta(x) d \right) dx$

$= -\int_{-\infty}^{\infty} P_\theta(x) \left( \frac{\nabla_\theta P_\theta(x)}{P_\theta(x)} \right)^T d \, dx - \frac{1}{2} d^T \int_{-\infty}^{\infty} P_\theta(x) \frac{P_\theta(x) \nabla_\theta^2 P_\theta(x) - (\nabla_\theta P_\theta(x))(\nabla_\theta P_\theta(x))^T}{P_\theta(x)^2} d \, dx$

$= -\int_{-\infty}^{\infty} \nabla_\theta P_\theta(x)^T d \, dx - \frac{1}{2} d^T \int_{-\infty}^{\infty} \nabla_\theta^2 P_\theta(x) d \, dx + \frac{1}{2} d^T \int_{-\infty}^{\infty} P_\theta(x) \left( \frac{\nabla_\theta P_\theta(x)}{P_\theta(x)} \right)\left( \frac{\nabla_\theta P_\theta(x)}{P_\theta(x)} \right)^T d \, dx$

$= -\left( \frac{d}{d\theta} \underbrace{\int_{-\infty}^{\infty} P_\theta(x) dx}_{1} \right)^T d - \frac{1}{2} d^T \left( \nabla_\theta^2 \underbrace{\int_{-\infty}^{\infty} P_\theta(x) dx}_{1} \right) d + \frac{1}{2} d^T \left( \int_{-\infty}^{\infty} P_\theta(x) (\nabla_\theta \log P_\theta(x))(\nabla_\theta \log P_\theta(x))^T dx \right) d$

$= -0 - 0 + \frac{1}{2} d^T \underbrace{\left( \int_{-\infty}^{\infty} P_\theta(x) (\nabla_\theta \log P_\theta(x))(\nabla_\theta \log P_\theta(x))^T dx \right)}_{E(\nabla_\theta \log P_\theta(x) \, \nabla_\theta \log P_\theta(x)^T) = I(\theta)} d$

$= \frac{1}{2} d^T I(\theta) d$


(e) $d^* = \arg\max_d L(\theta+d) \quad \text{c.t.} \quad D_{KL}(P_\theta \| P_{\theta+d}) = c$

$\Rightarrow L(d,\lambda) = L(\theta+d) - \lambda( D_{KL}(P_\theta \| P_{\theta+d}) - c)$

$\approx L(\theta) + d^T \nabla_\theta L(\theta')\big|_{\theta'=\theta} - \frac{1}{2}\lambda d^T I(\theta) d + \lambda c$

$\nabla_d L(d,\lambda) = \nabla_\theta L(\theta')\big|_{\theta'=\theta} - \lambda I(\theta) d = 0$

$\lambda I(\theta) d = \nabla_\theta L(\theta')\big|_{\theta'=\theta}$

$d = \frac{1}{\lambda} (I(\theta))^{-1} \nabla_\theta L(\theta)$

## Problem 4

(a) $l_{sup}(\theta) = \sum_{i=1}^{\hat{m}} \log P(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta)$

$l_{unsup}(\theta) = \sum_{i=1}^{m} \log P(x^{(i)}; \theta) = \sum_{i=1}^{m} \log \sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta)$

let $\mathcal{L}$ be $\sum_{i=1}^{m} \left( \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i^{(t)}(z^{(i)})} + \propto \left( \sum_{i=1}^{\hat{m}} \log P(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta) \right) \right)$

then $\mathcal{L}(\theta^{(t+1)}) \geq \underbrace{\sum_{i=1}^{m} \left( \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})} \right)}_{\text{part ①}} + \underbrace{\propto \sum_{i=1}^{\hat{m}} \log P(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta^{(t)})}_{\text{part ②}}$

$\downarrow$ because $\theta^{(t+1)}$ maximize $\mathcal{L}$, then $\mathcal{L}(\theta^{(t+1)}) \geq \mathcal{L}(\theta) \ \forall \theta$, including $\theta^{(t)}$

then part ① is just $\sum_{i=1}^{m} E_{z^{(i)} \sim Q_i} \log \frac{P(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})} \geq \sum_{i=1}^{m} \log E_{z^{(i)} \sim Q_i} \frac{P(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})}$

By Jensen's inequality $= \sum_{i=1}^{m} \log \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \frac{P(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})}$

$= \sum_{i=1}^{m} \log P(x^{(i)}; \theta^{(t)}) = l_{unsup}(\theta^{(t)})$

and part ② is $l_{sup}(\theta^{(t)})$, so $\mathcal{L}(\theta^{(t+1)}) \geq l_{sup}(\theta^{(t)}) + \propto l_{unsup}(\theta^{(t)}) = \mathcal{L}(\theta^{(t)})$

(i.e. $l_{semi-sup}(\theta^{(t+1)}) \geq l_{semi-sup}(\theta^{(t)})$ )

(b) $\mu, \Sigma, \phi$

$$w_j^{(i)} = P(z^{(i)} = j \mid x^{(i)}; \phi, \mu, \Sigma) = \frac{P(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)}{\sum_{l=1}^{k} P(x^{(i)} \mid z^{(i)} = l) P(z^{(i)} = l)}$$

$$= \frac{\frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\right) \cdot \phi_j}{\sum_{l=1}^{k} \frac{1}{(2\pi)^{n/2} |\Sigma_l|^{1/2}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_l)^T \Sigma_l^{-1} (x^{(i)} - \mu_l)\right) \cdot \phi_l}$$

only $z^{(i)}$'s are the latent variables that need to be re-estimated (unlabeled)

(c) $\mathcal{l} = \sum_{i=1}^{m} \sum_{j=1}^{k} w_j^{(i)} \log \frac{\frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\right) \cdot \phi_j}{w_j^{(i)}} +$

$\alpha \sum_{i=1}^{\tilde{m}} \sum_{j=1}^{k} \log\left(\frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2}(\tilde{x}^{(i)} - \mu_j)^T \Sigma_j^{-1} (\tilde{x}^{(i)} - \mu_j)\right) \cdot \mathbb{1}\{\tilde{z}^{(i)} = j\} \cdot \phi_j\right)$

$\nabla_{\mu_l} \mathcal{l} = \underbrace{\sum_{i=1}^{m} w_l^{(i)} \left(\Sigma_l^{-1} x^{(i)} - \Sigma_l^{-1} \mu_l\right)}_{\text{from lecture note}} + \nabla_{\mu_l} \alpha \sum_{i=1}^{\tilde{m}} -\frac{1}{2}(\tilde{x}^{(i)} - \mu_l)^T \Sigma_j^{-1} (\tilde{x}^{(i)} - \mu_l) \cdot \mathbb{1}\{\tilde{z}^{(i)} = l\}$

$\Downarrow$

$\frac{1}{2} \alpha \sum_{i=1}^{\tilde{m}} \nabla_{\mu_l} \left(2\mu_l^T \Sigma_l^{-1} \tilde{x}^{(i)} - \mu_l^T \Sigma_l^{-1} \mu_l\right) \cdot \mathbb{1}\{\tilde{z}^{(i)} = l\}$

$= \alpha \sum_{i=1}^{\tilde{m}} \left(\Sigma_l^{-1} \tilde{x}^{(i)} - \Sigma_l^{-1} \mu_l\right) \cdot \mathbb{1}\{\tilde{z}^{(i)} = l\}$

then $\nabla_{\mu_l} \mathcal{l} = 0 \Rightarrow \sum_{i=1}^{m} w_l^{(i)} \left(\Sigma_l^{-1} x^{(i)} - \Sigma_l^{-1} \mu_l\right) + \alpha \sum_{i=1}^{\tilde{m}} \left(\Sigma_l^{-1} \tilde{x}^{(i)} - \Sigma_l^{-1} \mu_l\right) \cdot \mathbb{1}\{\tilde{z}^{(i)} = l\} = 0$

$$\mu_l = \frac{\sum_{i=1}^{m} w_l^{(i)} x^{(i)} + \alpha \sum_{i=1}^{\tilde{m}} \mathbb{1}\{\tilde{z}^{(i)} = l\} \tilde{x}^{(i)}}{\left(\sum_{i=1}^{m} w_l^{(i)}\right) + \alpha \sum_{i=1}^{\tilde{m}} \mathbb{1}\{\tilde{z}^{(i)} = l\}}$$

$$\nabla_{\Sigma_i} \mathcal{L} = -\frac{1}{2} \sum_{i=1}^{m} w_i^{(i)} \left( \Sigma_i^{-1} - \Sigma_i^{-1}(x^{(i)}-\mu_i)(x^{(i)}-\mu_i)^T \Sigma_i^{-1} \right) +$$

$$\nabla_{\Sigma_i} \propto \sum_{i=1}^{\tilde{m}} \underbrace{\left( \log \frac{1}{(2\pi)^{n/2}|\Sigma_i|^{1/2}} - \frac{1}{2}(\hat{x}^{(i)}-\mu_i)^T \Sigma_j^{-1}(\hat{x}^{(i)}-\mu_i) \right) \cdot \mathbb{I}\{\hat{z}^{(i)}=i\}}_{} = 0.$$

$$-\frac{1}{2} \propto \sum_{i=1}^{\tilde{m}} \left( \Sigma_i^{-1} - \Sigma_i^{-1}(\hat{x}^{(i)}-\mu_i)(\hat{x}^{(i)}-\mu_i)^T \Sigma_i^{-1} \right) \cdot \mathbb{I}\{\hat{z}^{(i)}=i\}$$

$$\therefore \nabla_\Sigma \mathcal{L} = +\frac{1}{2} \sum_{i=1}^{m} w_i^{(i)} \left( \Sigma_i - (x^{(i)}-\mu_i)(x^{(i)}-\mu_i)^T \right) + \frac{\alpha}{2} \sum_{i=1}^{\tilde{m}} \mathbb{I}\{\hat{z}^{(i)}=i\} \left( \Sigma_i - (\hat{x}^{(i)}-\mu_i)(\hat{x}^{(i)}-\mu_i)^T \right) = 0$$

$$\Sigma_i := \frac{\sum_{i=1}^{m} w_i^{(i)}(x^{(i)}-\mu_i)(x^{(i)}-\mu_i)^T + \alpha \sum_{i=1}^{\tilde{m}} (\hat{x}^{(i)}-\mu_i)(\hat{x}^{(i)}-\mu_i)^T \mathbb{I}\{\hat{z}^{(i)}=i\}}{\sum_{i=1}^{m} w_i^{(i)} + \alpha \sum_{i=1}^{\tilde{m}} \mathbb{I}\{\hat{z}^{(i)}=i\}}$$

for $\phi_i$

maximize $\quad (\because \sum \phi_j = 1)$

$$\mathcal{L}' = \sum_{i=1}^{m} \sum_{j=1}^{k} w_j^{(i)} \log \phi_j + \beta_1 \left( \sum_{j=1}^{k} \phi_j - 1 \right) + \alpha \left( \sum_{i=1}^{\tilde{m}} \sum_{j=1}^{k} \mathbb{I}\{\hat{z}^{(i)}=j\} \log \phi_j + \beta_2 \left( \sum_{j=1}^{k} \phi_j - 1 \right) \right)$$

$$\nabla_{\phi_i} \mathcal{L}' = \sum_{i=1}^{m} w_i^{(i)}/\phi_i + \beta_1 + \alpha \left( \sum_{i=1}^{\tilde{m}} \mathbb{I}\{\hat{z}^{(i)}=i\} \right)/\phi_i + \beta_2 = 0$$

$$\Rightarrow \phi_i = \frac{\sum_{i=1}^{m} w_i^{(i)} + \alpha \sum_{i=1}^{\tilde{m}} \mathbb{I}\{\hat{z}^{(i)}=i\}}{-\beta_1 - \alpha\beta_2} \quad \text{and} \quad -\beta_1 - \alpha\beta_2 = \sum_{i=1}^{m} \sum_{j=1}^{k} w_j^{(i)} + \alpha \sum_{i=1}^{\tilde{m}} \sum_{j=1}^{k} w_j^{(i)}$$

$$= m + \alpha\tilde{m}$$

$$\therefore \phi_i = \frac{\sum_{i=1}^{m} w_i^{(i)} + \alpha \sum_{i=1}^{\tilde{m}} \mathbb{I}\{\hat{z}^{(i)}=i\}}{m + \alpha\tilde{m}}$$

(f) [i] unsupervised EM takes $^{much}$more iterations to converge (over 100 iterations) than SS EM (SS EM takes about 20-30 iterations)

[ii] Based on the plot... SS EM is more stable since 3 plots are almost the same, but unsupervised EM's plots differ each time (the assignment is different, in particular, the boundaries on the left 2 groups changes each time)

[iii] SS EM has better quality from plots, it clearly has 3 Gaussian dist. w/ low variance and 1 w/ high variance

$$\nabla_\lambda L(d, \lambda) = -\frac{1}{2} d^\top I(\theta) d + c$$

with $d = \frac{1}{\lambda} (I(\theta))^{-1} \nabla_\theta l(\theta)$, $\nabla_\lambda L(d, \lambda) = -\frac{1}{2} (\frac{1}{\lambda} (\nabla_\theta l(\theta))^\top (I(\theta))^{-\top} (I(\theta))^{-1} (I(\theta))^{-1} \nabla_\theta l(\theta)) + c = 0$

$$\Rightarrow 2c = \frac{1}{\lambda^2} (\nabla_\theta l(\theta))^\top (I(\theta))^{-\top} \nabla_\theta l(\theta)$$

$$\lambda = \left( \frac{1}{2c} (\nabla_\theta l(\theta))^\top (I(\theta))^{-\top} \nabla_\theta l(\theta) \right)^{1/2}$$

then $d^* = \left( \frac{1}{2c} \nabla_\theta l(\theta))^\top (I(\theta))^{-\top} \nabla_\theta l(\theta) \right)^{-\frac{1}{2}} (I(\theta))^{-1} \nabla_\theta l(\theta)$

**(f)** Natural Gradient direction $(I(\theta))^{-1} \nabla_\theta l(\theta)$

Newton's method direction $-(H(\theta))^{-1} \nabla_\theta l(\theta)$

$$I(\theta) = \mathop{E}_{y \sim p(y;\theta)} [ -\nabla_\theta^2 \log P(y;\theta) ]$$

$H(\theta) = \frac{\partial^2}{\partial \theta^2} a(\theta^\top x) \, x \, x^\top$ (from HW1)

$\qquad = \text{Var}(y|x, \theta) \, x \, x^\top$

$\nabla_\theta^2 l(\theta)$, $l(\theta) = \log [p(y;\theta)]$

$$(I(\theta))^{-1} \nabla l(\theta) = \left( \mathop{E}_y [ -\nabla_\theta^2 \, l(\theta) ] \right)^{-1} \nabla_\theta l(\theta)$$

$\qquad = \left[ -\mathop{E}_y [ \text{Var}(y|x;\theta) x x^\top ) ] \right]^{-1} \nabla_\theta l(\theta)$ ← from PS1, question 4

$\qquad = -\left( \mathop{E}_y [\text{Var}(y|x;\theta)] \, x x^\top \right)^{-1} \nabla_\theta l(\theta)$

$\qquad = -\left( \text{Var}(y|x;\theta) \, x x^\top \right)^{-1} \nabla_\theta l(\theta)$

$\qquad = -\left( \nabla_\theta^2 l(\theta) \right)^{-1} \nabla_\theta l(\theta)$

$\qquad = -(H(\theta))^{-1} \nabla_\theta l(\theta)$

# Problem 5.

(b)    Compression factor = $\dfrac{\text{previous \# of bytes to store image}}{\text{new \# of bytes to store image}}$

Previous: $512 \times 512 \times 3 \times 8 = 6,291,456$ bits $= 786,432$ bytes

Compressed: $512 \times 512 \times 4 = 1,048,576$ bits $= 131,072$ bytes

~~the number "8" comes from $2^8 = 256$ so we have 256 colors, then each byte has 8 bits~~

~~the number "4" comes from $2^4 = 16$ since we now have 16 colors, then each byte has 4 bits~~

The original image has 24 bits/pixel, but now for 16 colors, we only need $\log_2 16 = 4$ bits per pixel

So compression rate $= \dfrac{786,432}{131,072} = \dfrac{24}{4} = 6$.