

UCLA Biostatistics 285: Homework 2

Instructor: Dr. Michele Guindani

Soumyakanti Pan, e-mail: span18@ucla.edu

May 19, 2023

1 Problem 1

1.1 Part 1

The function `Sethu_jump` generates the jumps given a truncation option and a α . The `generate_DPH` uses a jump function (here we use the `Sethu_jump`) and takes input a base measure and its parameters along with α , truncation parameter K and number of samples to be generated. The final output is realizations of $DP(\alpha, \mathcal{N}(0, 1))$ approximated by finite truncation with 20 terms as described in Ishwaran and Zarepour (2002).

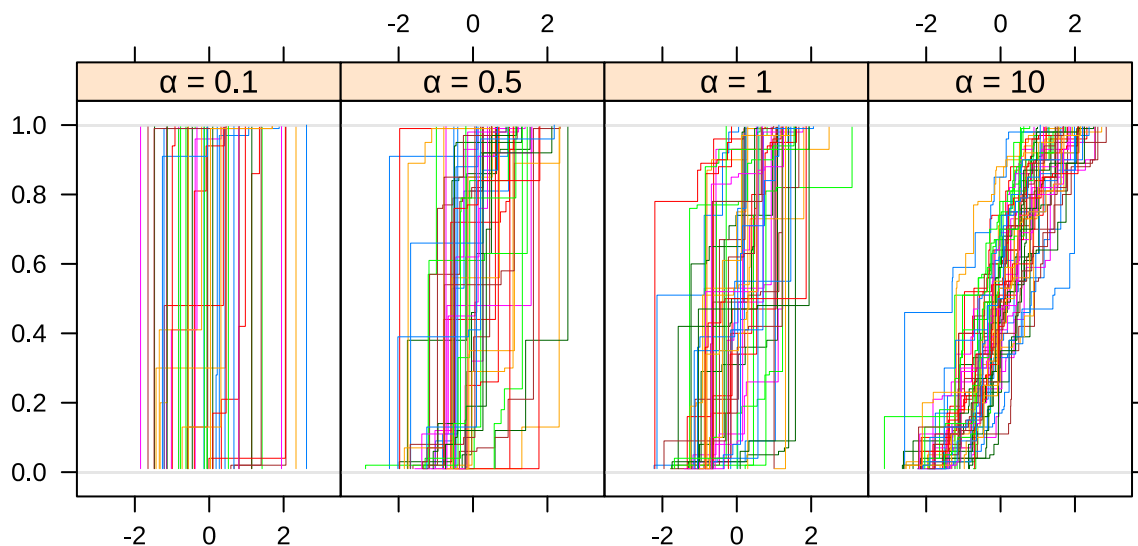


Figure 1: Prior c.d.f realizations of Dirichlet process with different base measures.

We can also get Monte Carlo estimates of the mean functional $\mu(G)$ and the variance functional $\sigma^2(G)$ from the prior realizations of G , drawn by assuming a truncation upto $K = 20$ terms.

We see in Figure 2, that naturally the mean functional $\mu(G)$ is centered around 0 since the base measure is centered around 0 but higher the value of α , more is the concentration of the mean functional around 0. In other words, the mean functional has higher kurtosis for higher α . On the other hand, the variance functional indicates that lower value of α indicates the lower dispersion of atoms in G , whereas, higher value of α indicates higher dispersion of atoms in G .

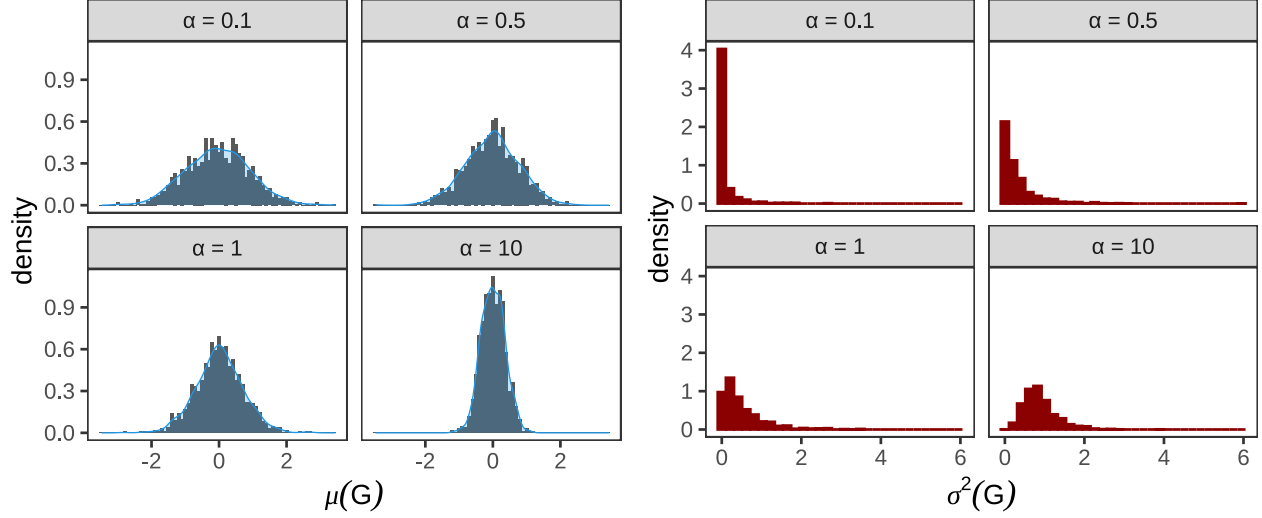


Figure 2: Prior distributions of mean and variance functionals from 1000 prior samples

1.2 Part 2

Following Weak Law of large Numbers (WLLN), we can estimate the expected number of nonempty clusters $E(M)$ by the mean number of unique atoms from each realization of c.d.f. sampled from the prior $DP(\alpha, \mathcal{N}(0, 1))$ with different values of α . Here, we have considered $\alpha = 0.1, 0.5, 1, 10$. Note that the sampled prior realizations are generated from an approximate prior achieved by truncation the infinite mixture to a mixture of 20 atoms. The red line in Figure 3 denotes the theoretical expected number of nonempty clusters as shown by Antoniak (1974) given by $E(M) = \alpha \log((\alpha + n)/\alpha)$. We see for smaller values of α , the approximation may have given reasonable estimates of number of non-empty clusters but I suspect that for a large value of α , the truncation approximation might not be reasonable.

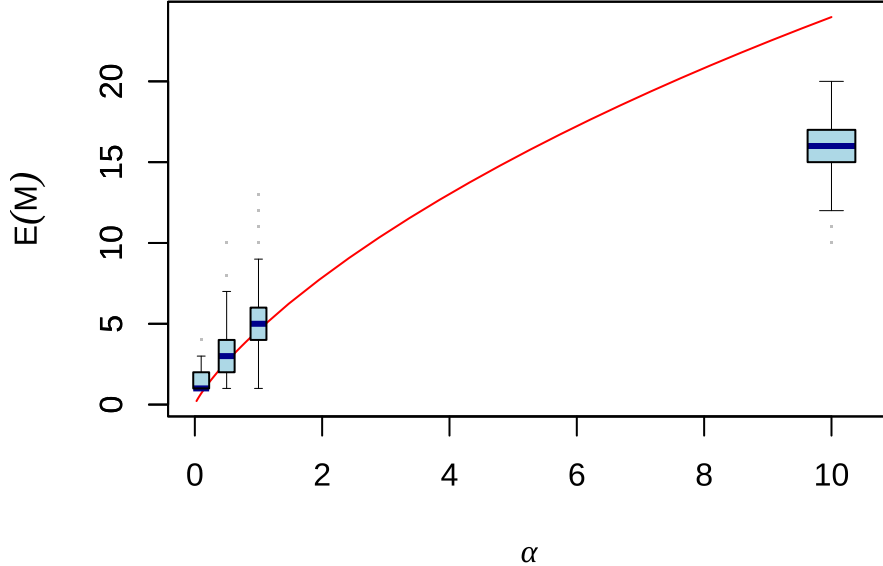


Figure 3: Comparison of theoretical expected number of non-empty clusters and estimated with finite truncation approximation based on 1000 samples

1.3 Part 4

We know from the conjugacy of a Dirichlet process, if we have the data generation model as $X_1, X_2, \dots, X_n \stackrel{\text{ind}}{\sim} G$ with $G \sim \text{DP}(\alpha G_0)$, $G_0 = \mathcal{N}(0, 1)$, then we can write the posterior as again a Dirichlet process with a different base measure.

$$G \mid X_1, X_2, \dots, X_n \sim \text{DP}\left(\alpha + n, \frac{1}{\alpha + n} \sum_{i=1}^n \delta_{X_i} + \frac{\alpha}{\alpha + n} G_0\right)$$

Hence, from Figure 4, we see that with increase in α , we see more posterior samples (in gray) tending towards the prior base measure (given in red dashed line) and for lesser values of α , we have more attenuation towards the empirical cdf of the given data (drawn in solid black line).

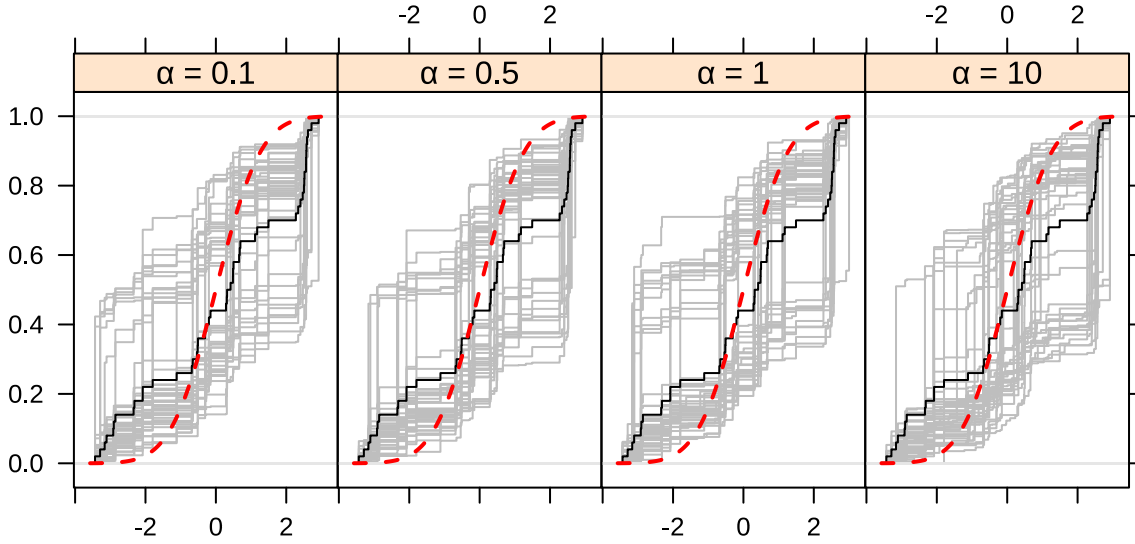


Figure 4: Posterior samples (gray) of the generative random measure compared with the empirical c.d.f. (black) and the original base measure of the DP prior (red).

References

- Antoniak, Charles E. 1974. “Mixtures of Dirichlet Processes with Applications to Bayesian Nonparametric Problems.” *The Annals of Statistics* 2 (6): 1152–74. <http://www.jstor.org/stable/2958336>.
- Ishwaran, Hemant, and Mahmoud Zarepour. 2002. “Exact and Approximate Sum Representations for the Dirichlet Process.” *Canadian Journal of Statistics* 30 (2): 269–83. <https://doi.org/https://doi.org/10.2307/3315951>.