Matt Dickenson `mcd31`
STA561/CS571 — Fall 2013    **Homework 9**    Due: 25 November, 2013

*Homework Notes:* I did not work with anyone else on this homework or refer to resources other than the course notes, textbook, and course Piazza page.

**Problem 1**

**A** The generative model for a continuous $\eta$, a base distribution $G_0$, concentration parameter $\alpha$, and the $\{B_1, ... B_K\}$ partitions ($K = \inf$), is:

$$
\begin{aligned}
(G(\eta \in B_1), ..., G(\eta \in B)K) &\sim \text{Dirich}(\alpha G_0(B_1), ... \alpha G_0(B_K)) \\
p(\eta_i \in B_j) &= \int p(\eta_i \in B_j | G) p(G|G_0) dG \\
&= \frac{\alpha G_0(B_j)}{\sum_K \alpha G_0(B_k)} \\
&\propto \alpha G_0(B_j)
\end{aligned}
$$

The posterior is

$$
G|\eta_{1:n}, \alpha, G_0 \sim DP(\alpha, G_0 + \sum_{i=1}^{n} \delta_{\eta_i}(\eta))
$$

A simple choice for the base distribution is $G_0$ is the Gamma distribution, due to the conjugacy of the Gamma distribution with the Gaussian distribution.

**B** For the cluster assignment step in the Gibbs sampler, we can exploit exchangeability.

```
1 cluster = function(x, centroids, alpha){
2   # do the restaurant process
3   table_counts = restaurant(x, alpha)
4   table_props = table_counts/sum(table_counts)
5   num_tables = length(table_counts)
6
7   # then exploit permutation
8   permuted_x = sample(x)
9   n = length(x)
10  table_assignments = rep(NA, n)
11
12  # pretend each x is last to arrive
13  for(i in 1:length(permuted_x)){
14    table_i = sample(c(1:num_tables), 1, prob=table_props)
15    table_assignments[i] = table_i
16  }
17 }
```

```
18
19 restaurant = function(x, alpha){
20   table_counts = c(1) # number of 'customers' at each 'table'
21                       # first customer at first table
22   for(m in 2:n){
23     tmp = c(table_counts, alpha)
24     table_props = tmp/sum(tmp)
25
26     # assign each 'customer' to a 'table' according to crp
27     table_m = sample(c(1:length(tmp)), 1, prob=table_props)
28     if(table_m==length(tmp)){ table_counts[table_m] = 1}
29     else{ table_counts[table_m] = table_counts[table_m] + 1}
30   }
31   return(table_counts) # sufficient statistic
32 }
```

**C** This algorithm does not discard empty clusters. Because the number of clusters is potentially infinite (in theory), it is possible for an observation to be assigned to a previously empty cluster at any iteration of the Gibbs sampler.

**D** Rather than discarding clusters with fewer than $\gamma$ points assigned to them, we allow the number of potential clusters to be infinite. Thus, this model better addresses the issue of not knowing the number of clusters *a priori*. If we examine the number of clusters at each iteration of the Gibbs sampler (or across multiple runs), we can even get a posterior distribution over the number of clusters.

**Problem 2**