

Clustering for Gritches

Shuheï Mano

The Institute of Statistical Mathematics, Japan

Data characterization meeting

November 25, 2013

Harmonic Analysis

Schuster, A. (1897) On Lunar and Solar Periodicities of Earthquakes.
Proc Roy Soc Lond 61: 455-465.

Number of events at periods $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N$, where $\mathbf{N} = 2q + 1$.

$$\mathbf{z}_t = \mathbf{a}_0 + \sum_{i=1}^q (\mathbf{a}_i \cos 2\pi f_i t + \mathbf{b}_i \sin 2\pi f_i t) + \epsilon_t, \quad f_i = \frac{i}{N}.$$

Least Squares Estimates:

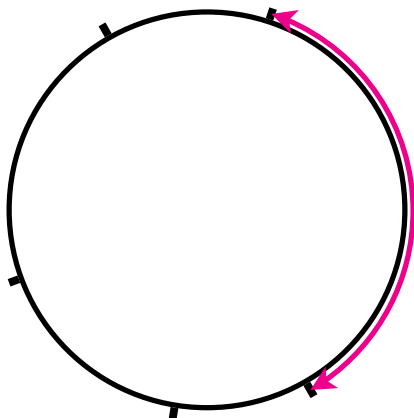
$$\hat{\mathbf{a}}_i = \frac{2}{N} \sum_{t=1}^N \mathbf{z}_t \cos 2\pi f_i t, \quad \hat{\mathbf{b}}_i = \frac{2}{N} \sum_{t=1}^N \mathbf{z}_t \sin 2\pi f_i t.$$

Periodgram:

$$I_i = \frac{N}{2} (\hat{\mathbf{a}}_i^2 + \hat{\mathbf{b}}_i^2), \quad i = 1, 2, \dots, q.$$

Extreme Sizes in Random Partition

Fisher, R.A. (1929) Tests of Significance in Harmonic Analysis. Proc Roy Soc Lond 125: 54-59.



The test based on the distribution of $X_{(1)}$ in $(X_1, X_2, \dots, X_n) \sim \text{Dir}(1)$.

Dirichlet Process

Definition (Ferguson, 1973)

Let μ be a finite measure on $(\mathcal{X}, \mathcal{B})$. A random measure \mathbf{D} on \mathcal{X} is called a Dirichlet process if for every finite measurable partition $\{\mathbf{B}_1, \dots, \mathbf{B}_k\}$ of \mathcal{X} , $(\mathbf{D}(\mathbf{B}_1), \dots, \mathbf{D}(\mathbf{B}_k)) \sim \text{Dir}(\mu(\mathbf{B}_1), \dots, \mu(\mathbf{B}_k))$.

Theorem (Ferguson, 1973)

1. \mathbf{Y}_i , i.i.d. with a probability measure $\mu(\cdot)/\theta$ with $\theta = \mu(\mathcal{X})$.
2. Let ρ_t be the gamma process with $\rho_0 = \mathbf{0}$ whose increment follows $\text{Gamma}(\theta t, 1)$ and $(\mathbf{Z}_1, \mathbf{Z}_2, \dots)$ be the jump sizes up to $t = 1$.
3. $\mathbf{D}(\cdot) = \sum_{i=1}^{\infty} \frac{\mathbf{Z}_i}{\rho_1} \delta_{\mathbf{Y}_i}(\cdot)$ is a Dirichlet process with parameter μ .

Poisson-Dirichlet Distribution

Definition (Kingman, 1978; Pitman, 1995)

Let

$$P_1 = W_1, \quad P_i = W_i \prod_{j=1}^{i-1} (1 - W_j), \quad i = 2, 3, \dots,$$

where $W_i \sim \text{Beta}(1 - \alpha, \theta + i\alpha)$, i.i.d. with $0 \leq \alpha < 1$ and $\theta > -\alpha$, or $\alpha < 0$ and $\theta = -\alpha m$, $m \in \mathbb{N}$. The distribution of the ranked sequence of \mathbf{P} is called the 2-parameter Poisson-Dirichlet distribution $\mathbf{PD}(\alpha, \theta)$.

Remark

- ▶ For $\alpha = 0$, $(P_i) \stackrel{d}{=} (\rho_1^{-1} Z_i)$ and $\mathbf{D}(\cdot) = \sum_{i=1}^{\infty} P_i \delta_{Y_i}(\cdot)$ is the 2-parameter generalization of the Dirichlet process.
- ▶ For $\alpha < 0$, $\mathbf{PD}(\alpha, -\alpha m)$ is the m -dimensional symmetric $\mathbf{Dir}(-\alpha)$.
- ▶ Most general such that \mathbf{P} is invariant under size-biased permutation.

Ewens-Pitman Random Partition

A partition of $n \in \mathbb{N}$ by integers is identified by multiplicity \mathbf{c} (size index) such that

$$\|\mathbf{s}\| := \sum_{i=1}^n s_i = k_n, \quad |\mathbf{s}| := \sum_{i=1}^n i s_i = n.$$

Example

$10 = 5 + 2 + 2 + 1$ gives $k_{10} = 4$, $s_5 = s_1 = 1$, $s_2 = 2$.

Definition (Ewens,1972; Pitman,1992)

An exchangeable random partition

$$P(\mathbf{S} = \mathbf{s}, K_n = k) = \frac{\left(\frac{\theta}{\alpha}\right)_k}{(\theta)_n} (-1)^{n-k} n! \prod_{j=1}^n \binom{\alpha}{j}^{s_j} \frac{1}{s_j!},$$

where $0 \leq \alpha < 1$ and $\theta > -\alpha$, or $\alpha < 0$ and $\theta = -\alpha m$, $m = 1, 2, \dots$

Chinese restaurant process

The Ewens-Pitman random partition is the **sampling distribution** from $PD(\alpha, \theta)$. Deonote the ranked sizes $\{i : C_i > 0\}$ by $L_{(1)}^{(n)}, L_{(2)}^{(n)}, \dots$

$$n^{-1}(L_{(1)}^{(n)}, L_{(2)}^{(n)}, \dots) \xrightarrow{d} (P_{(1)}, P_{(2)}, \dots), \quad n \rightarrow \infty.$$

Suppose n person occupy k tables (cluster) and n_i people sit at table i .
The next person

- ▶ sits at an empty table with probability $\frac{\theta + k\alpha}{\theta + n}$,
- ▶ sits at the table i with probability $\frac{n_i - \alpha}{\theta + n}$.

The random partition is numbers of people at each table.

Dirichlet Process Mixture Distribution

Assume σ^2 and σ_0 are known for simplicity. A gaussian mixture distribution is

$$f(\mathbf{x}|\mathbf{P}, \mu) = \sum_{i=1}^m P_i \phi(\mathbf{x}|\mu_i, \sigma^2)$$

with $\mu_i \sim \mathbf{N}(\mathbf{0}, \sigma_0^2)$, $P_i \sim \mathbf{Dir}(-\alpha)$.

How to choose the number of clusters?

- ▶ Try each and choose optimal one by some criteria (information criterion, cross validation,...)
- ▶ Dirichlet process with $\alpha \geq 0$: $\mathbf{D} \sim \mathbf{DP}(\alpha, \theta; \mathbf{N}(\mathbf{0}, \sigma_0^2))$.

Assignment

By Bayes' rule,

$$\begin{aligned} P(\mathbf{C}_i | \mathbf{C}_{-i}, \mathbf{X}) &\propto P(X_i | \mathbf{C}, \mathbf{X}_{-i}) P(\mathbf{C}_i | \mathbf{C}_{-i}, \mathbf{X}_{-i}) \\ &\hat{=} \phi(X_i | \hat{\mu}_i(\mathbf{C}_{-i}, \mathbf{X}_{-i}), \sigma^2) P(\mathbf{C}_i | \mathbf{C}_{-i}) \end{aligned}$$

The second factor is given by CRP. Sampling from $P(\mathbf{C} | \mathbf{X})$ is possible by using the Gibbs sampler.

For a gritch clustering, \mathbf{X} is not a position but a wave form. Constructing the likelihood is the key issue.