

1. (a)

I did this homework with Luning Zhao.

We work on our homework separately, but we discuss when meeting problems.

Homework is fine.

(b) I certify that all solutions are entirely in my words. and that I have not looked at another student's solutions. I have credited all external sources in this write up.

Siyao Jia

$$2. (a) \text{ when } k=n, \quad M_i = \frac{1}{|C_i|} \sum_{x_j \in C_i} x_j = \frac{1}{|C_i|} x_i = x_i$$

$$\sum_{i=1}^k \sum_{x_j \in C_i} \|x_j - M_i\|_2^2 = \sum_{i=1}^n (M_i - x_i)^T (M_i - x_i) = 0$$

So the minimum value of object is 0 when $k=n$.

$$(b) \quad f(M_i) = \lambda \|M_i\|^2 + \sum_{x_j \in C_i} \|x_j - M_i\|^2 \\ = \lambda M_i^T M_i + \sum (x_j - M_i)^T (x_j - M_i)$$

$$\frac{\partial f(M_i)}{\partial M_i} = 2\lambda M_i^T + \sum -2x_j^T + 2M_i^T = 2\lambda M_i^T + 2|C_i| M_i^T + \sum -2x_j^T$$

$$\text{Let } \frac{\partial f(M_i)}{\partial M_i} = 0 \Rightarrow M_i = \frac{\sum x_j}{\lambda + |C_i|}$$

(c) Assume the k meet locations are M_i where $i \in [1:k]$

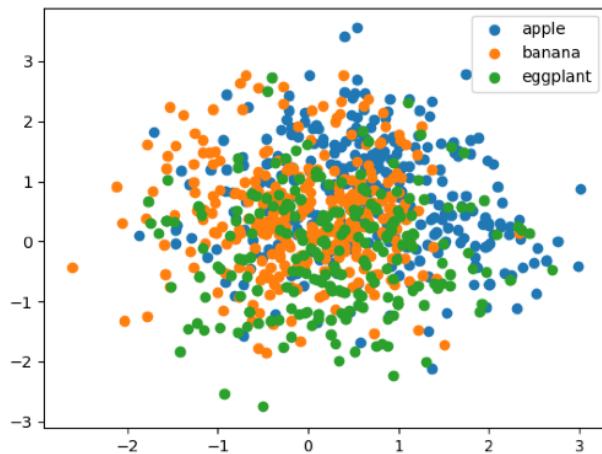
The distance for students near M_i to drive to M_i is $\sum_{x_j \in C_i} \|x_j - M_i\|^2$

The distance for vehicle to drive from M_i to $(0,0)$ is $\|M_i\|^2$.

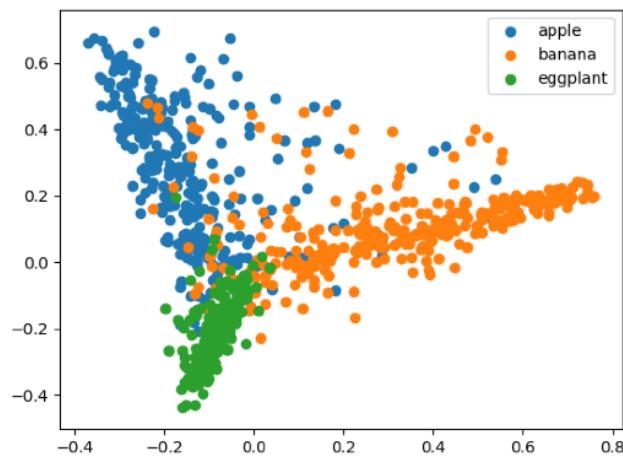
So the objective function is $\min_{C_1, \dots, C_k} \sum_{i=1}^k (\lambda \|M_i\|^2 + \sum_{x_j \in C_i} \|x_j - M_i\|^2)$

(d) Set $i(j) = \operatorname{argmin}_i \frac{1}{|C_i|} \sum_{x_j \in C_i} k(x_j, x_i)$

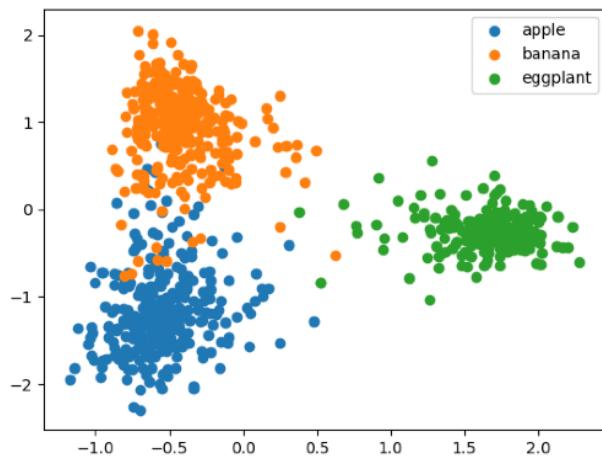
3 (a)



(b)



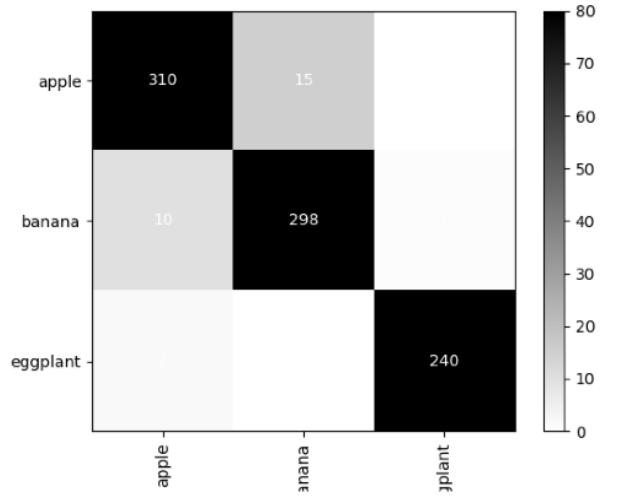
(c)



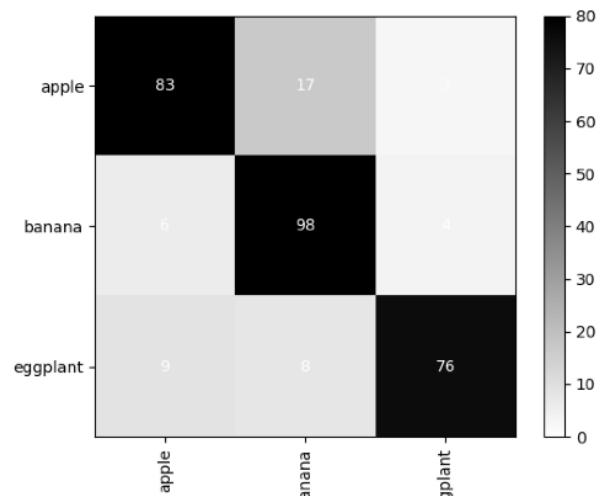
CCA is the best for separation among classes. Random is the worst. Because CCA finds the dimension that correlates most with the class.

(d)

training:

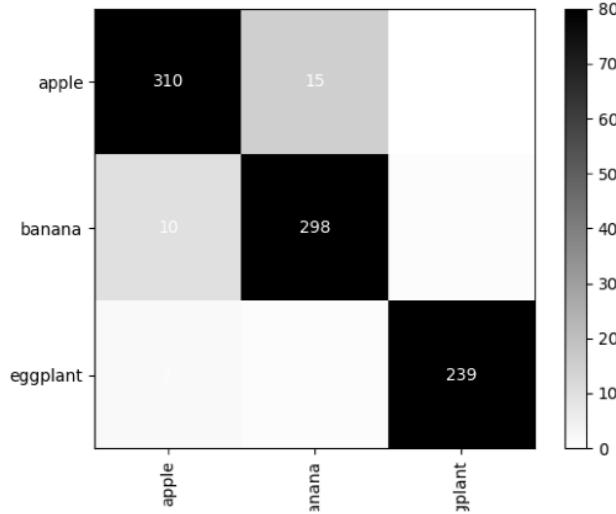


validation:

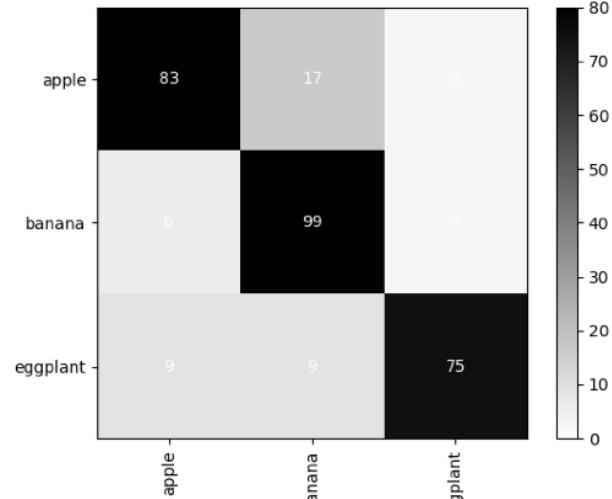


(e)

training:



validation:



```

class LDA_Model():

    def __init__(self, class_labels):
        #####SCALE AN IDENTITY MATRIX BY THIS TERM AND ADD TO COMPUTED COVARIANCE
        self.reg_cov = 0.001
        self.NUM_CLASSES = len(class_labels)

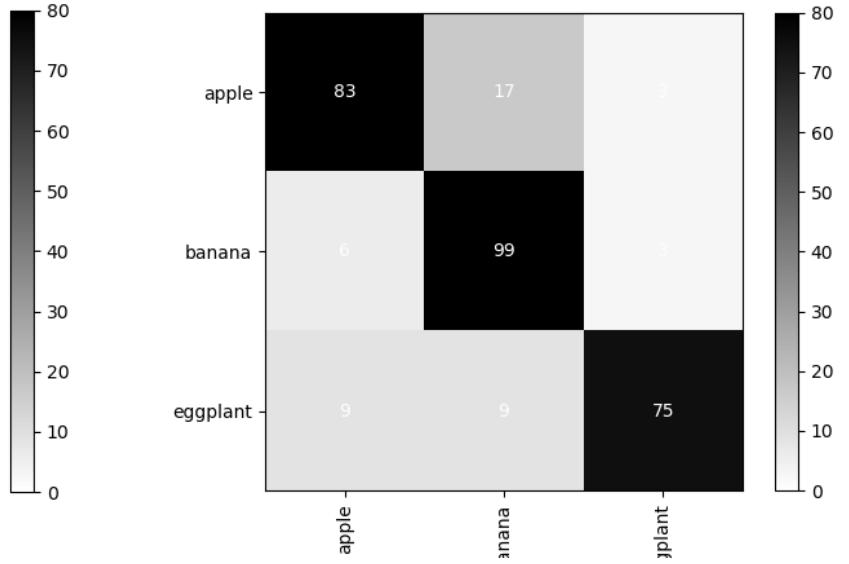
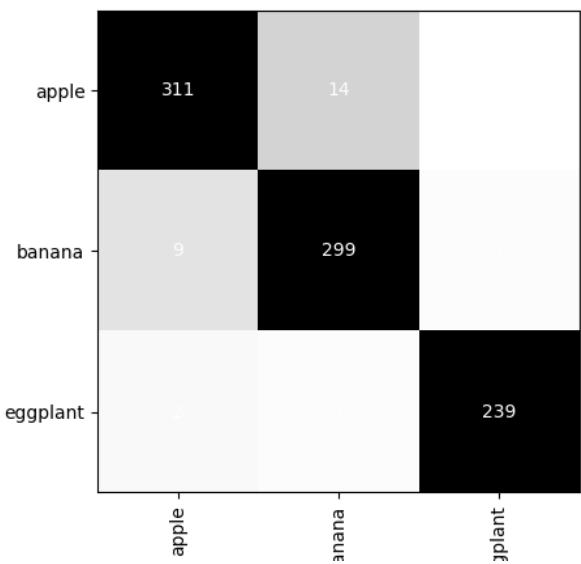
    def train_model(self, X, Y):
        """
        FILL IN CODE TO TRAIN MODEL
        MAKE SURE TO ADD HYPERPARAMTER TO MODEL
        """
        self.u = []
        for i in range(self.NUM_CLASSES):
            idx = np.where(np.array(Y)==i)[0]
            self.u.append(np.mean(np.array(X)[idx], axis=0))
        self.sig = np.cov(np.array(X).T)
        self.sig += np.identity(len(self.sig))*self.reg_cov

    def eval(self, x):
        """
        Fill in code to evaluate model and return a prediction
        Prediction should be an integer specifying a class
        """
        y = []
        for i in range(self.NUM_CLASSES):
            y.append((x-self.u[i]) @ inv(self.sig) @ (x-self.u[i]).T)
        pdb.set_trace()
        return np.argmin(y)
    
```

(f)

training:

validation:



```

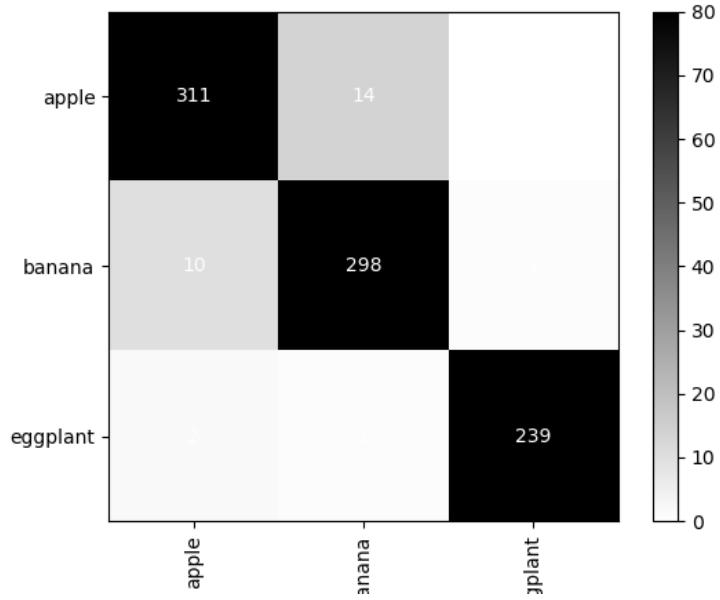
class QDA_Model():
    def __init__(self, class_labels):
        #####SCALE AN IDENTITY MATRIX BY THIS TERM AND ADD TO COMPUTED COVARIANCE MATRIX TO PREVENT IT BEING SINGULAR
        self.reg_cov = 0.01
        self.NUM_CLASSES = len(class_labels)

    def train_model(self, X, Y):
        """
        FILL IN CODE TO TRAIN MODEL
        MAKE SURE TO ADD HYPERPARAMETER TO MODEL
        """
        self.u = []
        self.sig = []
        for i in range(self.NUM_CLASSES):
            idx = np.where(np.array(Y)==i)[0]
            self.u.append(np.mean(np.array(X)[idx], axis=0))
            sig = np.cov(np.array(X)[idx].T)
            sig += np.identity(len(sig))*self.reg_cov
            self.sig.append(sig)

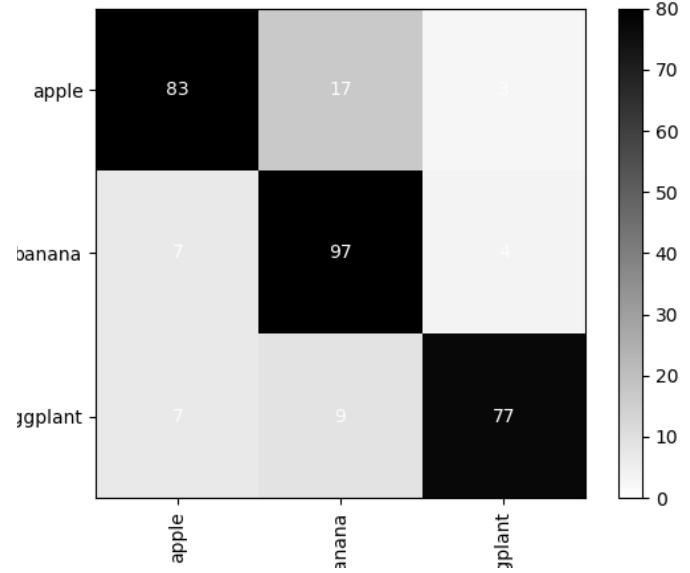
    def eval(self, x):
        """
        Fill in code to evaluate model and return a prediction
        Prediction should be an integer specifying a class
        """
        y = []
        for i in range(self.NUM_CLASSES):
            y.append((x-self.u[i]) @ inv(self.sig[i]) @ (x-self.u[i]).T + np.log(det(self.sig[i])))
        return np.argmax(y)
    
```

(g)

training:

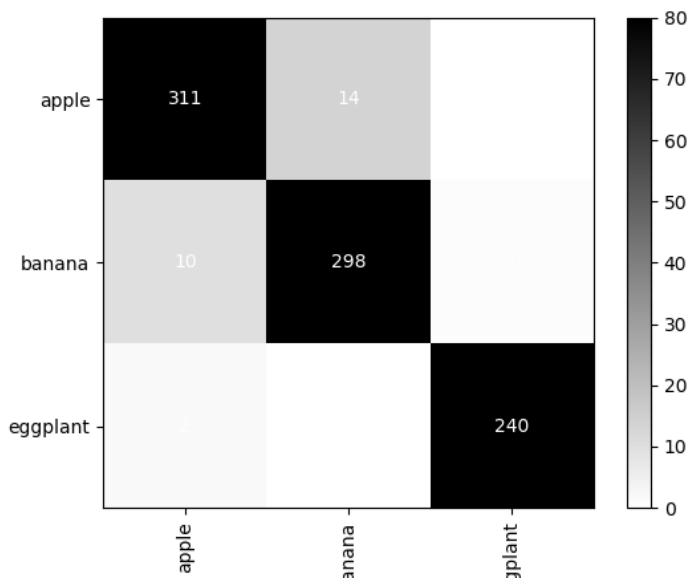


validation:

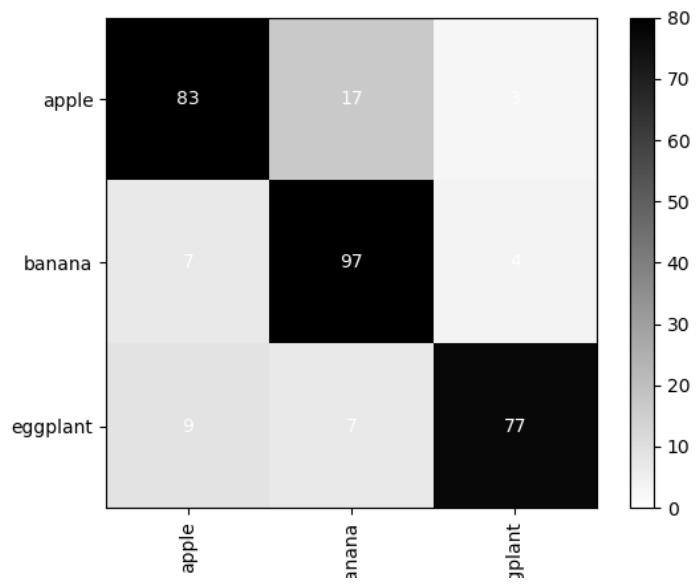


(h)

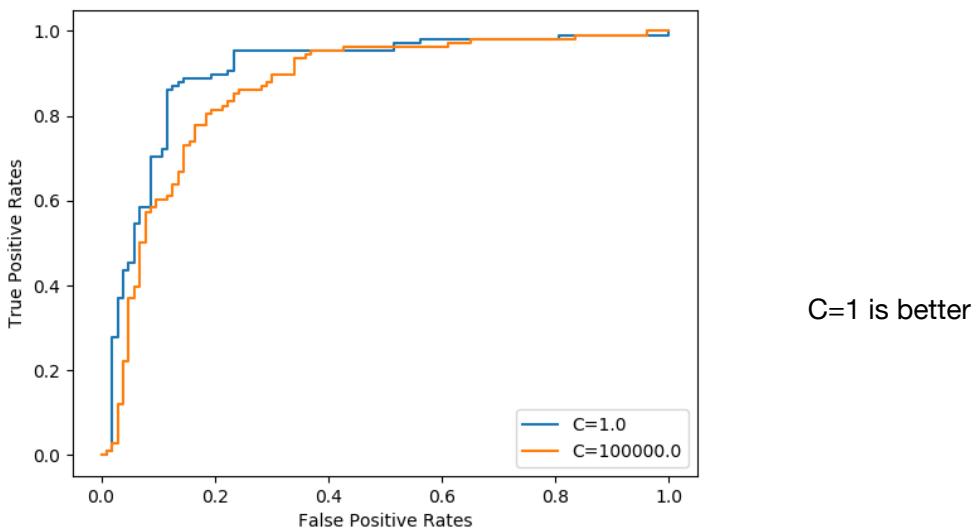
training:



validation:



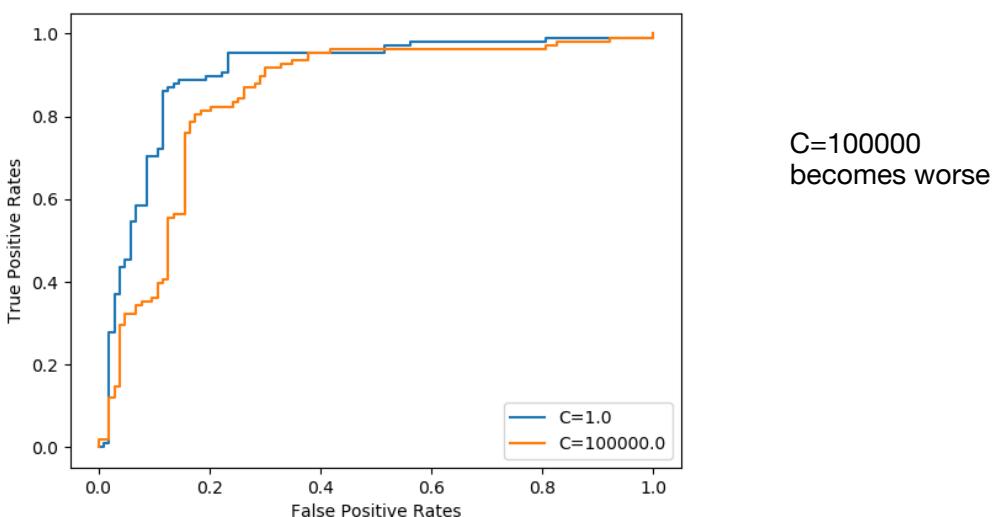
(i)



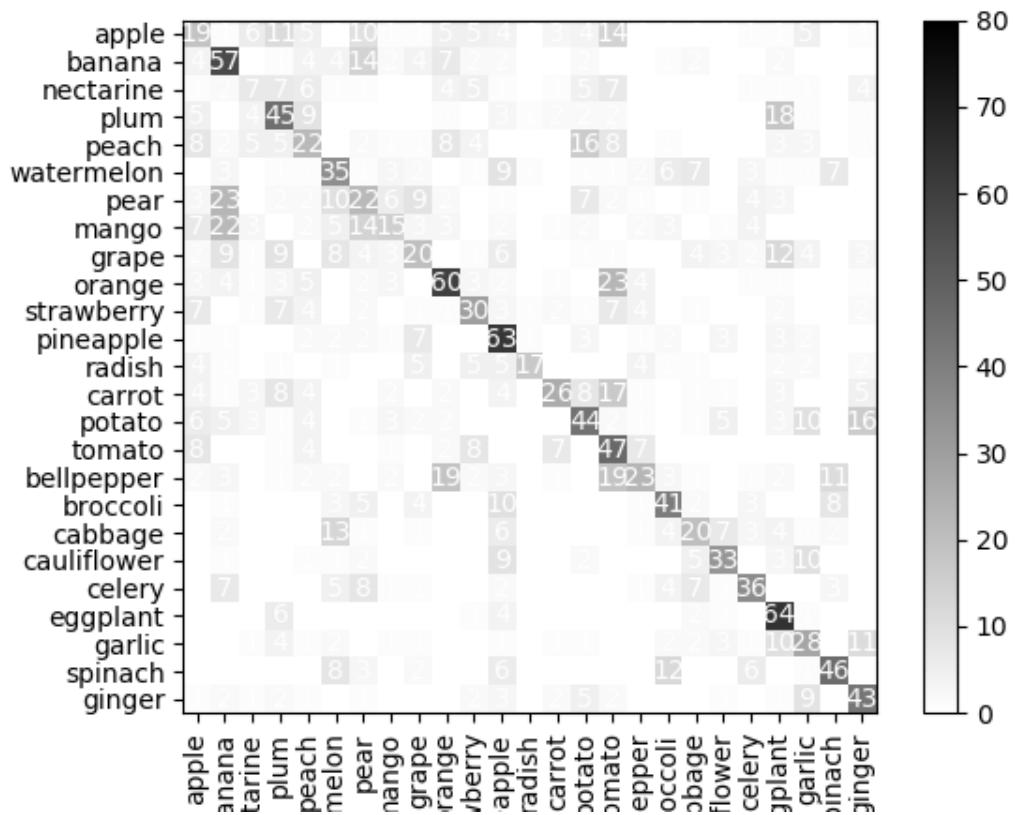
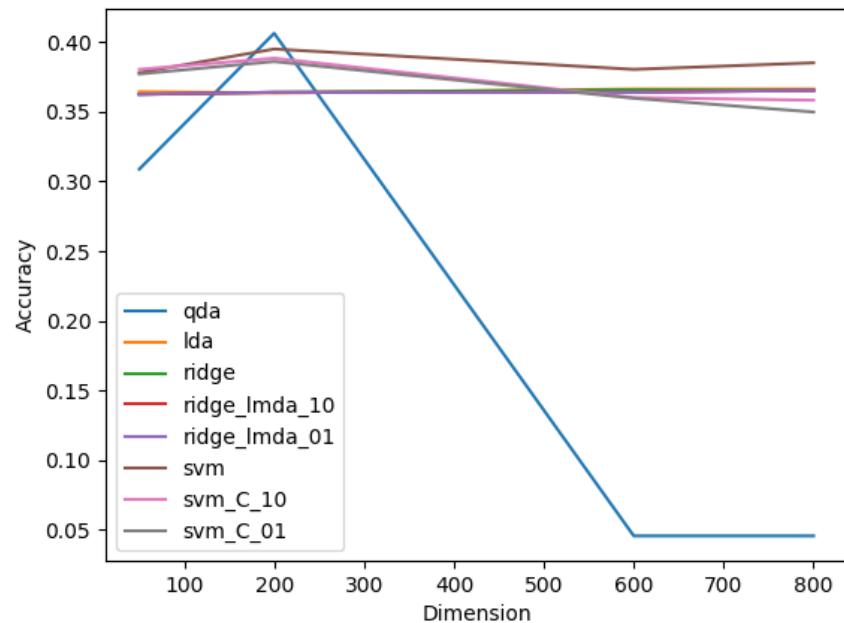
```
def ROC(scores, labels):
    thresholds = sorted(np.unique(scores))
    thresholds = [-float("Inf")] + thresholds + [float("Inf")]
    tps = []
    fps = []

    # student code start here
    # TODO: Your code
    # student code end here
    for thres in thresholds:
        tpr, fpr = compute_tp_fp(thres, scores, labels)
        tps.append(tpr)
        fps.append(fpr)
    return tps, fps
```

(j)



(k)



$$\begin{aligned}
4.(a) \quad & P(X=x, Z=1|\theta) = \frac{1}{2\pi\sigma_1} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} \times 0.5 \\
& P(X=x, Z=2|\theta) = \frac{1}{2\pi\sigma_2} e^{-\frac{(x-\mu_2)^2}{2\sigma_2^2}} \times 0.5 \\
& P(X=x; \theta) = \frac{1}{2\pi\sigma_1} \left(\frac{1}{\sigma_1} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} + \frac{1}{\sigma_2} e^{-\frac{(x-\mu_2)^2}{2\sigma_2^2}} \right) \\
& \ell(X=x; \theta) = -\log \sigma_1 \cdot \frac{(x-\mu_1)^2}{2\sigma_1^2} - \log \sigma_2 \cdot \frac{(x-\mu_2)^2}{2\sigma_2^2} - \log(2\pi) \\
(b) \quad & \ell(X_1=x_1, \dots, X_n=x_n) = \sum_{i=1}^n \ell(X_i=x_i) \\
& = \sum_{i=1}^n \left[-\log \sigma_1 \cdot \frac{(x_i-\mu_1)^2}{2\sigma_1^2} - \log \sigma_2 \cdot \frac{(x_i-\mu_2)^2}{2\sigma_2^2} \right] + c.
\end{aligned}$$

$$\begin{aligned}
(c) \quad & \ell(x_i; \theta) = -\log P(x_i; \theta) \\
& = -\log \left(\sum_{k=1}^2 P(x_i, k; \theta) \right) \\
& = -\log \left(\sum_{k=1}^2 \frac{P(x_i, k; \theta)}{q_i(k)} q_i(k) \right) \\
& \geq \sum_{k=1}^2 q_i(k) \log \frac{P(x_i, k; \theta)}{q_i(k)} \\
& = \sum_{k=1}^2 q_i(k) \log P(x_i, k; \theta) + \sum_{k=1}^2 q_i(k) \log \left(\frac{1}{q_i(k)} \right) = F_i(\theta; q_i)
\end{aligned}$$

$$\begin{aligned}
\ell(\{x_i\}; \theta) &= \sum_i \ell(x_i; \theta) \\
&\geq \sum_{i=1}^n F_i(\theta; q_i) = F(\theta; q)
\end{aligned}$$

$$\begin{aligned}
(d) \quad & q(Z_1=z_1, \dots, Z_n=z_n) = \prod_{i=1}^n P(Z=z_i | X=x_i; \theta^t) = \prod_{i=1}^n f_i(z_i; z_i) \\
& \Rightarrow q_i(z_i=z_i) = P(Z=z_i | X=x_i; \theta^t) \\
& F_i(\theta^t; q_i) = \sum_{k=1}^2 q_i(k) \log \frac{P(x_i, k; \theta^t)}{q_i(k)}, \\
& = \sum_{k=1}^2 q_i(k) \log \frac{P(x_i, k; \theta^t)}{P(Z=z_i | X=x_i, \theta^t)} \\
& = \sum_{k=1}^2 q_i(k) P(X=x_i | \theta^t) = \ell(x_i | \theta)
\end{aligned}$$

$$\zeta_0 \quad \ell(\{x_i\}; \theta^t) = \sum_i F_i(\theta^t; q_i) = F(\theta^t; q)$$

$\zeta_0 \quad q^{t+1}(Z=z_1, \dots, Z_n=z_n) = \prod_{i=1}^n P(Z=z_i | X=x_i; \theta^t)$ is a valid maximizer for E-step.

$$\begin{aligned}
 (e) \quad q_i^{t+1} (z_i=1) &= p(z=1 | x=x_i; \theta^t) \\
 &= \frac{p(z=1 | x=x_i; \theta^t)}{p(z=1 | x=x_i; \theta^t) + p(z=2 | x=x_i; \theta^t)} \\
 &= \frac{\frac{1}{\sigma_1} \exp(-\frac{(x_i - \mu_1)^2}{2\sigma_1^2})}{\frac{1}{\sigma_1} \exp(-\frac{(x_i - \mu_1)^2}{2\sigma_1^2}) + \frac{1}{\sigma_2} \exp(-\frac{(x_i - \mu_2)^2}{2\sigma_2^2})}
 \end{aligned}$$

Similarly $q_i^{t+1} (z_i=2) = \frac{\frac{1}{\sigma_2} \exp(-\frac{(x_i - \mu_2)^2}{2\sigma_2^2})}{\frac{1}{\sigma_1} \exp(-\frac{(x_i - \mu_1)^2}{2\sigma_1^2}) + \frac{1}{\sigma_2} \exp(-\frac{(x_i - \mu_2)^2}{2\sigma_2^2})}$

Intuitively, we want the distribution of $p(z=k)$ proportional to the probability $p(x|z, \theta^t)$. So this update makes sense

$$\begin{aligned}
 (f) \quad L &((\mu_1, \mu_2, \sigma_1, \sigma_2; \theta^t; q^{t+1})) \\
 &= \sum_{i=1}^n L(x_i; \theta, q_i^{t+1}) \\
 &= \sum_{i=1}^n \left(q_i^{t+1} \log p(x_i; k=1, \theta) + (1-q_i^{t+1}) \log p(x_i; k=2, \theta) \right) \\
 &= \sum_{i=1}^n \left(q_i^{t+1} \left(-\log \sigma_1 - \frac{(x_i - \mu_1)^2}{2\sigma_1^2} \right) + (1-q_i^{t+1}) \left(-\log \sigma_2 - \frac{(x_i - \mu_2)^2}{2\sigma_2^2} \right) \right) + C \\
 &= C - \sum_{i=1}^n \left[q_i^{t+1} \left(\frac{(x_i - \mu_1)^2}{2\sigma_1^2} + \log \sigma_1 \right) + (1-q_i^{t+1}) \left(\frac{(x_i - \mu_2)^2}{2\sigma_2^2} + \log \sigma_2 \right) \right]
 \end{aligned}$$

$$(g) \quad \frac{\partial L}{\partial \mu_1} = -\sum_{i=1}^n q_i^{t+1} \frac{2(x_i - \mu_1)(-1)}{2\sigma_1^2} = -\frac{\sum_{i=1}^n q_i^{t+1} (\mu_1 - x_i)}{\sigma_1^2}$$

Similarly, $\frac{\partial L}{\partial \mu_2} = -\frac{\sum_{i=1}^n (1-q_i^{t+1})(\mu_2 - x_i)}{\sigma_2^2}$

$$\frac{\partial L}{\partial \sigma_1} = -\sum_{i=1}^n q_i^{t+1} \left(\frac{(x_i - \mu_1)^2 \cdot (-2)}{2\sigma_1^3} + \frac{1}{\sigma_1} \right) = \frac{\sum_{i=1}^n q_i^{t+1} (x_i - \mu_1)^2}{\sigma_1^3} - \frac{\sum_{i=1}^n q_i^{t+1}}{\sigma_1}$$

Similarly, $\frac{\partial L}{\partial \sigma_2} = \frac{\sum_{i=1}^n (1-q_i^{t+1})(x_i - \mu_2)^2}{\sigma_2^3} - \frac{\sum_{i=1}^n (1-q_i^{t+1})}{\sigma_2}$

$$(b) \frac{\partial L}{\partial \mu_1} = -\frac{\sum_{i=1}^n q_i^{t+1} (\mu_1 - x_i)}{\sigma_1^2} = 0 \Rightarrow \sum_{i=1}^n q_i^{t+1} (\mu_1 - x_i) = 0$$

$$\Rightarrow \mu_1^{t+1} = \frac{\sum_{i=1}^n q_i^{t+1} x_i}{\sum_{i=1}^n q_i^{t+1}}$$

Similarly, $\frac{\partial L}{\partial \mu_2} = 0 \Rightarrow \mu_2^{t+1} = \frac{\sum_{i=1}^n (1-q_i^{t+1}) x_i}{\sum_{i=1}^n (1-q_i^{t+1})}$

$$\frac{\partial L}{\partial \sigma_1} = \frac{\sum_{i=1}^n q_i^{t+1} (x_i - \mu_1)^2}{\sigma_1^3} = \frac{\sum_{i=1}^n q_i^{t+1}}{\sigma_1} = 0 \Rightarrow \frac{\sum_{i=1}^n q_i^{t+1} (x_i - \mu_1)^2}{\sigma_1^2} = \frac{\sum_{i=1}^n q_i^{t+1}}{\sigma_1^2}$$

$$\Rightarrow (\sigma_1^2)^{t+1} = \frac{\sum_{i=1}^n q_i^{t+1} (x_i - \mu_1^{t+1})^2}{\sum_{i=1}^n q_i^{t+1}}$$

Similarly, $\frac{\partial L}{\partial \sigma_2} = 0 \Rightarrow (\sigma_2^2)^{t+1} = \frac{\sum_{i=1}^n (1-q_i^{t+1}) (x_i - \mu_2^{t+1})^2}{\sum_{i=1}^n (1-q_i^{t+1})}$

Intuitively, the new μ_i is the weighted mean of all the points, and σ_i is the weighted variance of all points, so it makes sense.

$$(i) \ell(\{x_i\}; \theta) = \sum_{i=1}^n \log(p(x_i; \theta))$$

$$= \sum_{i=1}^n \log\left(\frac{1}{2} \frac{e^{-\frac{1}{2}(x_i - \mu)^2}}{\sqrt{2\pi}} + \frac{1}{2} \frac{e^{-\frac{1}{2}(x_i - \mu)^2}}{\sqrt{2\pi}}\right)$$

$$(j) q_i^{t+1}(z_1=z_1, \dots, z_n=z_n) = \prod_{i=1}^n P(z_i=z_i | X=x_i; \theta^t) \quad \text{Eq (4)}$$

$$\Rightarrow q_i^{t+1}(z=1) = P(z=1 | X=x_i; \theta^t) = \frac{P(z=1 | X=x_i; \theta^t)}{P(z=1 | X=x_i; \theta^t) + P(z=2 | X=x_i; \theta^t)}$$

$$= \frac{\exp(-\frac{(x_i - \mu)^2}{2})}{\exp(-\frac{(x_i - \mu)^2}{2}) + \exp(-\frac{(x_i + \mu)^2}{2})}$$

$$(k) L(\mu; q^{th}) = C - \sum_{i=1}^n \left[q_i^{th} \left(\frac{(x_i - \mu)^2}{2\sigma_1^2} + \log \sigma_1 \right) + (1 - q_i^{th}) \left(\frac{(x_i - \mu)^2}{2\sigma_2^2} + \log \sigma_2 \right) \right]$$

plug in $\mu_1 = \mu$, $\mu_2 = -\mu$, $\sigma_1 = \sigma_2 = 1$

$$L(\mu; q^{th}) = C - \sum_{i=1}^n \left[q_i^{th} \frac{(x_i - \mu)^2}{2} + (1 - q_i^{th}) \frac{(x_i + \mu)^2}{2} \right]$$

$$\begin{aligned} \frac{\partial L(\mu; q^{th})}{\partial \mu} &= - \sum_{i=1}^n \left[q_i^{th} \frac{2(x_i - \mu)(-1)}{2} + (1 - q_i^{th}) \frac{2(x_i + \mu)}{2} \right] \\ &= - \sum_{i=1}^n \left[-q_i^{th}(x_i - \mu) + (1 - q_i^{th})(x_i + \mu) \right] \\ &= - \sum_{i=1}^n \left[-q_i^{th}x_i + q_i^{th}\mu + (1 - q_i^{th})x_i + (1 - q_i^{th})\mu \right] \\ &= - \sum_{i=1}^n \left[(1 - 2q_i^{th})x_i + \mu \right] = 0 \end{aligned}$$

$$\Rightarrow \mu^{th} = \frac{1}{n} \sum_{i=1}^n (2q_i^{th} - 1)x_i$$

$$\begin{aligned} (l) \frac{d(\ell(\{x_i\}; \mu)/n)}{d\mu} &= \frac{1}{n} \frac{d(\sum_{i=1}^n \log(\frac{1}{2} \frac{e^{-\frac{1}{2}(x_i - \mu)^2}}{\sqrt{2\pi}} + \frac{1}{2} \frac{e^{-\frac{1}{2}(x_i + \mu)^2}}{\sqrt{2\pi}}))}{d\mu} \\ &= \frac{1}{n} \sum_{i=1}^n \frac{e^{-\frac{1}{2}(x_i - \mu)^2} (-1)(x_i - \mu) + e^{-\frac{1}{2}(x_i + \mu)^2} (1)(x_i + \mu)}{e^{-\frac{1}{2}(x_i - \mu)^2} + e^{-\frac{1}{2}(x_i + \mu)^2}} \\ &= \frac{1}{n} \sum_{i=1}^n \frac{-(e^{-\frac{1}{2}(x_i - \mu)^2} + e^{-\frac{1}{2}(x_i + \mu)^2})\mu + x_i(e^{-\frac{1}{2}(x_i - \mu)^2} - e^{-\frac{1}{2}(x_i + \mu)^2})}{e^{-\frac{1}{2}(x_i - \mu)^2} + e^{-\frac{1}{2}(x_i + \mu)^2}} \\ &= \frac{1}{n} \sum_{i=1}^n \left(\frac{e^{-\frac{1}{2}(x_i - \mu)^2} - e^{-\frac{1}{2}(x_i + \mu)^2}}{e^{-\frac{1}{2}(x_i - \mu)^2} + e^{-\frac{1}{2}(x_i + \mu)^2}} x_i - \mu \right) \\ &= \left[\frac{1}{n} \sum_{i=1}^n (2w_i - 1)x_i \right] - \mu. \quad \text{where } w_i(\mu) = \frac{e^{-\frac{(x_i - \mu)^2}{2}}}{e^{-\frac{(x_i - \mu)^2}{2}} + e^{-\frac{(x_i + \mu)^2}{2}}} \end{aligned}$$

$$\mu_t^{GA} = \mu_t^{GA} + \alpha \frac{d(\ell(\{x_i\}; \mu))}{d\mu} \Big|_{\mu = \mu_t^{GA}}$$

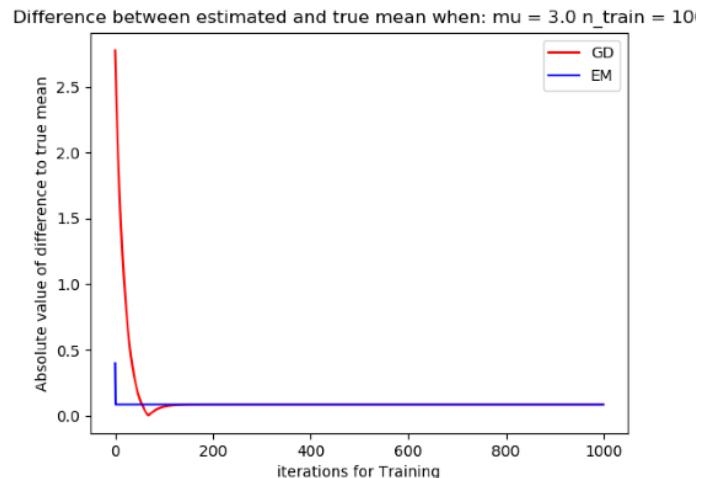
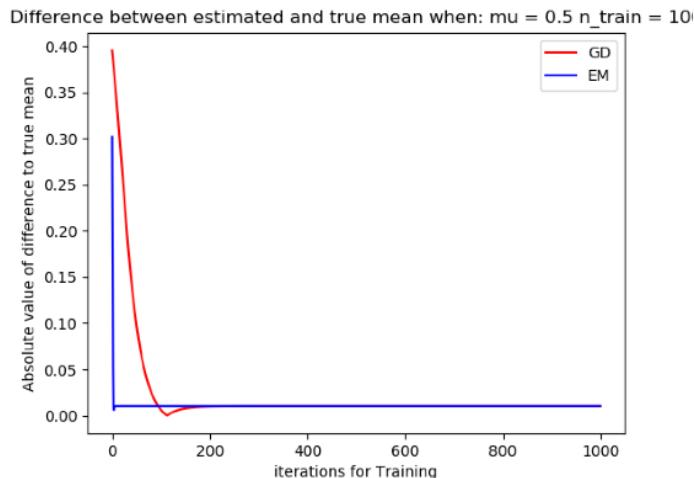
$$= \mu_t^{GA} + \alpha \left[\left[\frac{1}{n} \sum_{i=1}^n ((2w_i(\mu_t^{GA}) - 1)x_i \right] - \mu_t^{GA} \right]$$

$$= (1 - \alpha)\mu_t^{GA} + \alpha \left[\left[\frac{1}{n} \sum_{i=1}^n (2w_i(\mu_t^{GA}) - 1)x_i \right] \right]$$

4 (m)

EM update doesn't depend on u_t directly, while GA does depend on u_t .

EM and GA all depend on ξ_i linearly, where EM has coefficients of $2q_i$ and GA has coefficients of $2\alpha^* w_i$.

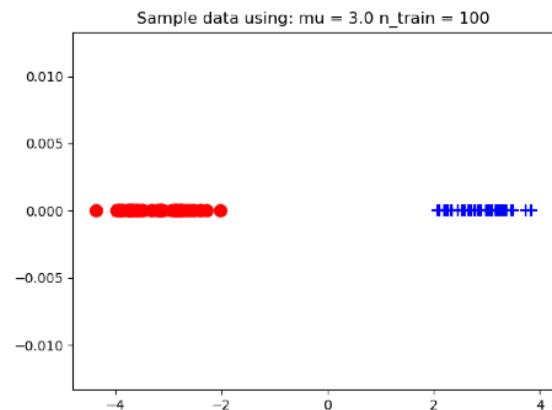
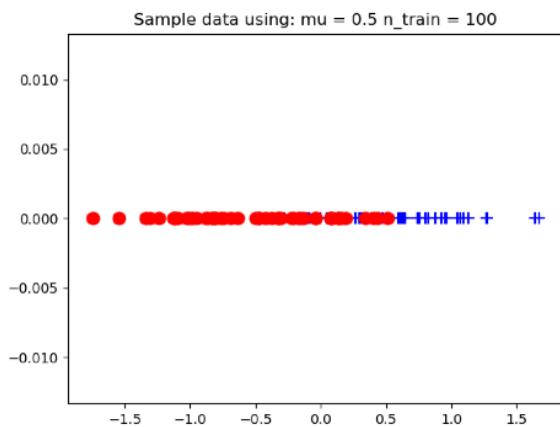


EM converges faster than GD in both cases.

Yes, they matches the updates derived previously.

4(n)

K means will work if two groups are well separated when u are large.
However, it does depend on if u is small or large.

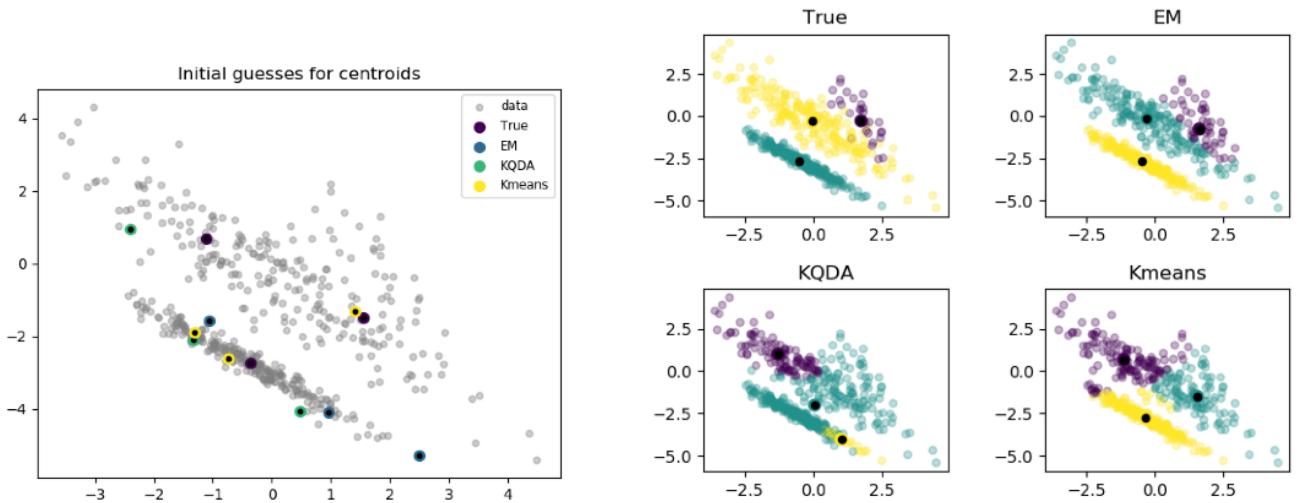


True mean: 0.500, GA (final) estimate: 0.510, EM (final) estimate: 0.510, K-Means (final) estimate: 0.591

True mean: 3.000, GA (final) estimate: 3.085, EM (final) estimate: 3.085, K-Means (final) estimate: 3.085

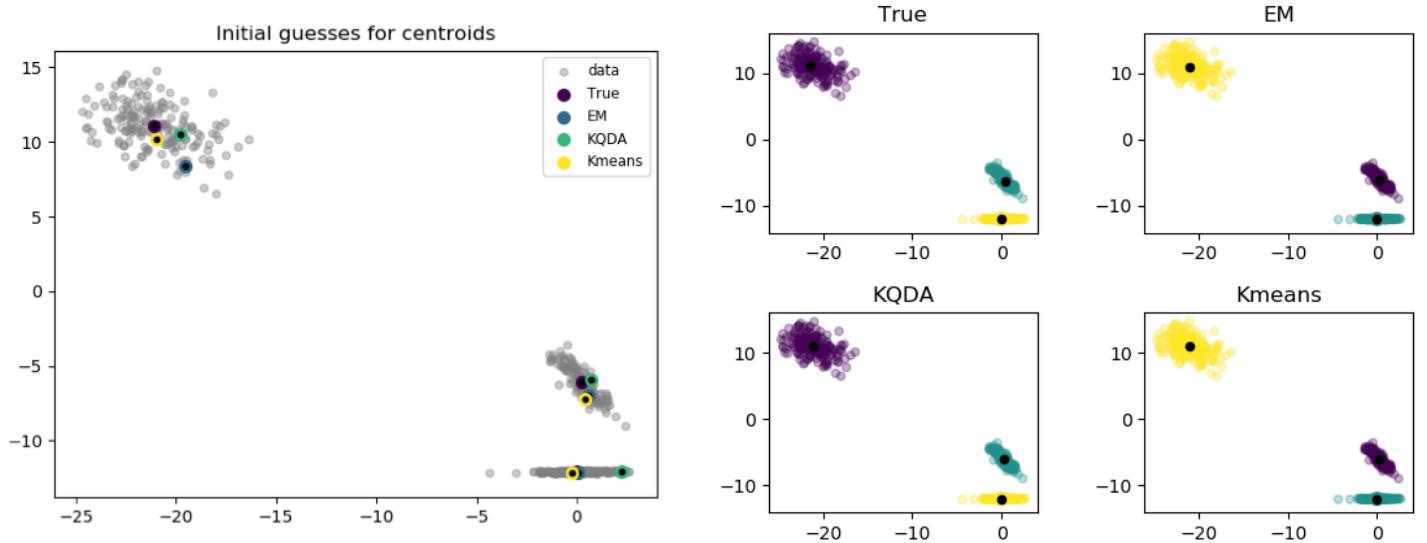
if u is small, two groups will overlap, so k means will not work well.

5. (a)



There is some overlap between groups, so Means and KQDA doesn't work very well, but EM gives a pretty satisfying classification.

(b)



When groups are well separated, all methods work pretty well.

6. Q: Suppose we want to maximize $f(x,y) = x+y$, subject to constraint $x^2+y^2=1$.
 Find the optimal x and y using Lagrangian multiplier.

A: $\max f(x,y) = x+y$, given $g(x,y) = x^2+y^2-1 = 0$

$$\begin{aligned} L(x,y,\lambda) &= f(x,y) + \lambda g(x,y) \\ &= x+y + \lambda(x^2+y^2-1) \end{aligned}$$

$$\frac{\partial L}{\partial x} = 1+2\lambda x = 0 \Rightarrow x = -\frac{1}{2\lambda}$$

$$\frac{\partial L}{\partial y} = 1+2\lambda y = 0 \Rightarrow y = -\frac{1}{2\lambda}$$

$$\frac{\partial L}{\partial \lambda} = x^2+y^2-1=0 \Rightarrow \left(-\frac{1}{2\lambda}\right)^2 + \left(-\frac{1}{2\lambda}\right)^2 + 0 \Rightarrow \lambda = \pm \frac{1}{\sqrt{2}}$$

when $\lambda = \frac{1}{\sqrt{2}}$ $x = -\frac{1}{\sqrt{2}}$, $y = -\frac{1}{\sqrt{2}}$, $f(x,y) = \sqrt{2}$

when $\lambda = -\frac{1}{\sqrt{2}}$ $x = \frac{1}{\sqrt{2}}$, $y = \frac{1}{\sqrt{2}}$, $f(x,y) = \sqrt{2}$.

so $f(x)_{\max} = \sqrt{2}$ when $x=y = \pm \frac{1}{\sqrt{2}}$.