

HOMEWORK 4

Name: Yang Bai UNI: yb2356

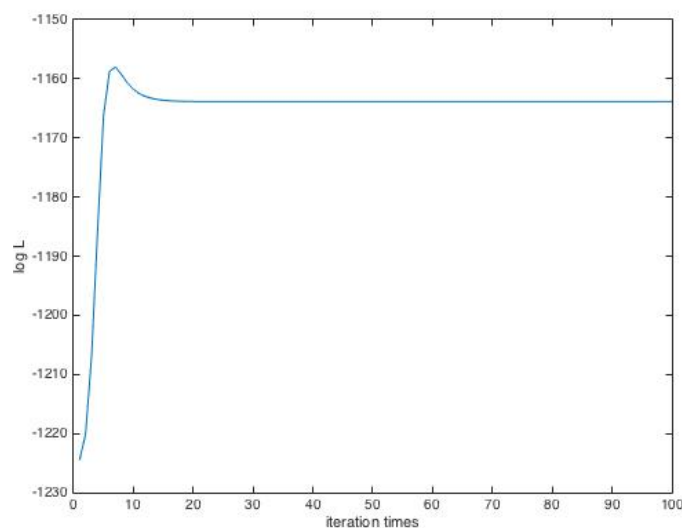
- **Problem1**

a) Implement the EM-GMM algorithm and run it for 100 iterations on the data provided for $K = 2, 4, 8, 10$.

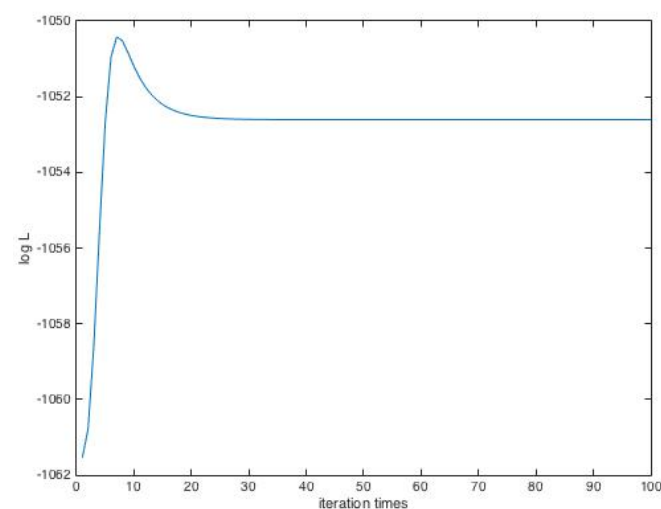
When $K = 8, 10$, I set the iteration number to 200 because 100 seems not enough. If we initialize the parameter randomly, the result usually becomes strange. So I used some initialization algorithm for gaussian mixture model.

b) For each K , plot the log likelihood over the 100 iterations. What pattern do you observe and why might this not be the best way to do model selection?

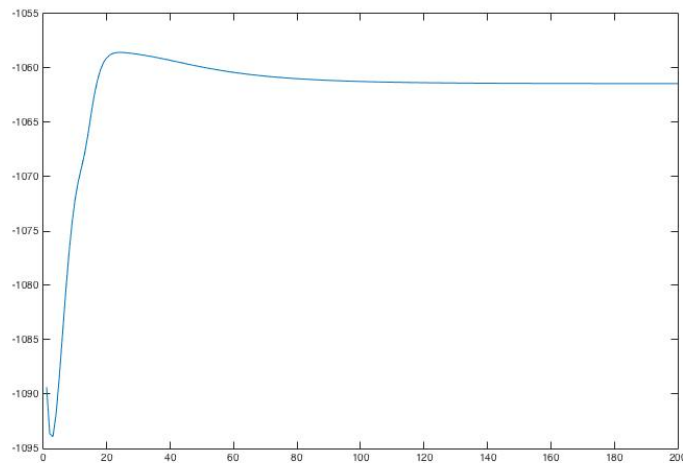
$K = 2$:



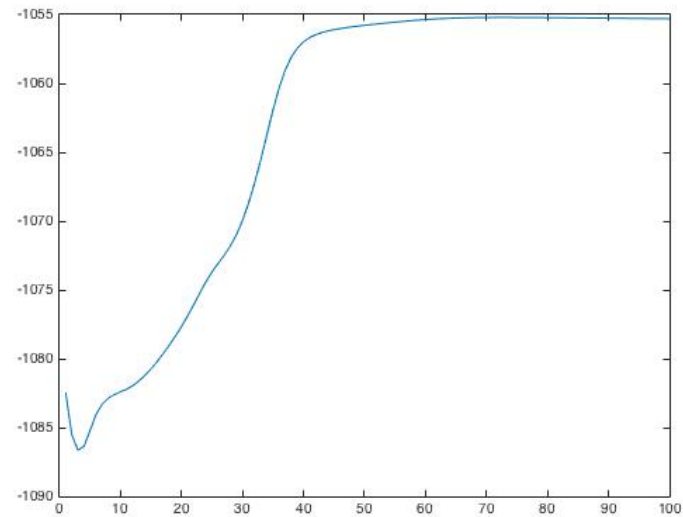
$K = 4$:



K =8:



K =10

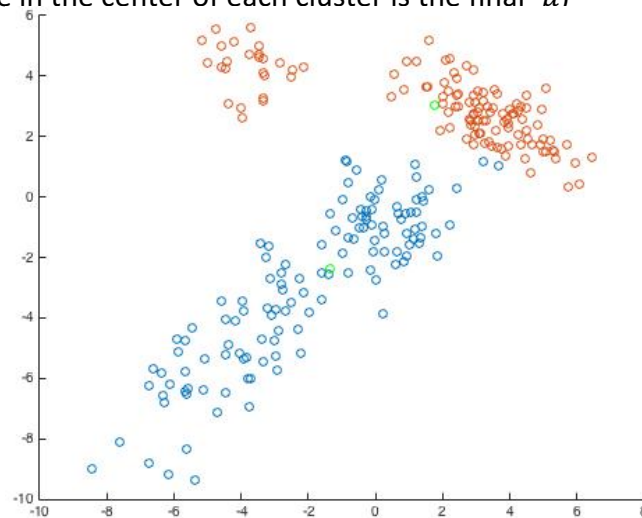


Log likelihoods all converge after several hundred iterations.

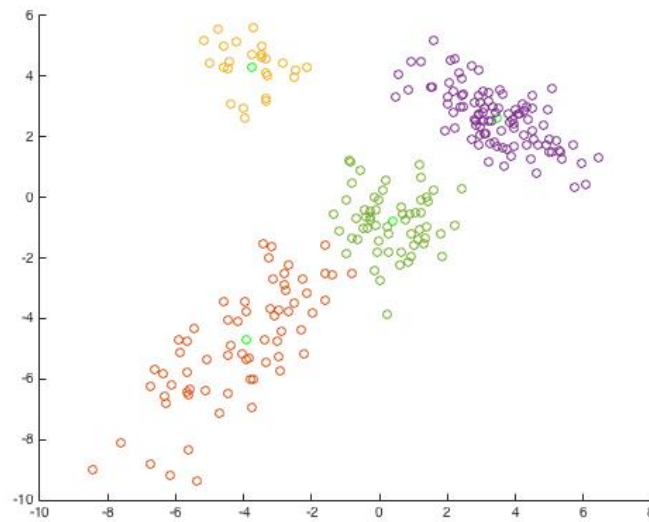
c) For the final iteration of each model, plot the data and indicate the most probable cluster of each observation according to $q(c_i)$ by a cluster-specific symbol. What do you notice about these plots as a function of K?

(The green circle in the center of each cluster is the final μ)

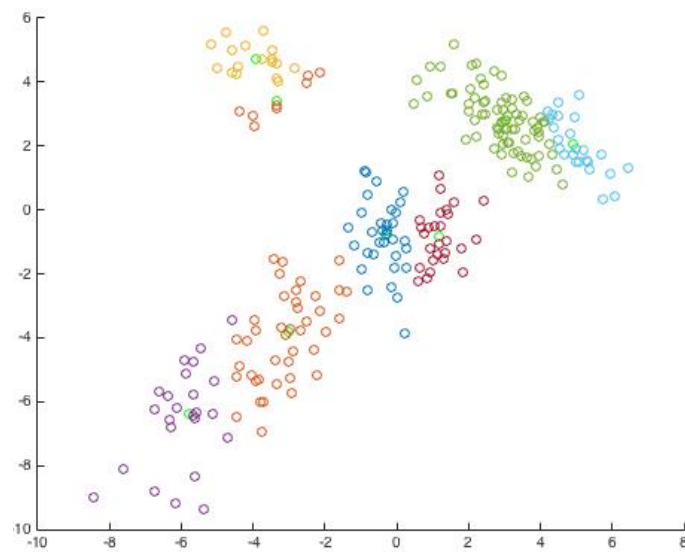
K=2



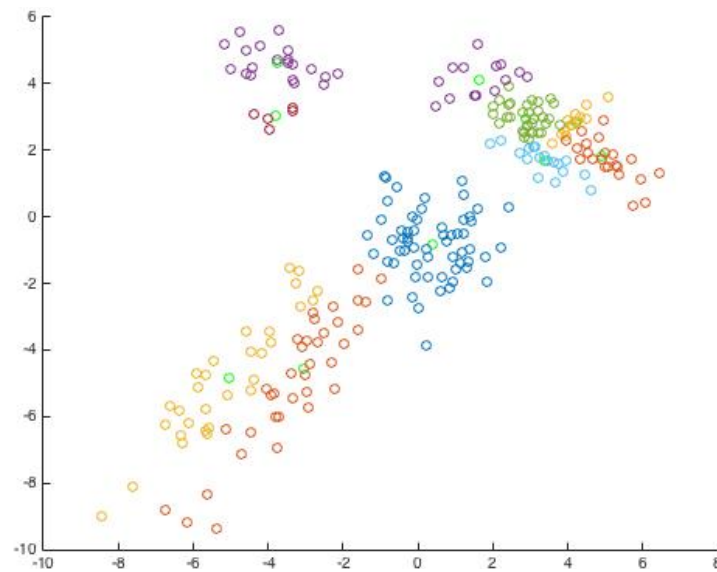
K=4



K=8



K=10



When $k = 8, 10$, the clustering result looks mixed up. The probable explanation is that the data does not have so many distinct clusters, so when K is too big, the algorithm can not return a correct result.

• **Problem 2**

variational objective function

$$\mathcal{L} = \mathbb{E}[\ln p(x, c, \pi, \mu, \Lambda)] - \mathbb{E}[\ln q]$$

The first part of L:

$$\mathbb{E}[\ln p(x, c, \pi, \mu, \Lambda)]$$

$$= \sum_{i=1}^n \mathbb{E}[\ln p(x_i | c_i, \mu_{ci}, \Lambda_{ci}^{-1})] + \sum_{i=1}^n \mathbb{E}[\ln p(c_i | \pi)] + \mathbb{E}[\ln p(\pi)]$$

$$+ \sum_{j=1}^K \mathbb{E}[\ln p(\mu_j)] + \sum_{j=1}^K \mathbb{E}[\ln p(\Lambda_j)]$$

$$\mathbb{E}[\ln p(x_i | c_i, \mu_{ci}, \Lambda_{ci}^{-1})] = \frac{1}{2} \mathbb{E}[\ln |\Lambda_{ci}|] - \frac{1}{2} \mathbb{E}[(x_i - \mu_{ci})^T \Lambda_{ci} (x_i - \mu_{ci})] + -\frac{d}{2} \ln 2\pi$$

$$= -\frac{d}{2} \ln 2\pi + \frac{1}{2} \sum_{j=1}^K \phi_i(j) \mathbb{E}[\ln |\Lambda_j|]$$

$$- \frac{1}{2} \sum_{j=1}^K \phi_i(j) \left[(x_i - \mathbb{E}[\mu_j])^T \mathbb{E}[\Lambda_j] (x_i - \mathbb{E}[\mu_j]) - \text{trace}(\mathbb{E}[\Lambda_j] \Sigma_j') \right]$$

$$\mathbb{E}[\ln p(c_i | \pi)] = \sum_{j=1}^K \phi_i(j) \mathbb{E}[\ln \pi_j] = \sum_{j=1}^K \phi_i(j) \left(\psi(\alpha_j') - \psi\left(\sum_k \alpha_k'\right) \right)$$

$$\mathbb{E}[\ln p(\pi)] = \mathbb{E} \left[\ln \left(\frac{\Gamma(\sum_{j=1}^K \alpha_j)}{\prod_{j=1}^K \Gamma(\alpha_j)} \prod_{j=1}^K \pi_j^{\alpha_j-1} \right) \right]$$

$$= \ln \Gamma\left(\sum_{j=1}^K \alpha_j\right) - \sum_{j=1}^K \ln \Gamma(\alpha_j) + \sum_{j=1}^K (\alpha_j - 1) \mathbb{E}[\ln \pi_j]$$

$$= \sum_{j=1}^K (\alpha_j - 1) \left(\psi(\alpha_j') - \psi\left(\sum_k \alpha_k'\right) \right) + \ln \Gamma\left(\sum_{j=1}^K \alpha_j\right) - \sum_{j=1}^K \ln \Gamma(\alpha_j)$$

$$\mathbb{E}[\ln p(\mu_j)] = -\frac{1}{2c} (m_j'^T m_j' + \text{trace}(\Sigma_j')) + \text{const}$$

$$\mathbb{E}[\ln p(\Lambda_j)] = \frac{a-d-1}{2} \mathbb{E}[\ln |\Lambda_j|] - \frac{1}{2} \text{trace}(B \cdot \mathbb{E}[\Lambda_j]) + \frac{ad}{2} \ln 2 - \ln \Gamma_d\left(\frac{a}{2}\right)$$

$$\mathbb{E}[\ln q(c_i | \pi)] = \sum_{j=1}^K \phi_i(j) \ln \phi_i(j)$$

$$\begin{aligned} \mathbb{E}[\ln q(\Lambda_j)] &= \frac{a'_j - d - 1}{2} \mathbb{E}[\ln |\Lambda_j|] - \frac{1}{2} \text{trace}(B'_j \cdot \mathbb{E}[\Lambda_j]) - \frac{a'_j d}{2} \ln 2 + \frac{a'_j}{2} \ln |B'_j| \\ &\quad - \ln \Gamma_d \left(\frac{a'_j}{2} \right) \end{aligned}$$

$$\begin{aligned} \mathbb{E}[\ln q(\pi)] &= \mathbb{E} \left[\ln \left(\frac{\Gamma(\sum_{j=1}^K \alpha'_j)}{\prod_{j=1}^K \Gamma(\alpha'_j)} \prod_{j=1}^K \pi_j^{\alpha'_j - 1} \right) \right] \\ &= \sum_{j=1}^K (\alpha'_j - 1) \mathbb{E}[\ln \pi_j] + \ln \Gamma \left(\sum_{j=1}^K \alpha'_j \right) - \sum_{j=1}^K \ln \Gamma(\alpha'_j) \\ &= \sum_{j=1}^K (\alpha'_j - 1) \left(\psi(\alpha'_j) - \psi \left(\sum_k \alpha'_k \right) \right) + \ln \Gamma \left(\sum_{j=1}^K \alpha'_j \right) - \sum_{j=1}^K \ln \Gamma(\alpha'_j) \\ \mathbb{E}[\ln q(\mu_j)] &= -\frac{1}{2} \ln |\Sigma'_j| + \text{const} \end{aligned}$$

And we can calculate that

$$\mathbb{E}[\Lambda_j] = a'_j B'^{-1}_j$$

$$\mathbb{E}[\ln |\Lambda_j|] = \frac{d(d-1)}{4} \ln \pi + d \ln 2 - \ln |B'_j| + \sum_{k=1}^d \psi \left(\frac{a'_j}{2} + (1-k)/2 \right)$$

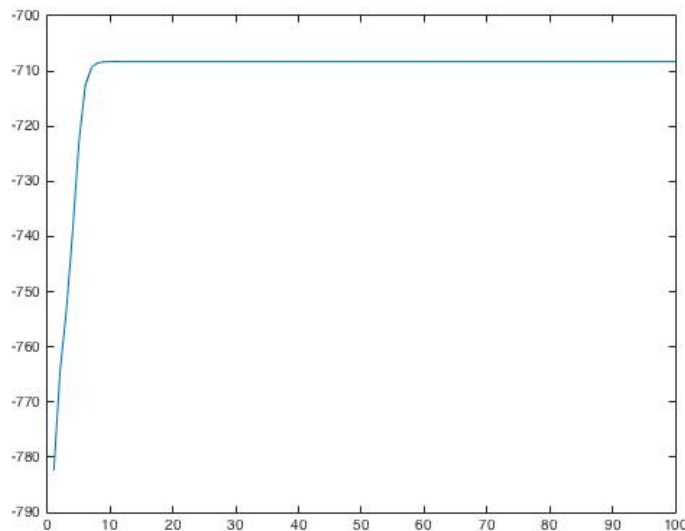
$$\mathbb{E}[\mu_j] = m'_j$$

a) Implement the variational inference algorithm discussed in class and in the notes for $K = 2, 4, 10, 25$ and 100 iterations each.

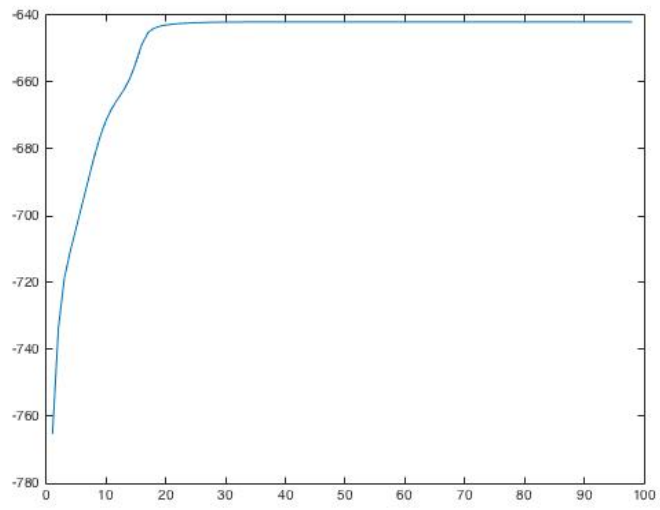
b) For each K , plot the variational objective function over the 100 iterations. What pattern do you observe?

When calculating L in the code, I did not add any constant in the above formula.

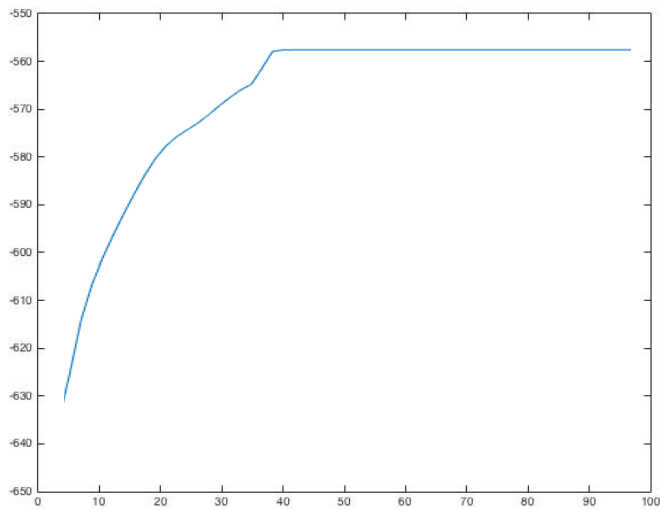
$K = 2$:



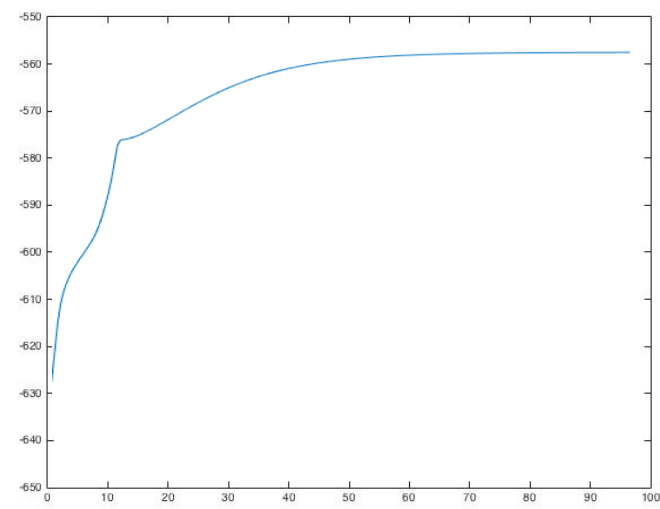
K =4



K =10



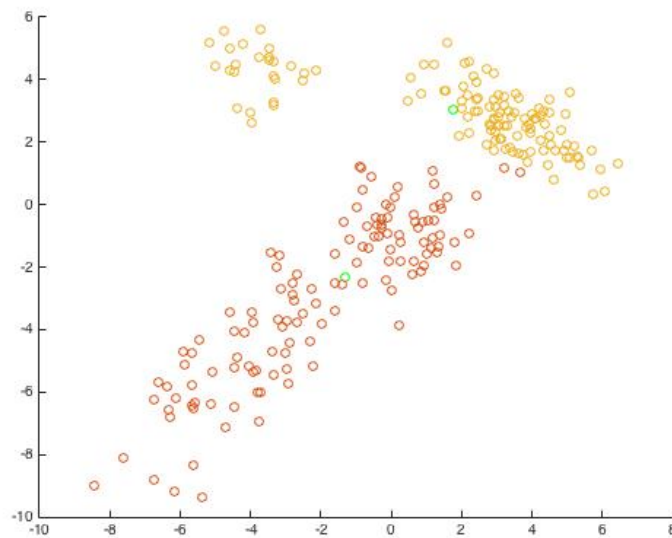
K =25



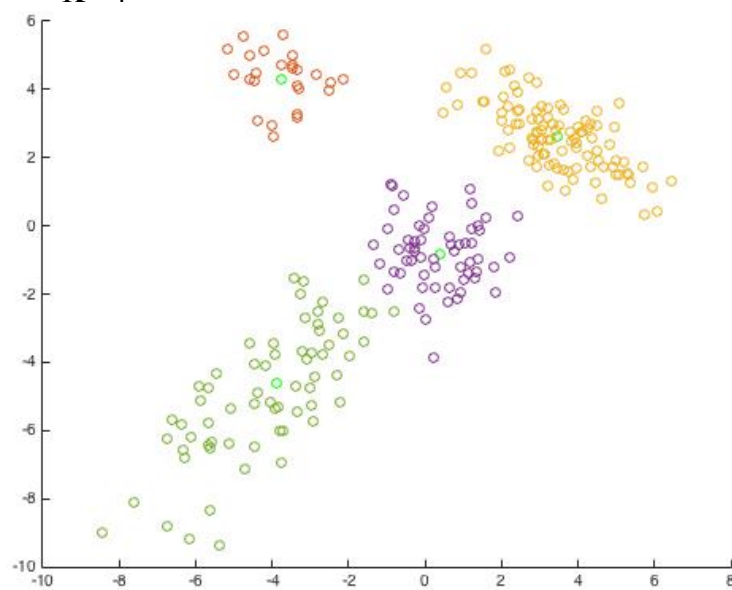
c) For the final iteration of each model, plot the data. What do you notice about these plots as a function of K?

(The green circle in the center of each cluster is the final μ)

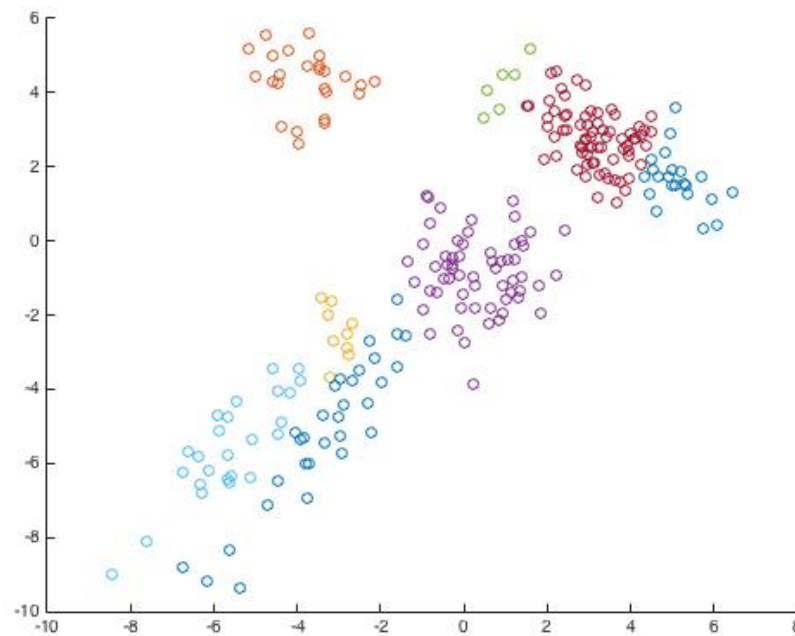
K = 2



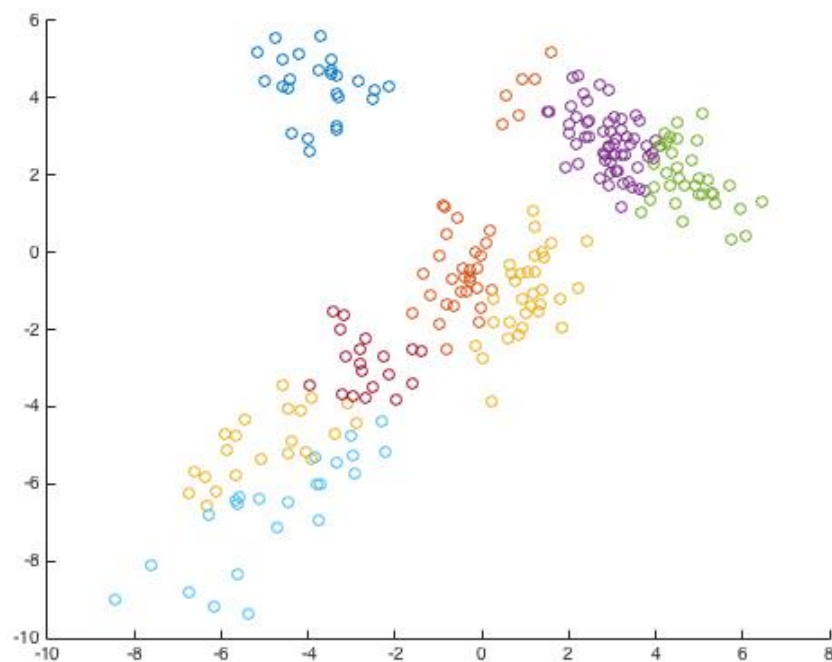
K = 4



K=10



K=25

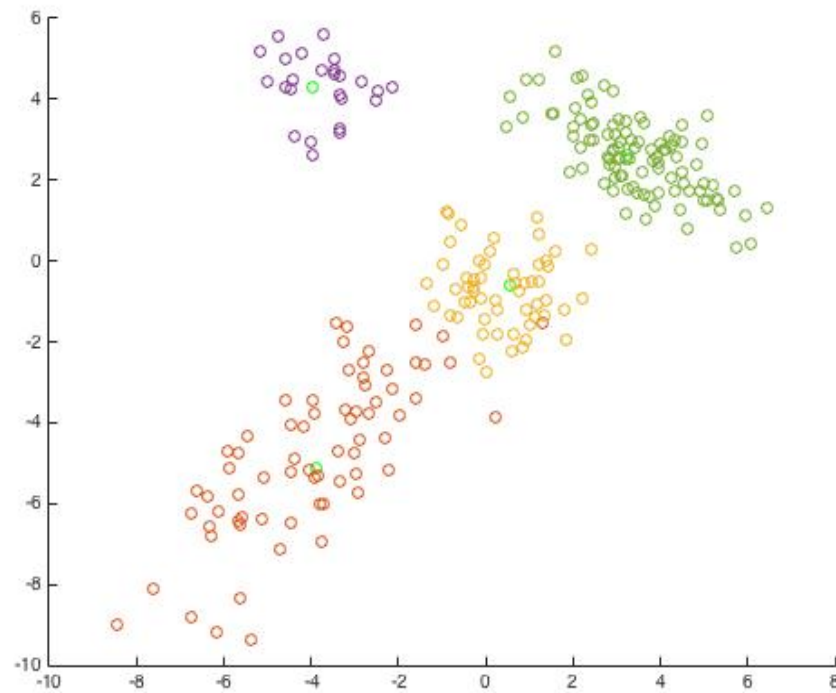


The clustering result of K =10, 25 varies. Generally, the algorithm can not cluster the data into 10 or 25 clusters. The number of clusters are mostly smaller than 10.

- **Problem 3.**

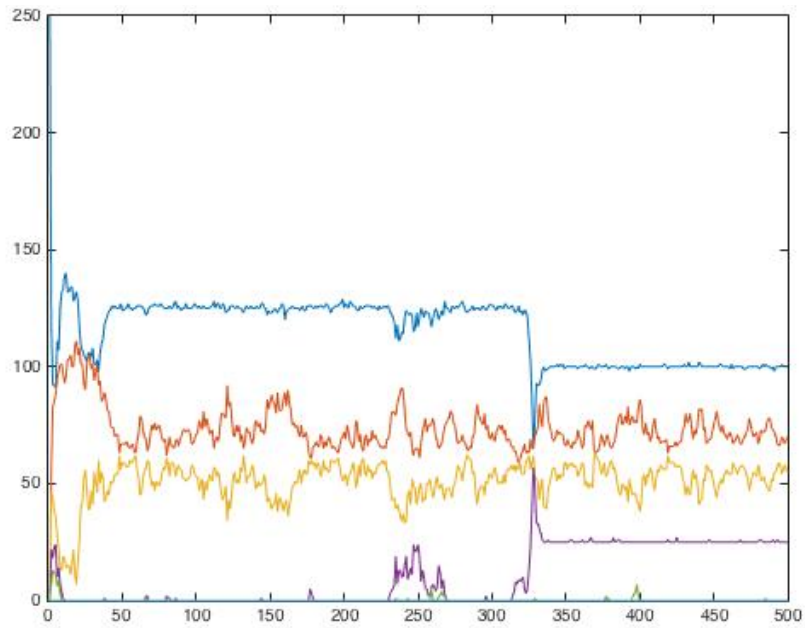
a) Implement the above-mentioned Gibbs sampling algorithm discussed in class and described in the notes. Run your algorithm on the data provided for 500 iterations.

I tested the program for many times and after 500 iterations, the number of clusters is 4.

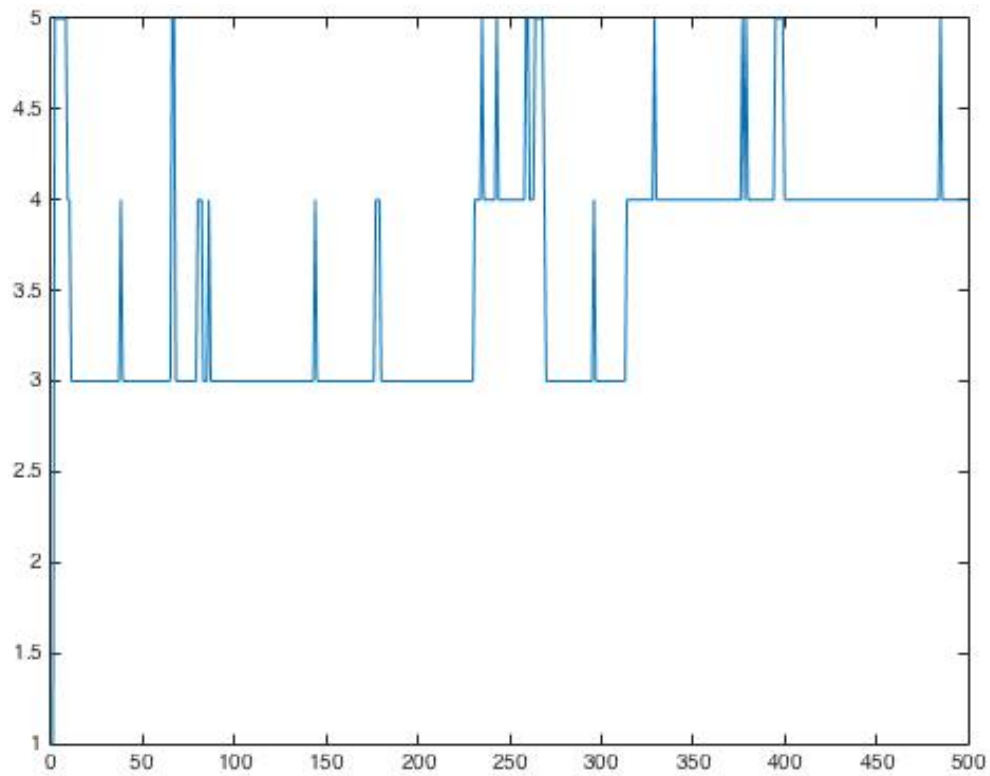


b) Plot the number of observations per cluster as a function of iteration for the six most probable clusters. These should be shown as lines that never cross; for example the i th value of the “second” line will be the number of observations in the second largest cluster after completing the i th iteration. If there are fewer than six clusters then set the remaining values to zero.

The number of data points in each cluster always changes a little although the iteration time is set to 500.



c) Plot of the total number of clusters that contain data as a function of iteration.



There are still some peaks with value 5 when the result becomes stable. This may be caused by the sampling process of ci, μ, Λ in every iteration.