

Parsimonious Hierarchical Modeling Using Repulsive Distributions

José J. Quinlan[†], Garritt L. Page[‡] and Fernando A. Quintana[†]

[†] Pontificia Universidad Católica de Chile

[‡] Brigham Young University

Abstract

The use of nonparametric mixtures for density estimation has become routine in statistical practice, the Dirichlet Process Mixture (DPM) model being one of its most popular versions. The popularity of models such as the DPM is largely based on their flexibility and tractability. A common problem of fitting DPM-like models to data though is that they tend to produce a large number of (sometimes) redundant clusters. In this work we propose a method that is able to produce parsimonious mixture models (i.e. involving fewer clusters than other alternatives), essentially without sacrificing flexibility or model fit. This method is based on the idea of repulsion, that is, that any two mixture components are encouraged to be well separated. We propose a family of multivariate probability densities dominated by the n -fold Lebesgue measure on the measurable space $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ whose d dimensional coordinates tend to repel each other in a smooth way. The induced probability measure, called $\text{Rep}_{n,d}(f_0, C_0, \rho)$, has a close relation with Gibbs measures, Graph theory and Point Processes. We investigate the global properties of the $\text{Rep}_{n,d}(f_0, C_0, \rho)$ class. We also explore their use in the context of mixture models for density estimation, and discuss computational implementation. Finally, we illustrate this method with some well-known data sets.

Key Words: Mixture Models, Gibbs measures, Graph Theory, Point Processes.