

Bayesian Nonparametric Density Estimation

version: 2017-11-16 · 19:32:10

Basic Idea

► Parametric Density Estimation:

$$Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} F(y | \theta), \quad \theta \sim \pi(\theta).$$

Bayesian inference is on $p(\theta | \mathbf{Y}) \propto \mathcal{L}(\theta | \mathbf{Y}) \times \pi(\theta)$.

► Nonparametric Density Estimation:

$$Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} F(y), \quad F \sim \pi(F).$$

- $\pi(F)$ is a distribution on CDFs.
- Bayesian inference is on $p(F | \mathbf{Y}) \propto \mathcal{L}(F | \mathbf{Y}) \times \pi(F)$.
- **Convenient Prior:** A **conjugate** prior for $L(F | \mathbf{Y})$ is the **Dirichlet Process**, $F \sim \text{DP}(F_0, \alpha)$.

The Dirichlet Process

Dirichlet Distribution:

- ▶ Let $X_k \stackrel{\text{ind}}{\sim} \text{Gamma}(\alpha \rho_k, \beta)$ for $k = 1, \dots, K$, where $\alpha, \rho_k > 0$ and $\sum_{k=1}^K \rho_k = 1$. Then

$$\mathbf{Y} = \left(\frac{X_1}{\sum_{i=1}^K X_i}, \dots, \frac{X_K}{\sum_{i=1}^K X_i} \right) \sim \text{Dirichlet}(\boldsymbol{\rho}, \alpha).$$

- ▶ \mathbf{Y} is a probability vector: $Y_k > 0$ and $\sum_{k=1}^K Y_k = 1$.
- ▶ Mean and variance:

$$E[\mathbf{Y}] = \boldsymbol{\rho}, \quad \text{var}(\mathbf{Y}) = \frac{\text{diag}(\boldsymbol{\rho}) - \boldsymbol{\rho}\boldsymbol{\rho}'}{\alpha + 1}.$$

The Dirichlet Process

Dirichlet Process:

- **Notation:** Let $F(y)$ be an arbitrary CDF and $B \subseteq \mathbb{R}$. Then $F(B) := \Pr(Y \in B)$, where $Y \sim F(y)$.
- **Definition:** Let F_0 be a CDF and $\alpha > 0$. Then $F \sim \text{DP}(F_0, \alpha)$ is said to follow a **Dirichlet Process** if for any finite partition $B_1 \amalg B_2 \amalg \dots \amalg B_K = \mathbb{R}$,

$$(F(B_1), \dots, F(B_K)) \sim \text{Dirichlet}(\boldsymbol{\rho}, \alpha), \quad \boldsymbol{\rho} = (F_0(B_1), \dots, F_0(B_K)).$$

The Dirichlet Process

- ▶ **DP:** $F \sim \text{DP}(F_0, \alpha) \iff F(\mathbf{B}) \sim \text{Dirichlet}(F_0(\mathbf{B}), \alpha), \forall \mathbf{B} = \coprod_{k=1}^K B_k = \mathbb{R}.$
- ▶ **Representation:** Turns out that $F \sim \text{DP}(F_0, \alpha)$ is a discrete distribution with countably many atoms:

$$F(y) = \sum_{k=1}^{\infty} w_k \delta_{y_k}(y).$$

- ▶ **Sampling:** Draw $F \sim \text{DP}(F_0, \alpha)$ by **stick-breaking** procedure:
 1. Draw $y_1, y_2, \dots \stackrel{\text{iid}}{\sim} F_0(y).$
 2. Draw $\beta_1, \beta_2, \dots \stackrel{\text{iid}}{\sim} \text{Beta}(1, \alpha)$ and let $w_k = \beta_k \prod_{i=1}^{k-1} (1 - \beta_i).$
 3. $F(y) = \sum_{k=1}^{\infty} w_k \delta_{y_k}(y)$ is a draw from $\text{DP}(F_0, \alpha).$

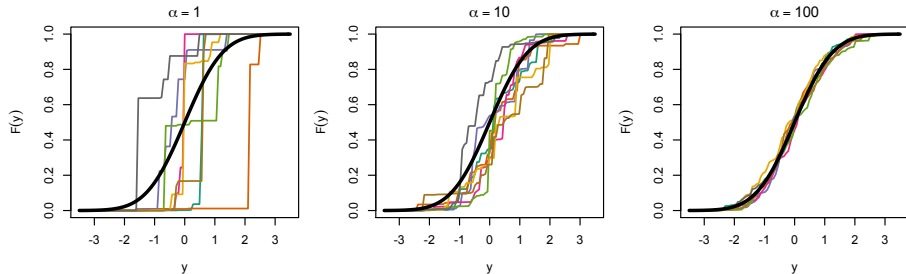
The Dirichlet Process

- ▶ **DP:** $F(y) = \sum_{k=1}^{\infty} w_k \delta_{y_k}(y) \sim \text{DP}(F_0, \alpha) \iff F(\mathbf{B}) \sim \text{Dirichlet}(F_0(\mathbf{B}), \alpha)$.
- ▶ **Sampling:** Draw $F \sim \text{DP}(F_0, \alpha)$ by **stick-breaking** procedure:
 1. Draw $y_1, y_2, \dots \stackrel{\text{iid}}{\sim} F_0(y)$.
 2. Draw $\beta_1, \beta_2, \dots \stackrel{\text{iid}}{\sim} \text{Beta}(1, \alpha)$ and let $w_k = \beta_k \prod_{i=1}^{k-1} (1 - \beta_i)$.
 3. $F(y) = \sum_{k=1}^{\infty} w_k \delta_{y_k}(y)$ is a draw from $\text{DP}(F_0, \alpha)$.
- ▶ **In practice:**
 - ▶ Can't do this exactly because we can't store an infinite sequence in memory.
 - ▶ Instead, draw $F(y) = \sum_{k=1}^K w_k \delta_{y_k}(y)$, where K is predetermined by memory allocation, or e.g., note that $w_1 > w_2 > \dots$, and

$$E[w_k] = \alpha^{k-1} / (1 + \alpha)^k,$$

and use this to bound expectation as a function of K . (or use a `while`-loop and stop when $1 - \sum_{k=1}^K w_k < \epsilon$, if dynamic memory allocation is not a concern.)

Example



$F_1, \dots, F_8 \stackrel{\text{iid}}{\sim} \text{DP}\{\mathcal{N}(0, 1), \alpha\}$ for different values of α .

The Dirichlet Process

- ▶ **DP:** $F(y) = \sum_{k=1}^{\infty} w_k \delta_{y_k}(y) \sim \text{DP}(F_0, \alpha) \iff F(\mathbf{B}) \sim \text{Dirichlet}(F_0(\mathbf{B}), \alpha)$.
- ▶ **CDF Sampling:** Draw $F(y) \approx \sum_{k=1}^K w_k \delta_{y_k}(y)$ by **stick-breaking** procedure.
- ▶ **Marginal Sampling:** Can **exactly** sample Y_1, \dots, Y_n from

$$\begin{aligned} Y_1, \dots, Y_n &\stackrel{\text{iid}}{\sim} F(y), \quad F \sim \text{DP}(F_0, \alpha) \\ \iff Y_1, \dots, Y_n &\stackrel{\text{iid}}{\sim} \int F(y) \times \pi(F | F_0, \alpha) d\{F\} \end{aligned}$$

by **Chinese Restaurant Process**:

1. Draw $Y_1 \sim F_0(y)$.
2. Let $\hat{F}_i(y)$ denote the empirical distribution of Y_1, \dots, Y_i . Draw

$$Y_{i+1} \sim \frac{\alpha}{\alpha + i} F_0(y) + \frac{i}{\alpha + i} \hat{F}_i(y).$$

The Y_i drawn this way are not iid. However, they are exchangeable, i.e., order in which we draw doesn't matter.

The Dirichlet Process

- ▶ **DP:** $F(y) = \sum_{k=1}^{\infty} w_k \delta_{y_k}(y) \sim \text{DP}(F_0, \alpha) \iff F(\mathbf{B}) \sim \text{Dirichlet}(F_0(\mathbf{B}), \alpha)$.
- ▶ **CDF Sampling:** Draw $F(y) \approx \sum_{k=1}^K w_k \delta_{y_k}(y)$ by **stick-breaking** procedure.
- ▶ **Marginal Sampling:**

Can **exactly** sample $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} F(y), \quad F \sim \text{DP}(F_0, \alpha)$

by **Chinese Restaurant Process**:

1. Draw $Y_1 \sim F_0(y)$.
2. Let $\hat{F}_i(y)$ denote the empirical distribution of Y_1, \dots, Y_i . Draw

$$Y_{i+1} \sim \frac{\alpha}{\alpha + i} F_0(y) + \frac{i}{\alpha + i} \hat{F}_i(y).$$

CRP analogy: Customer $i + 1$ enters restaurant, sits at new table with probability $\alpha/(\alpha + i)$ and orders dish $Y_{i+1} \sim F_0(y)$. Otherwise, randomly chooses among existing tables proportionally to how many people are sitting there, and eats whatever dish is at the table.

The Dirichlet Process

► **DP:** $F(y) = \sum_{k=1}^{\infty} w_k \delta_{y_k}(y) \sim \text{DP}(F_0, \alpha) \iff F(\mathbf{B}) \sim \text{Dirichlet}(F_0(\mathbf{B}), \alpha).$

► **Marginal Sampling:** Draw $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} F(y), \quad F \sim \text{DP}(F_0, \alpha)$

with **CRP**:

► $Y_{i+1} \sim \frac{\alpha}{\alpha+i} F_0(y) + \frac{i}{\alpha+i} \hat{F}_i(y),$ where $\hat{F}_i(y)$ is the ECDF of $Y_1, \dots, Y_i.$

► **Marginal Distribution:**

► Assume that the PDF $f_0(y)$ exists.

► Let $\tilde{Y}_1, \dots, \tilde{Y}_K$ denote the unique values of $\mathbf{Y} = (Y_1, \dots, Y_n),$ and $\mathbf{n} = (n_1, \dots, n_K)$ denote number of Y_i having each value.

► The marginal distribution is
$$p(\mathbf{Y} | F_0, \alpha) = \frac{\alpha^K \prod_{k=1}^K f_0(\tilde{Y}_k)}{\prod_{i=1}^n (\alpha + i)} \times g(\mathbf{n}),$$

where $g(\mathbf{n})$ doesn't depend on F_0 or $\alpha.$ (combinatorics result on which observations have ties, but interestingly doesn't depend on exact order of ties.)

Bayesian Inference

► Model and Prior:

$$Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} F(y)$$

$$F(y) \sim \text{DP}\{F_0(y | \boldsymbol{\theta}), \alpha\}$$

$$(\boldsymbol{\theta}, \alpha) \sim \pi(\boldsymbol{\theta}, \alpha).$$

► Conditional Distribution:

$$F(y) | \boldsymbol{\theta}, \alpha, \mathbf{Y} \sim \text{DP} \left\{ \frac{\alpha}{\alpha + n} F_0(y | \boldsymbol{\theta}) + \frac{n}{\alpha + n} \underbrace{\hat{F}_n(y)}_{\text{ECDF of } \mathbf{Y}}, \alpha + n \right\}.$$

- **Marginal Likelihood:** Assume PDF $f_0(y | \boldsymbol{\theta})$ exists and let $\tilde{Y}_1, \dots, \tilde{Y}_K$ denote the unique values of \mathbf{Y} . Then

$$\mathcal{L}(\boldsymbol{\theta}, \alpha | \mathbf{Y}) \propto \frac{\alpha^K \Gamma(\alpha)}{\Gamma(\alpha + n)} \prod_{k=1}^K f_0(\tilde{Y}_k | \boldsymbol{\theta}).$$

(Note that $\prod_{i=1}^n (\alpha + i) = \Gamma(\alpha + n) / \Gamma(\alpha)$.)

Dirichlet Process Mixture Model:

- DPM Model:

$$\begin{aligned}Y_i | \theta_i &\stackrel{\text{iid}}{\sim} f(y | \theta_i) \\ \theta_1, \dots, \theta_n &\stackrel{\text{iid}}{\sim} G(\theta) \\ G(\theta) &\sim \text{DP}\{G_0(\theta | \eta), \alpha\}\end{aligned}$$

- By writing $G(\theta) = \sum_{k=1}^{\infty} w_k \delta_{\theta_k}(\theta)$, can view

$$Y | G \sim \sum_{k=1}^{\infty} w_k f(y | \theta_k)$$

as an infinite-component mixture model.

Dirichlet Process Mixture Model

- DPM Model:

$$\begin{aligned}Y_i | \theta_i &\stackrel{\text{ind}}{\sim} f(y | \theta_i) \\ \theta_1, \dots, \theta_n &\stackrel{\text{iid}}{\sim} G(\theta) \\ G(\theta) &\sim \text{DP}\{G_0(\theta | \eta), \alpha\}\end{aligned}$$

- Clustering:

- $\Theta = (\theta_1, \dots, \theta_n)$ takes on $1 \leq K \leq n$ distinct values $\tilde{\theta}_1, \dots, \tilde{\theta}_K$

- *Cluster allocation*:

$$\mathcal{C} = \mathcal{C}(\mathbf{Y}) = \Pi_{k=1}^K \mathcal{S}_k, \quad \mathcal{S}_k = \{Y_i : \theta_i = \tilde{\theta}_k\}.$$

\implies don't need to prespecify number of clusters.

- *Cluster Probability*:

$$\Pr(\text{cluster allocation is } \mathcal{C}_0 | \mathbf{Y}, \eta, \alpha) = \Pr(\mathcal{C}(\mathbf{Y}) = \mathcal{C}_0 | \mathbf{Y}, \eta, \alpha)$$

Dirichlet Process Mixture Model

► **DPM Model:**

$$\begin{aligned}Y_i | \theta_i &\stackrel{\text{iid}}{\sim} f(y | \theta_i) \\ \theta_1, \dots, \theta_n &\stackrel{\text{iid}}{\sim} G(\theta) \\ G(\theta) &\sim \text{DP}\{G_0(\theta | \eta), \alpha\}\end{aligned}$$

- **Mixing Kernel:** The most common choice is $Y | \mu, \sigma \sim \mathcal{N}(\mu, \sigma^2)$. But to simplify calculations, let

$$Y \sim \text{NEF}(\theta) \quad \Longleftrightarrow \quad \begin{aligned}f(y | \theta) &= \exp\{\mathbf{T}'\theta - \Phi(\theta) + h(y)\} \\ \mathbf{T} &= \mathbf{T}(y)\end{aligned}$$

be a **Natural Exponential Family**. (Normal is an NEF but using non-standard parametrization.)

Bayesian Inference

► Model and Prior:

$$\begin{aligned}Y_i | \theta_i &\stackrel{\text{iid}}{\sim} f(y | \theta_i) \\ \theta_1, \dots, \theta_n &\stackrel{\text{iid}}{\sim} G(\theta) \\ G(\theta) &\sim \text{DP}\{G_0(\theta | \eta), \alpha\} \\ (\eta, \alpha) &\sim \pi(\eta, \alpha).\end{aligned}$$

► **Mixing Kernel:** $Y \sim \text{NEF}(\theta)$. Usually $Y \sim \mathcal{N}(\mu, \sigma^2)$ with $\theta = (\mu, \sigma^2)$.

► MCMC Sampling:

- $p(\eta, \alpha | \Theta, \mathbf{Y}) = p(\eta, \alpha | \Theta)$, where $\Theta = (\theta_1, \dots, \theta_n)$, and this is just regular inference for DP.
- So only need to draw $p(\Theta | \eta, \alpha, \mathbf{Y})$ to implement a Gibbs sampler.

Bayesian Inference

► DPM Model and Prior:

$$\begin{aligned} Y_i | \theta_i &\stackrel{\text{ind}}{\sim} f(y | \theta_i) & G(\theta) &\sim \text{DP}\{G_0(\theta | \eta), \alpha\} \\ \theta_1, \dots, \theta_n &\stackrel{\text{iid}}{\sim} G(\theta) & (\eta, \alpha) &\sim \pi(\eta, \alpha). \end{aligned}$$

► Prior distribution: Componentwise we have

$$\theta_i | \Theta_{-i}, \eta, \alpha \sim \frac{\alpha}{\alpha + n - 1} G_0(\theta | \eta) + \frac{1}{\alpha + n - 1} \sum_{j \neq i} \delta_{\theta_j}(\theta).$$

(Recall that θ_i from CRP are exchangeable, i.e., order of sampling doesn't matter.)

Bayesian Inference

► DPM Model and Prior:

$$Y_i | \theta_i \stackrel{\text{iid}}{\sim} f(y | \theta_i) \quad \theta_1, \dots, \theta_n \stackrel{\text{iid}}{\sim} G(\theta) \quad G(\theta) \sim \text{DP}\{G_0(\theta | \boldsymbol{\eta}), \alpha\} \quad (\boldsymbol{\eta}, \alpha) \sim \pi(\boldsymbol{\eta}, \alpha).$$

► Prior distribution:

$$\theta_i | \boldsymbol{\Theta}_{-i}, \boldsymbol{\eta}, \alpha \sim \frac{\alpha}{\alpha + n - 1} G_0(\theta | \boldsymbol{\eta}) + \frac{1}{\alpha + n - 1} \sum_{j \neq i} \delta_{\theta_j}(\theta).$$

► Posterior distribution:

$$\theta_i | \boldsymbol{\Theta}_{-i}, \boldsymbol{\eta}, \alpha, \mathbf{Y} \sim r \cdot p(\theta_i | \boldsymbol{\eta}, \alpha, Y_i) + \sum_{j \neq i} q_j \cdot \delta_{\theta_j}(\theta), \quad \text{where}$$

$$q_j \propto \frac{f(Y_i | \theta_j)}{\alpha + n - 1}$$

(for each $\theta_j \in \boldsymbol{\Theta}_{-i}$: posterior \propto prior \times likelihood)

$$r \propto \int f(Y_i | \theta) \frac{\alpha \cdot g_0(\theta | \boldsymbol{\eta})}{\alpha + n - 1} d\theta$$

$$\left(\begin{array}{l} r = \Pr(\theta_i \notin \boldsymbol{\Theta}_{-i} | \boldsymbol{\eta}, \alpha, \mathbf{Y}) \\ \propto p(Y_i | \boldsymbol{\eta}) \times \alpha / (\alpha + n - 1) \\ = \frac{\alpha}{\alpha + n - 1} \int f(Y_i | \theta) g_0(\theta | \boldsymbol{\eta}) d\theta \end{array} \right)$$

$$r + \sum_{j \neq i} q_j = 1.$$

Bayesian Inference

► **DPM Posterior distribution:**

$$\theta_i | \Theta_{-i}, \eta, \alpha, \mathbf{Y} \sim r \cdot p(\theta_i | \eta, Y_i) + \sum_{j \neq i} q_j \cdot \delta_{\theta_j}(\theta),$$
$$\begin{aligned} q_j &\propto f(Y_i | \theta_j) \\ r &\propto \alpha \int f(Y_i | \theta) g_0(\theta | \eta) d\theta \\ r + \sum_{j \neq i} q_j &= 1 \end{aligned}$$

► **Calculation:** Difficult in most cases. But when

$$Y | \theta \sim \text{NEF}(\theta) \iff f(y | \theta) = \exp\{\mathbf{T}'\theta - \Phi(\theta) + h(y)\}$$

and g_0 is the conjugate prior

$$\theta | \eta \sim \text{cEF}(\Upsilon, \nu) \iff g_0(\theta | \eta) = \exp\{\Upsilon'\theta - \nu\Phi(\theta) + q(\Upsilon, \nu)\},$$

then $p(\theta | \eta, Y) = g_0(\theta | \Upsilon + \mathbf{T}, \nu + 1)$, and

$$\int f(Y | \theta) g_0(\theta | \eta) d\theta = p(Y | \eta) = \frac{f(Y | \theta) g_0(\theta | \eta)}{p(\theta | \eta, Y)} = \frac{\exp\{h(Y) + q(\Upsilon, \nu)\}}{\exp\{q(\Upsilon + \mathbf{T}, \nu + 1)\}}.$$

$\implies r$ can also be calculated explicitly.