BNPlib for density estimation:

A nonparametric C++ library (part 3)

> Bruno Guindani Elena Zazzetti

February 19th, 2020



https://github.com/poliprojects/BNPlib

Model

DP and DPM models

Having observed the iid sample $\{y_i\}_i$, i = 1, ..., n:

• Dirichlet process model (discrete):

$$y_i|G \stackrel{\mathsf{iid}}{\sim} G$$

$$G \sim DP(MG_0)$$

Dirichlet process mixture (DPM) model (continuous):

$$y_i|G \stackrel{\text{iid}}{\sim} f_G(\cdot) = \int_{\Theta} f(\cdot|\boldsymbol{\vartheta}) G(d\boldsymbol{\vartheta})$$

 $G \sim DP(MG_0)$

Equivalent formulations (1)

• (DPM) is equivalent to:

$$\begin{aligned} y_i | \vartheta_i &\overset{\text{ind}}{\sim} f(\cdot | \vartheta_i), \quad i = 1, \dots, n \\ \vartheta_i | G &\overset{\text{iid}}{\sim} G, \quad i = 1, \dots, n \\ G &\sim DP(MG_0) \end{aligned}$$

• State: ϑ_i latent variables (discrete)

Equivalent formulations (2)

• (DPM) is also equivalent to:

$$y_i|\boldsymbol{\phi}, c_i \stackrel{\text{ind}}{\sim} f(\cdot|\phi_{c_i}), \quad i = 1, \dots, n$$

$$c_i|\mathbf{p} \stackrel{\text{iid}}{\sim} \sum_{j=1}^K p_j \delta_j(\cdot), \quad i = 1, \dots, n$$

$$\phi_c \stackrel{\text{iid}}{\sim} G_0, \quad c = 1, \dots, k$$

$$\mathbf{p} \sim \text{Dir}(M/K, \dots, M/K)$$

$$K \to +\infty$$

- State: c_i allocations to clusters
- ullet State: ϕ_{c_i} unique values for each cluster
- ullet Only the finitely many ϕ_c used are kept track of

Case study

• (DPM) with Normal Normal-InverseGamma (NNIG) prior:

$$y_i | \vartheta_i \overset{\text{ind}}{\sim} f(\cdot | \vartheta_i), \quad i = 1, \dots, n$$

 $\vartheta_i | G \overset{\text{iid}}{\sim} G, \quad i = 1, \dots, n$
 $G \sim DP(MG_0)$

$$\begin{split} f(y|\boldsymbol{\vartheta}) &= N(y|\mu,\sigma^2), \\ G_0(\boldsymbol{\vartheta}|\mu_0,\lambda_0,\alpha_0,\beta_0) &= N\left(\mu|\mu_0,\frac{\sigma^2}{\lambda_0}\right) \times \mathsf{Inv-Gamma}(\sigma^2|\alpha_0,\beta_0) \end{split}$$

ullet State: $oldsymbol{artheta}=(\mu,\sigma)$,

◄□▶◀圖▶◀불▶◀불▶ 불 쒸٩

Algorithms

General structure

```
void step(){
    sample_allocations();
    sample_unique_values();
}
void run(){
    initialize();
    unsigned int iter = 0;
    while(iter < maxiter){</pre>
        step();
        if(iter >= burnin){
             save_iteration(iter);
        }
        iter++;
```

aaa

aaa

Applications

Cluster estimation

unsigned int cluster_estimate();

$$\hat{k} = \arg\min_{k} \|D^{(k)} - \bar{D}\|_{F}^{2} = \arg\min_{k} \sum_{i,j} (D_{ij}^{(k)} - \bar{D}_{ij})^{2}$$

Density estimation

void eval_density(const std::vector<double> grid);

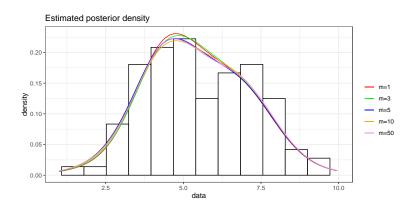
$$\hat{f}^{(k)}(x) = \sum_{j} \frac{n_{j}^{(k)}}{M+n} f\left(x|\phi_{j}^{(k)}\right) + \frac{M}{M+n} m(x)$$

$$\hat{m}(x) = \frac{1}{m} \sum_{h=0}^{m-1} f\left(x|\phi_{h}\right)$$

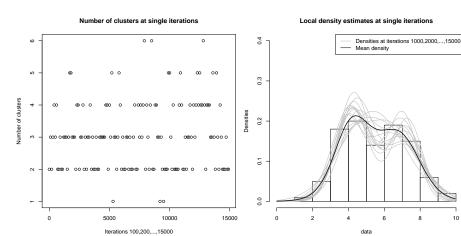
$$\Longrightarrow \hat{f}(x) = \frac{1}{K} \sum_{h=0}^{m-1} \hat{f}^{(k)}(x)$$

Results

Auxiliary parameters

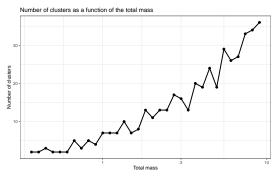


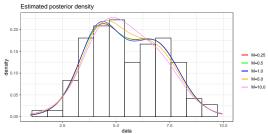
Oscillations



10

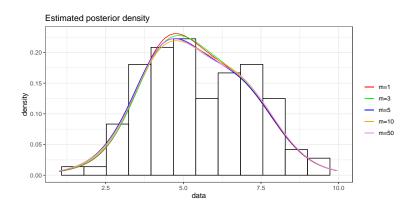
Total mass



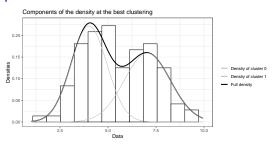


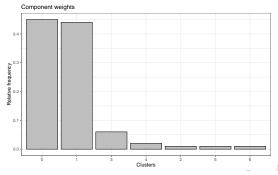
16 / 20

Auxiliary parameters

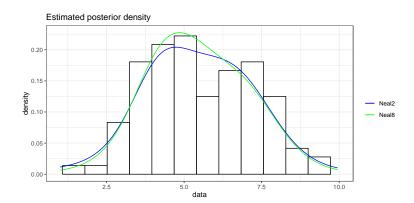


Density components





Neal2 vs Neal8



Bibliography

- 🦫 Muller, Quintana, Bayesian Nonparametric Data Analysis
- Neal (2000), Markov Chain Sampling Methods for Dirichlet Process Mixture Models
- 医 Ishwaran, James (2001), Gibbs Sampling Methods for Stick-Breaking Priors
- Nurphy (2007), Conjugate Bayesian analysis of the Gaussian distribution
- Protocol Buffers: https://developers.google.com/protocol-buffers/ docs/cpptutorial
- Stan: http://mc-stan.org/math
- Eigen: https://eigen.tuxfamily.org/dox
- GitHub codes of Mario Beraha and Riccardo Corradin for similar projects
- Course material for Bayesian Statistics: https://beep.metid.polimi.it/web/2019-20-bayesian-statistics-alessandra-guglielmi-/