

Homework 4 Bayesian Machine Learning

Problem Set-up

We are given observations $X = \{x_1, \dots, x_n\}$ where each $x_i \in \mathbb{R}^d$. We model this as being generated from a Gaussian mixture model of the form

$$x_i | c_i \sim \text{Normal}(\mu_{c_i}, \Lambda_{c_i}^{-1}), \quad c_i \stackrel{iid}{\sim} \text{Discrete}(\pi)$$

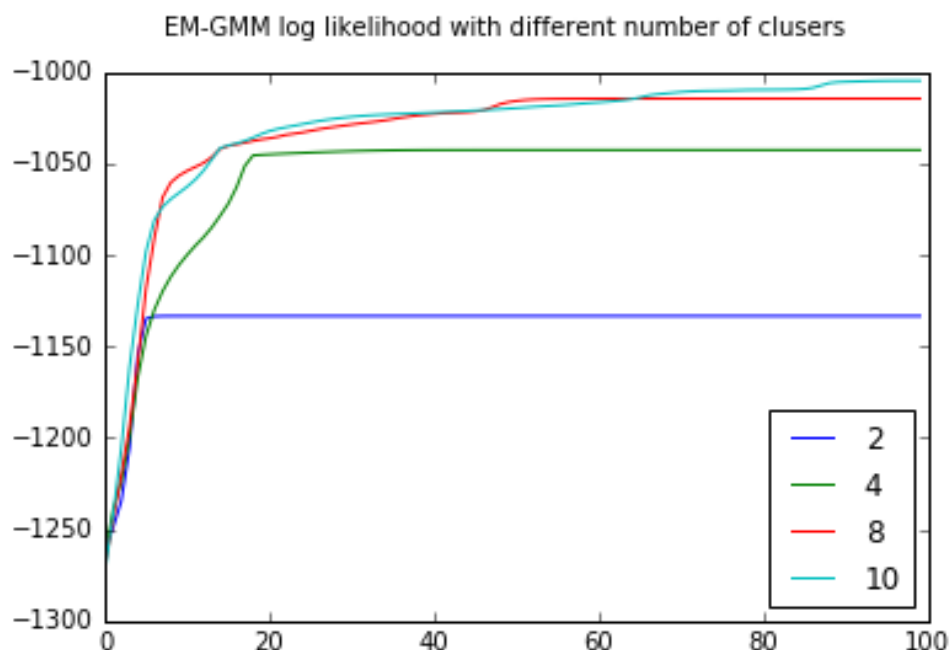
In this homework, you will implement three algorithms for learning this mixture model, one based on maximum likelihood EM, one on variational inference and one on Gibbs sampling. Use the data provided for all experiments.

Problem 1. (30 points)

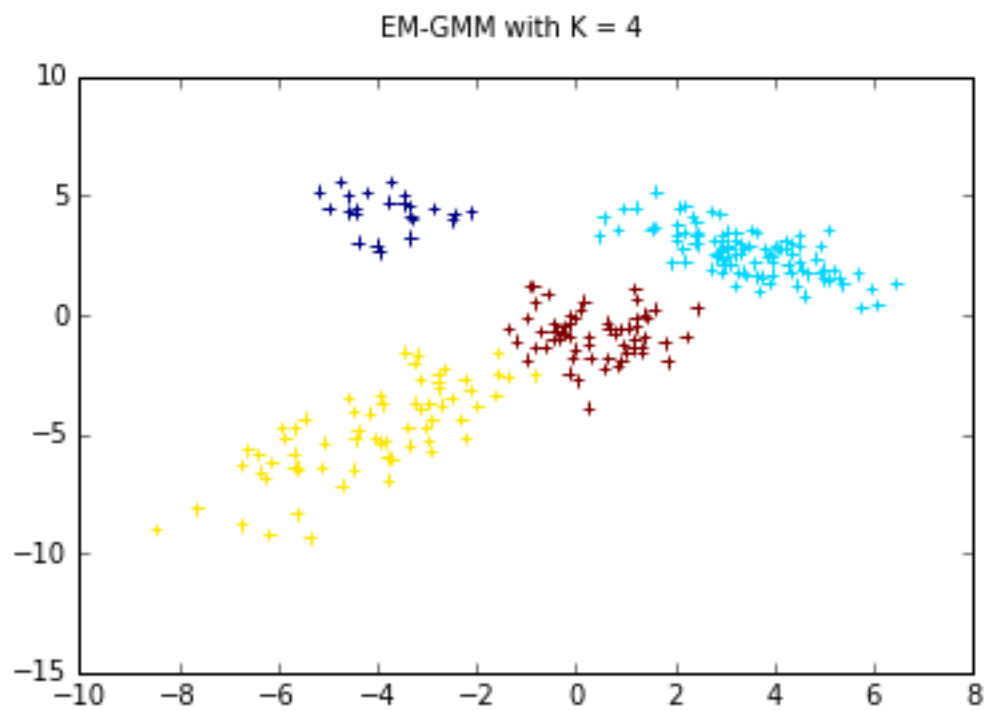
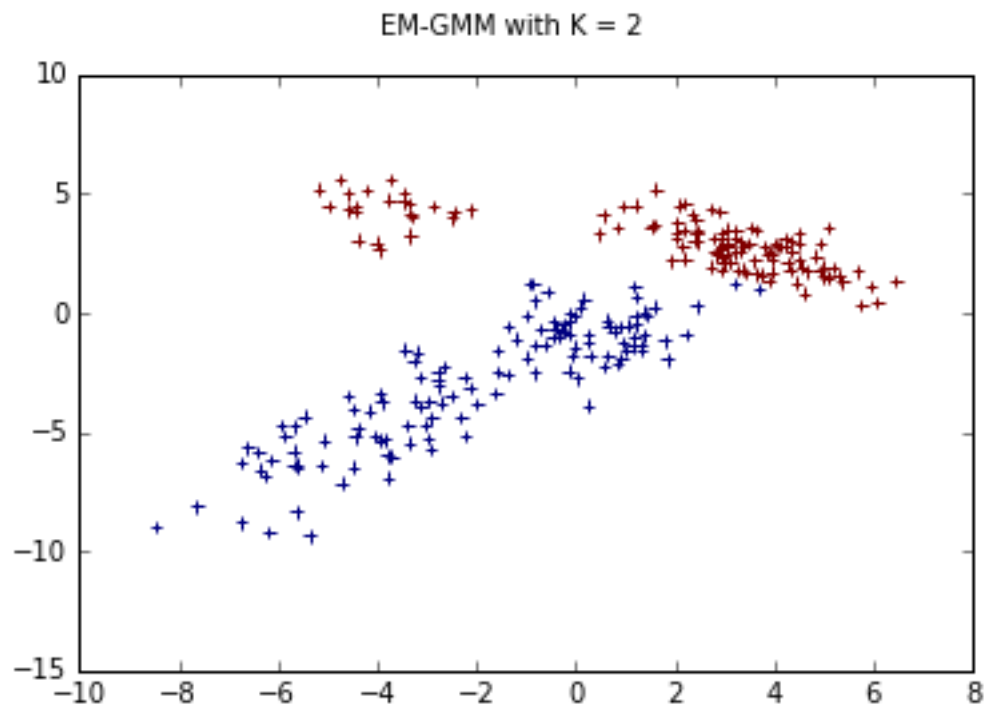
In this problem, you will implement the EM algorithm for learning maximum likelihood values of π and each (μ_j, Λ_j) for $j = 1, \dots, K$. The algorithm is given in the notes, and also in Section 9.2 of Bishop's book.

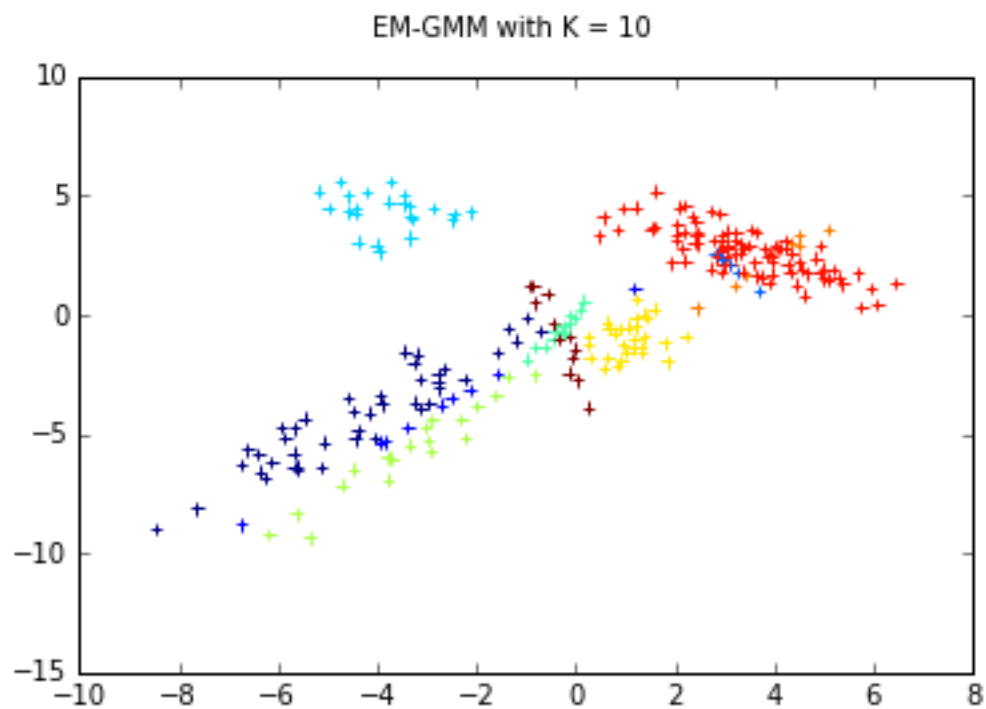
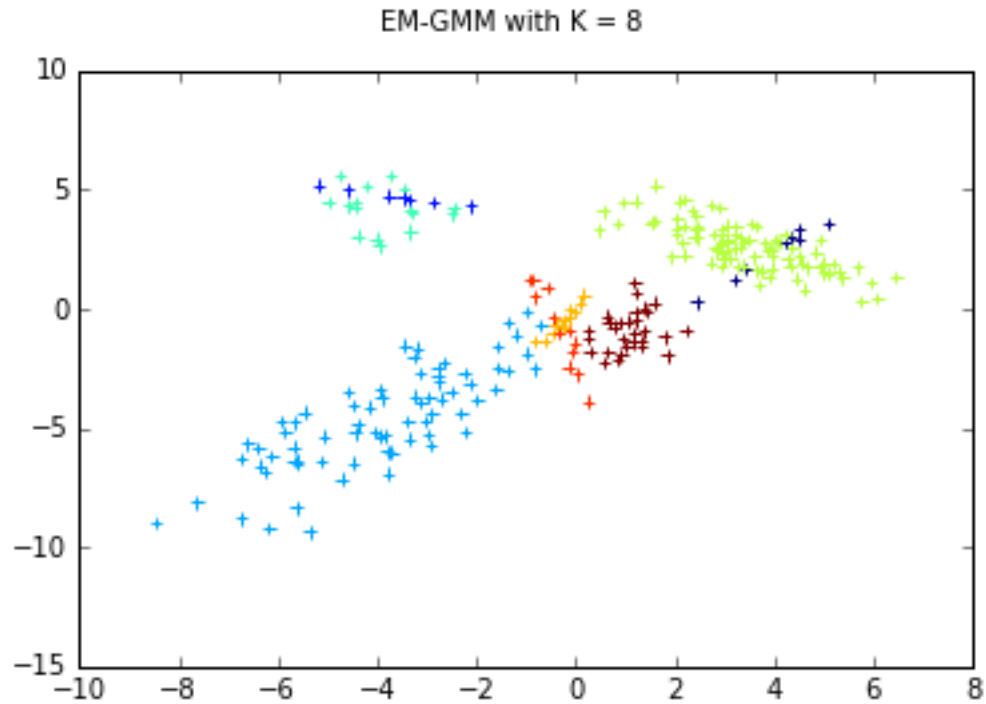
- Implement the EM-GMM algorithm and run it for 100 iterations on the data provided for $K = 2, 4, 8, 10$.
- For each K , plot the log likelihood over the 100 iterations. What pattern do you observe and why might this not be the best way to do model selection?
- For the final iteration of each model, 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 ?

B,



C,





As k increases, the number of non-empty clusters increase with it, which is the property of this type of classification. Originally I would have separated them into 4 clusters, but the higher number of K revealed some interesting structure too as it decomposed previous clusters into smaller one (like the brown one in the last plot) which indeed seems like it deserves its own distribution.

Problem 2. (35 points)

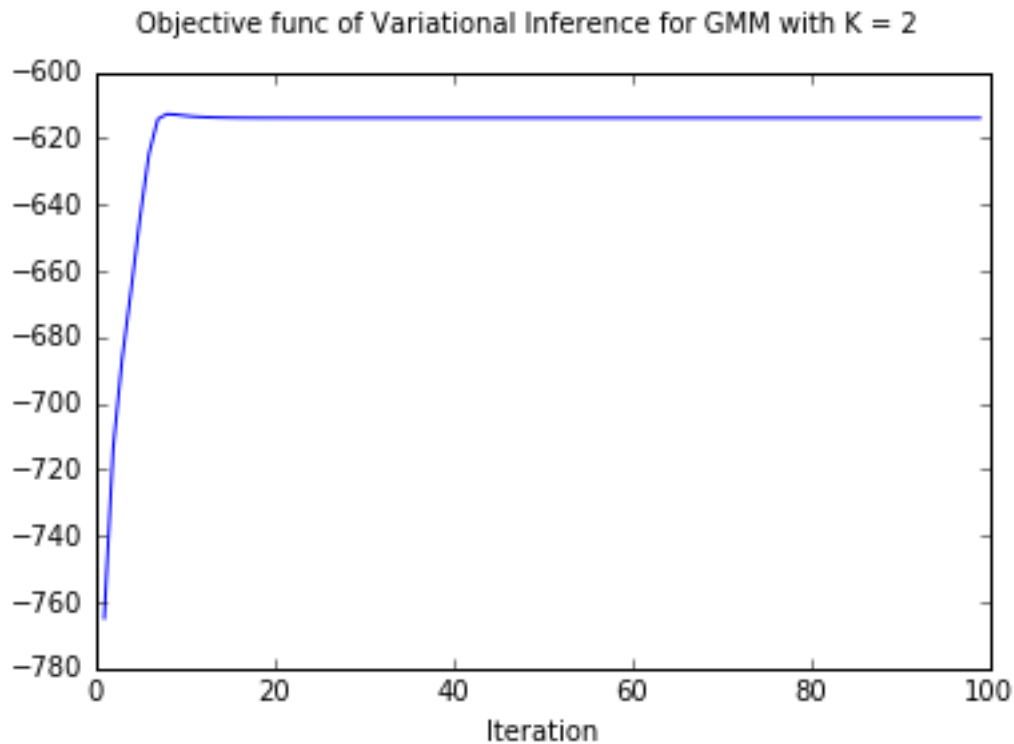
In this problem, you will implement a variational inference algorithm for approximating the posterior distribution of the GMM variables. We therefore require prior distributions on these variables. For this problem, we use

$$\pi \sim \text{Dirichlet}(\alpha), \quad \mu_j \sim \text{Normal}(0, cI), \quad \Lambda_j \sim \text{Wishart}(a, B)$$

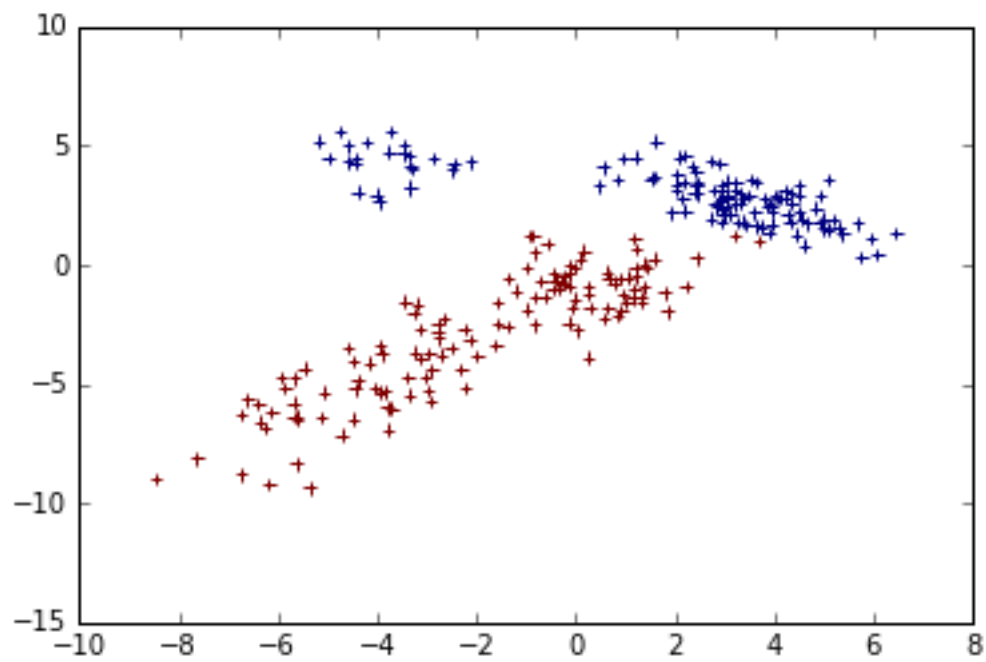
For this problem, set $\alpha = 1$, $c = 10$, $a = d$ and $B = \frac{d}{10}A$ where A is the empirical covariance of the data. Approximate the posterior distribution of these variables with q distributions factorized on π , and each μ_j , Λ_j and c_i as discussed in class.

- Implement the variational inference algorithm discussed in class and in the notes for $K = 2, 4, 10, 25$ and 100 iterations each.
- For each K , plot the variational objective function over the 100 iterations. What pattern do you observe?
- 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 ?

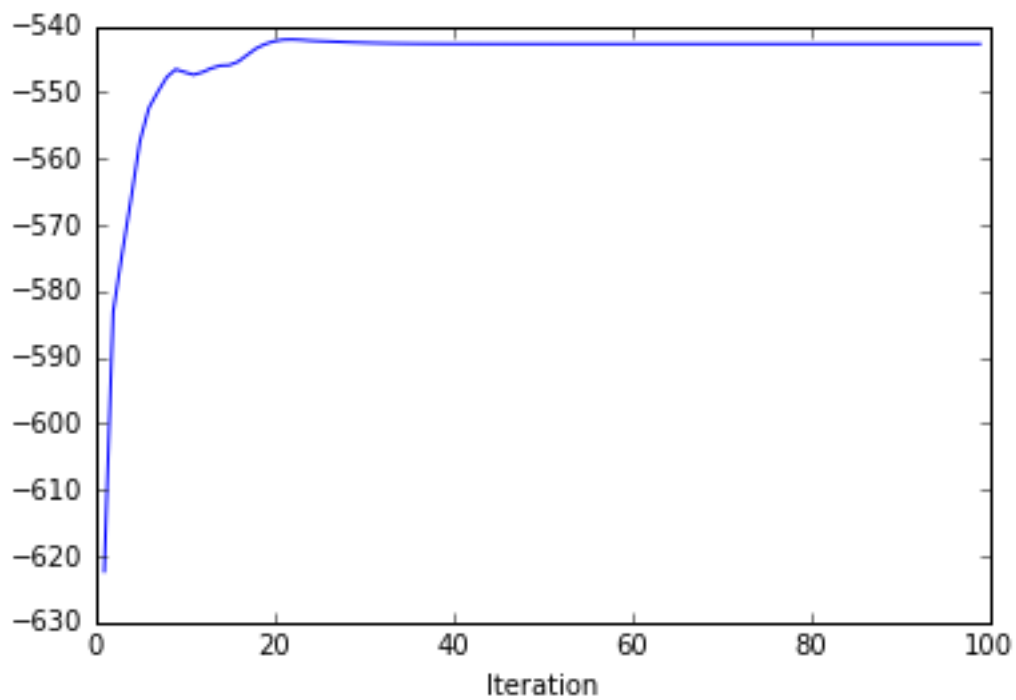
B & C,



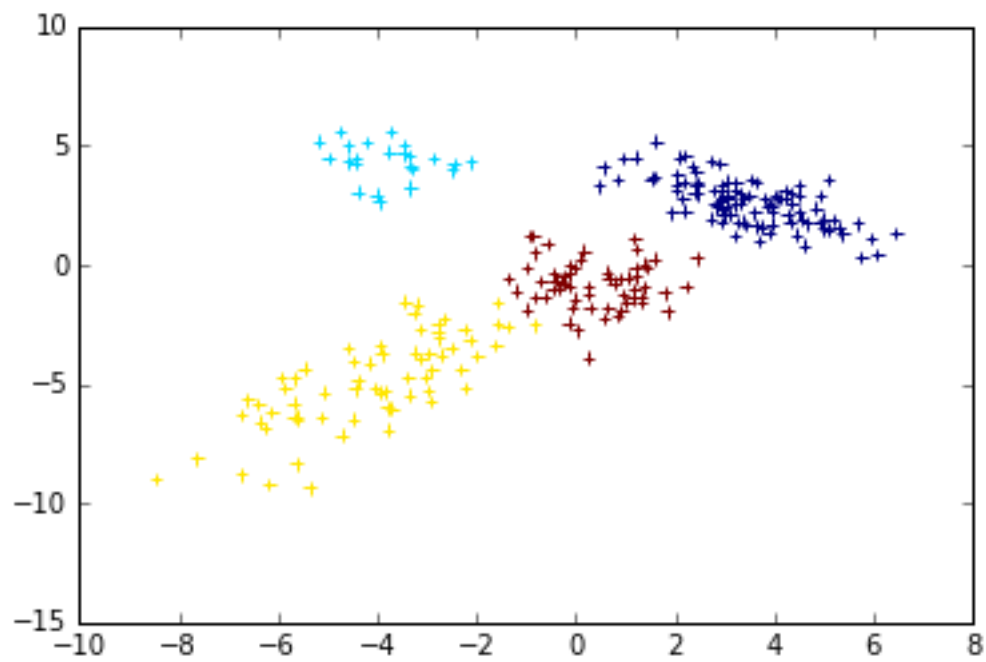
Variational Inference for GMM with $K = 2$



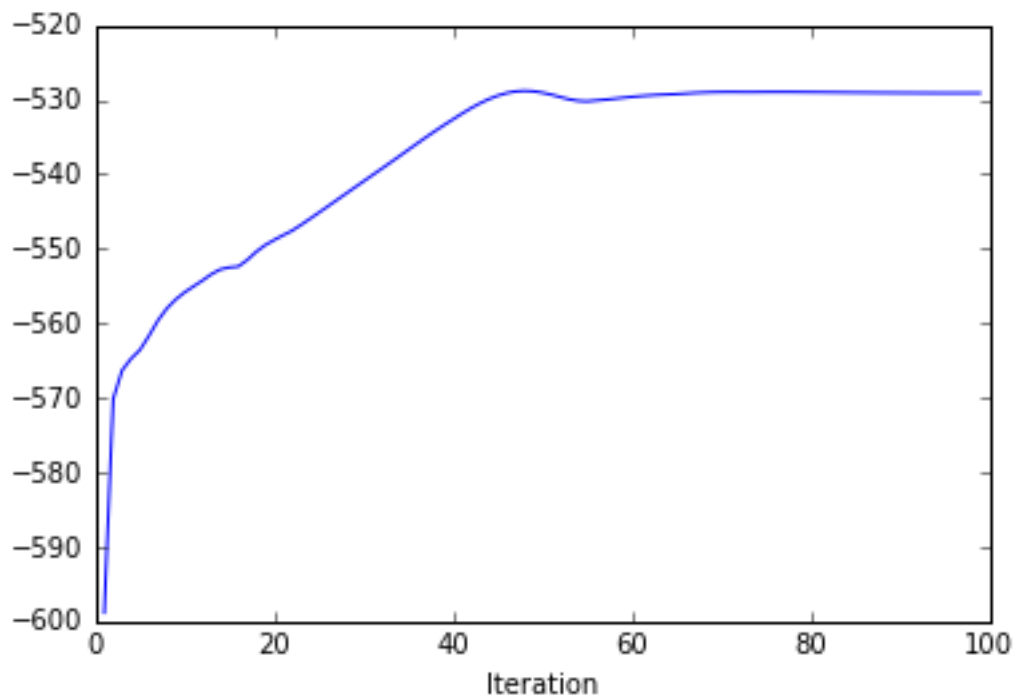
Objective func of Variational Inference for GMM with $K = 4$



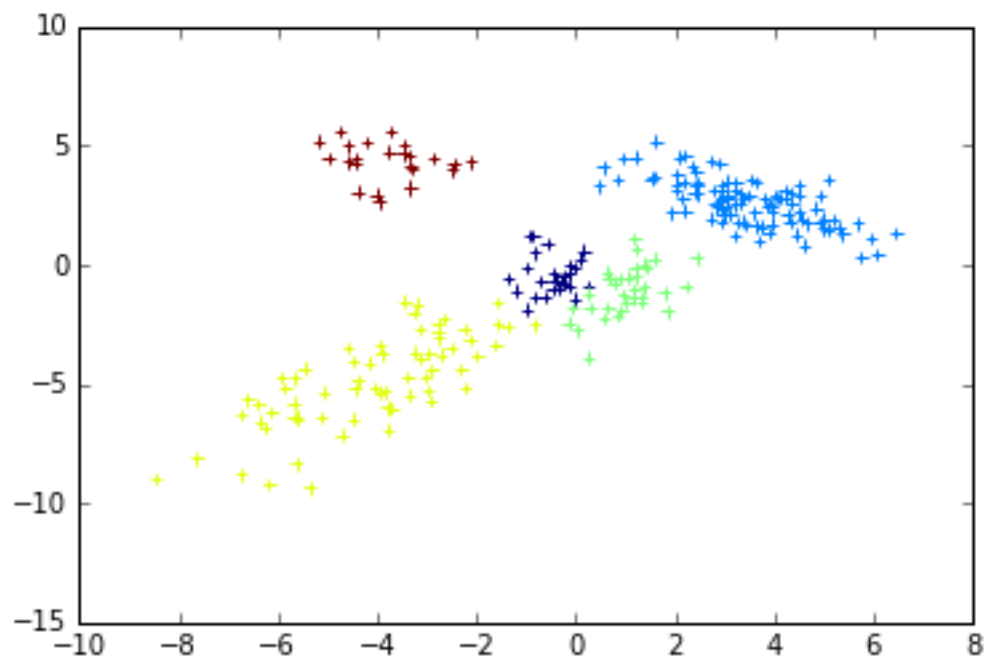
Variational Inference for GMM with $K = 4$



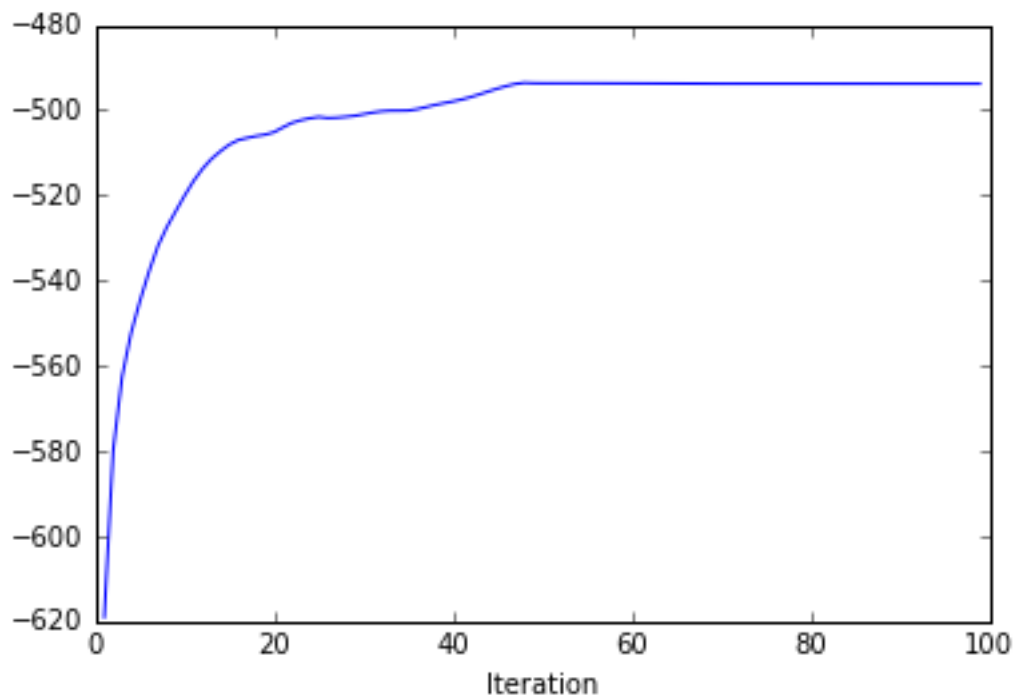
Objective func of Variational Inference for GMM with $K = 10$

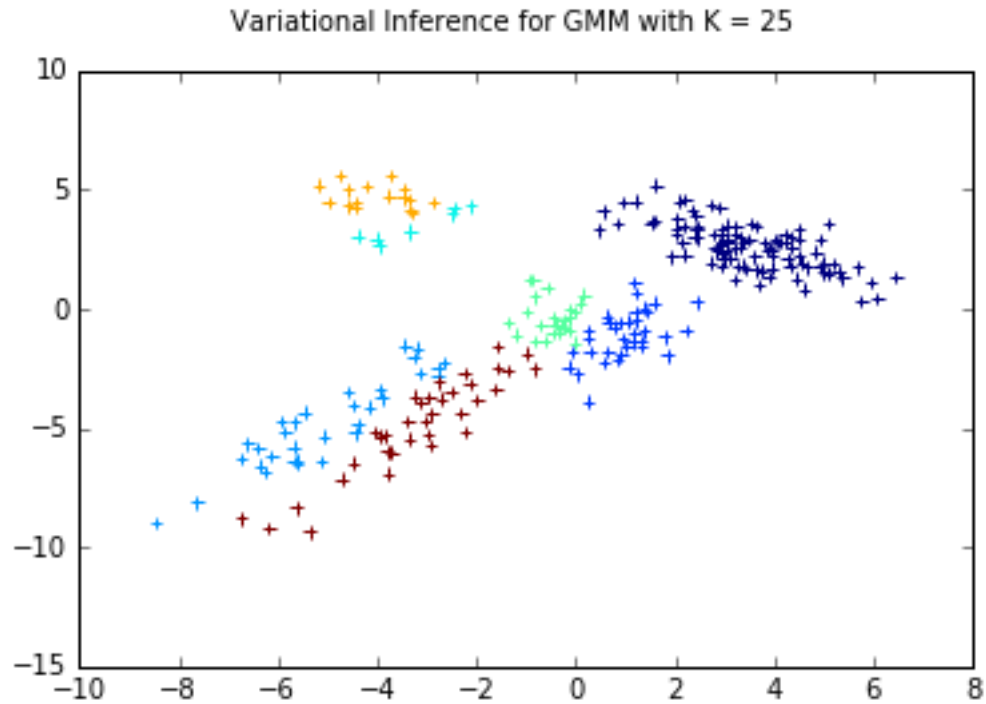


Variational Inference for GMM with $K = 10$



Objective func of Variational Inference for GMM with $K = 25$





For the objective functions, we can see that as the number of k increases, the objective function increases, as well as it takes more time to converge (in $K=2$ it converges pretty soon, and takes a lot more time for higher K).

As for the number of cluster used, we see that from 4 on, we don't actually end up having the max amount of clusters. For $K=10$ its 6, for $K=25$ its 7. That is because our dirichlet prior.

Problem 3. (35 points)

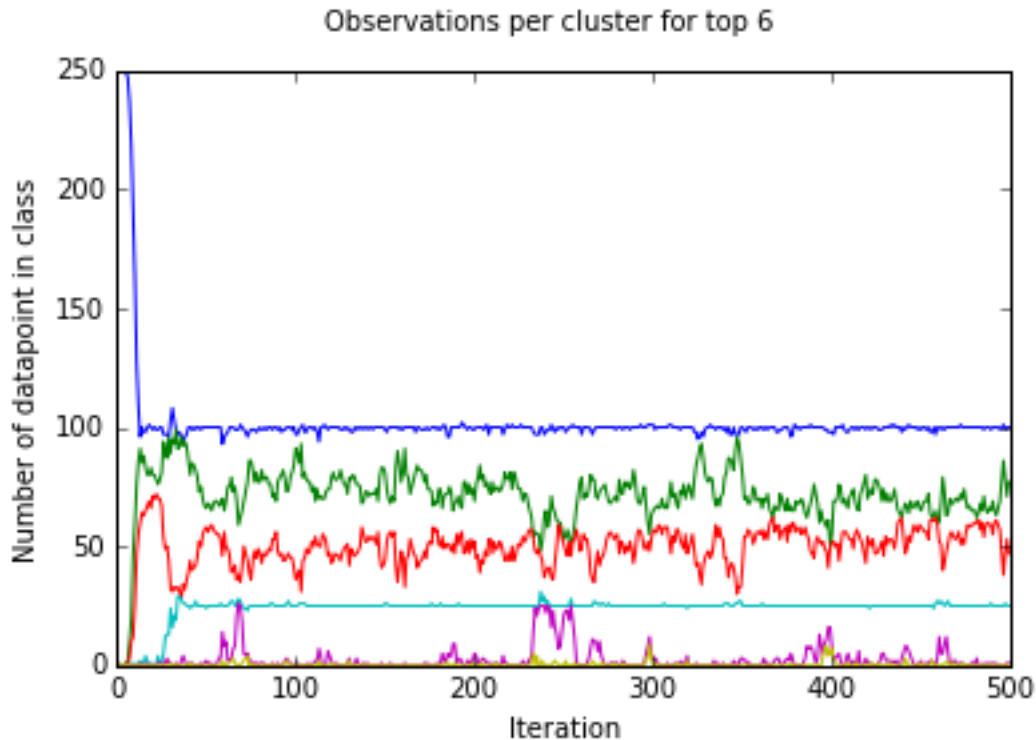
In this problem, you will implement a Bayesian nonparametric sampler for a marginalized version of the GMM. In contrast to Problem 2, in this problem we will use a joint prior on (μ_j, Λ_j) . This is done for computational convenience in calculating the marginal distribution of the data. Specifically, we use the prior distribution

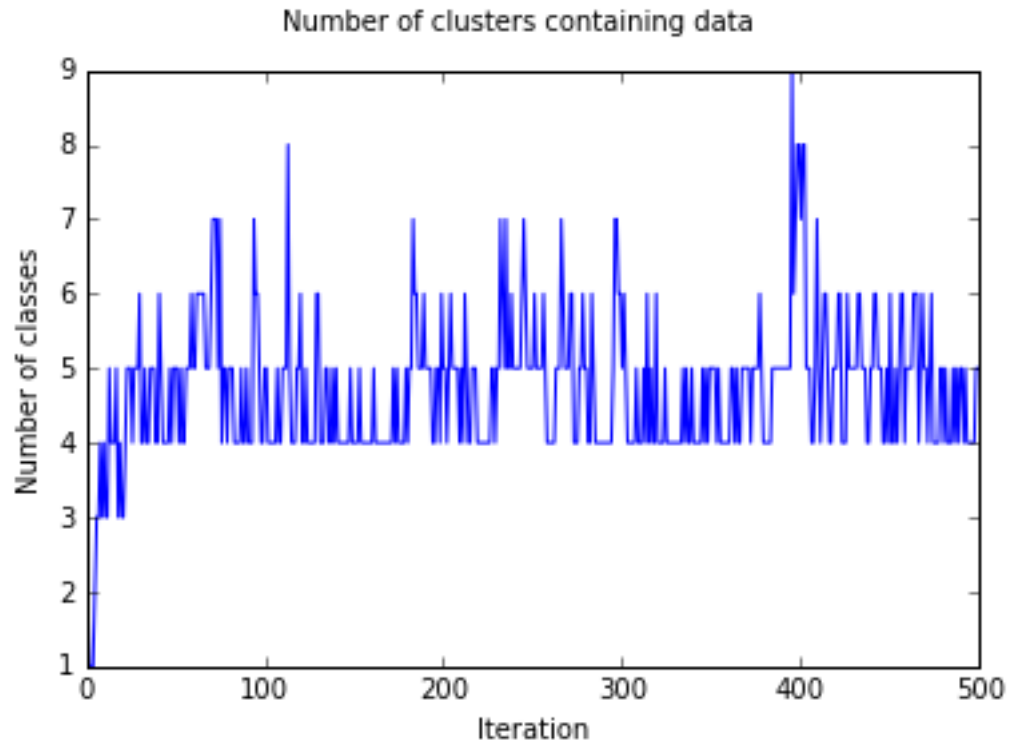
$$\mu_j | \Lambda_j \sim \text{Normal}(m, (c\Lambda)^{-1}), \quad \Lambda_j \sim \text{Wishart}(a, B)$$

as well as the limit of the prior $\pi \sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K)$ as $K \rightarrow \infty$.

In this problem you will implement the marginal sampler where π is integrated out. For this problem, set m to be the empirical mean of the data, $c = 1/10$, $a = d$ and $B = c \cdot d \cdot A$ where A is the empirical covariance of the data. For the “cluster innovation parameter” set $\alpha = 1$.

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





It can be seen that the number of clusters stabilizes around 4, occasionally a couple new is created (these are single or very few point clusters) which die soon after. Two of the clusters seems pretty stable, while the other two is constantly trading data points.