

A Gentle Introduction to Bayesian Nonparametrics



Nils Lid Hjort

Department of Mathematics, University of Oslo

Big Insight, 1/ii/17

Traditional Bayesian analysis: with data y from model with likelihood $L(\theta)$, the Bayes formula takes the **pre-data prior** $\pi(\theta)$ to the **post-data posterior**

$$\pi(\theta | \text{data}) = \frac{\pi(\theta)L(\theta)}{\int \pi(\theta')L(\theta') d\theta'}.$$

This requires well-defined densities (w.r.t. suitable measures), and in essence that θ is **finite-dimensional**.

Bayesian nonparametrics: about attempts to carry out such schemes, from pre-data to post-data, in **infinite- or very high-dimensional models**.

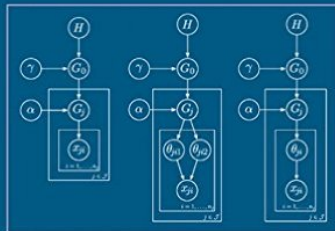
This is a **tall order**: conceptually; elicitation-wise; mathematically (distributions over very big spaces); operationally (there is no direct Bayes theorem); computationally (MCMC with 1000 parameters?).

Priors and posteriors for densities, **regression functions**, hazard rates, **big hierarchical models**, etc.

Plan & outline

- A What is it?
- B From the Dirichlet distribution to the Dirichlet process
- C The Beta process (e.g. for survival analysis)
- D Bernshteĭn–von Mises theorems (sometimes not)
- E Stationary time series
- F Bayesian nonparametrics for quantile analysis
- G Clustering models, ‘bigger things’
- H Other issues & concluding remarks

Cambridge Series in Statistical
and Probabilistic Mathematics



Bayesian Nonparametrics

Edited by Nils Lid Hjort, Chris Holmes,
Peter Müller and Stephen G. Walker

A: What is it?

Well, what is it^c ? Theorem: If the world is frequentist or Bayes, and parametric or nonparametric, then

$$IV = (I \cup II \cup III)^c.$$

| | | |
|---------------|-------------|-------|
| | frequentist | Bayes |
| parametric | I | II |
| nonparametric | III | IV |

I: Smallish finite models, estimation and inference for aspects of θ .

II: Smallish finite models, estimation and posterior inference, via prior $\pi(\theta)$ (this was all of Bayes inference, from c. 1774 to c. 1973).

III: Bigger models, density estimation, nonparametric regression, confidence bands, etc.

IV: Priors and posteriors for **random functions**, bigger structures, **hierarchies of hierarchies**, ...

Bayesian nonparametrics invites constructions for ‘approximately normal’, ‘approximately linear regression’, etc.

With ψ_1, ψ_2, \dots orthogonal functions on $[0, 1]$, like $\psi_j(u) = \sqrt{2} \cos(j\pi u)$, try

$$f(y) = f(y, \theta) \exp \left\{ \sum_{j=1}^{100} a_j \psi_j(F(y, \theta)) \right\} / c_{100}(a_1, \dots, a_{100}),$$

with a prior on θ along with $a_j \sim N(0, \tau^2/j^2)$. Data will tell us (via lots o’ MCMC) how close the real f is to the parametric start.

Approximately linear regression with approximately normal errors:

$$y_i = a + bx_i + \sum_{j=1}^{100} \gamma_j h_j(x_i) + \varepsilon_i \quad \text{for } i = 1, \dots, n,$$

with perhaps $\gamma_j \sim N(0, \kappa^2/j^2)$, where the ε_i are from $f \approx N(0, \sigma^2)$.

Fascinating and promising – but raises a long list of questions.

A de Finetti theorem

Why should we go for Bayesian nonparametrics? – Apart from ‘it works, and can solve big problems’, there’s a mathematical-probabilistic argument:

Consider y_1, y_2, \dots , a sequence of observations, with the **exchangeability property**:

$$(y_1, y_2, y_3, y_4, y_5) \sim (y_3, y_1, y_5, y_4, y_2),$$

i.e. have the same distribution – and similarly for all permutations, and all lengths. Then there is a **de Finetti measure** π , on the set of all distributions P , such that

$$\Pr\{y_1 \in A_1, \dots, y_n \in A_n\} = \int P(A_1) \cdots P(A_n) \pi(dP)$$

for all A_1, \dots, A_n , and all n .

So there is a **(nonparametric) prior** behind what you see (**whether you knew or not**), and y_1, y_2, \dots are **i.i.d. given P** .

B: The Dirichlet process: from finite to infinite

I begin with **two boxes** (1-or-0 measurements): With

$$y \sim \text{Bin}(n, p),$$

$$f(y | p) \propto p^y (1 - p)^{n-y},$$

and with $p \sim \text{Beta}(a, b)$,

$$p | y \propto p^{a-1} (1 - p)^{b-1} p^y (1 - p)^{n-y} = p^{a+y-1} (1 - p)^{b+n-y-1},$$

which means $p | \text{data} \sim \text{Beta}(a + y, b + n - y)$:

$$\hat{p} = \text{E}(p | y) = \frac{a + y}{a + b + n} = \frac{a + b}{a + b + n} p_0 + \frac{n}{a + b + n} \frac{y}{n},$$

$$\text{Var}(p | y) = \frac{1}{a + b + n} \hat{p}(1 - \hat{p}).$$

Thomas Bayes did this, with $(a, b) = (1, 1)$, i.e. a uniform prior – not in *Divine Benevolence, or an Attempt to Prove That the Principal End of the Divine Providence and Government is the Happiness of His Creatures* (1731), but in the other one (1763).

Then k boxes: from binomial to multinomial. With (y_1, \dots, y_k) the number of cases of types $1, \dots, k$, the likelihood is

$$\frac{n!}{y_1! \cdots y_k!} p_1^{y_1} \cdots p_{k-1}^{y_{k-1}} (1 - p_1 - \cdots - p_{k-1})^{y_k},$$

if the n trials are independent with the same probabilities p_1, \dots, p_k each time.

This calls on the Dirichlet distribution, $\text{Dir}(a_1, \dots, a_k)$:

$$\begin{aligned} \pi(p_1, \dots, p_{k-1}) &= \frac{\Gamma(a_1 + \cdots + a_k)}{\Gamma(a_1) \cdots \Gamma(a_k)} \\ &\quad \times p_1^{a_1-1} \cdots p_{k-1}^{a_{k-1}-1} (1 - p_1 - \cdots - p_{k-1})^{a_k-1} \end{aligned}$$

on the simplex of (p_1, \dots, p_{k-1}) .

Multiplying prior and likelihood:

$$(p_1, \dots, p_k) | (y_1, \dots, y_k) \sim \text{Dir}(a_1 + y_1, \dots, a_k + y_k).$$

So we can pass from prior $\text{Dir}(a_1, \dots, a_k)$ to posterior $\text{Dir}(a_1 + y_1, \dots, a_k + y_k)$ by simply adding the observed counts for the k boxes. We have

$$\hat{p}_j = E(p_j | \text{data}) = \frac{a_j + y_j}{a + n},$$
$$\hat{\sigma}_j^2 = \text{Var}(p_j | \text{data}) = \frac{1}{a + n + 1} \hat{p}_j (1 - \hat{p}_j),$$

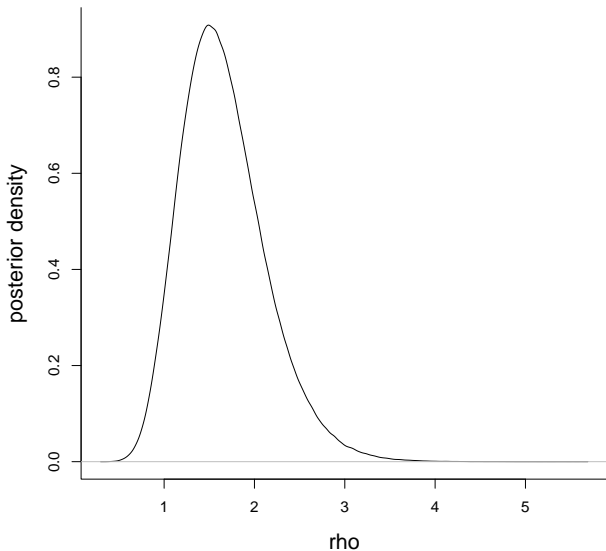
with $a = a_1 + \dots + a_k$.

Easy to use, via simulation:

$$(p_1, \dots, p_k) = \left(\frac{G_1}{G_1 + \dots + G_k}, \dots, \frac{G_k}{G_1 + \dots + G_k} \right),$$

with $G_j \sim \text{Gamma}(a_j, 1)$ (for the prior), or
 $G_j \sim \text{Gamma}(a_j + y_j, 1)$ (for the posterior).

Example: I throw my die 60 times and get 8, 8, 7, 13, 9, 15. Is p_6 bigger than it should be? With prior $\text{Dir}(2, 2, 2, 2, 2, 2)$ this is the posterior for $\rho = p_6 / (p_1 \cdots p_5)^{1/5}$, and $\Pr(\rho > 1 \mid \text{data}) = 0.947$:



Then from k boxes to the **infinite full-space process**: We're helped by this **collapsibility lemma**: If

$$(p_1, \dots, p_{10}) \sim \text{Dir}(a_1, \dots, a_{10}),$$

then

$$(p_1 + p_2, p_3 + p_4 + p_5, p_6, p_7 + p_8 + p_9 + p_{10}) \\ \sim \text{Dir}(a_1 + a_2, a_3 + a_4 + a_5, a_6, a_7 + a_8 + a_9 + a_{10}),$$

etc. With P_0 a distribution on the sample space S , we say that $P \sim \text{Dir}(aP_0)$, a Dirichlet process with parameter aP_0 , if

$$(P(A_1), \dots, P(A_k)) \sim \text{Dir}(aP_0(A_1), \dots, aP_0(A_k))$$

for each partition A_1, \dots, A_k .

Existence is non-trivial (Ferguson, 1973, Doksum 1974). We have

$$\mathbb{E} P(A) = P_0(A) \quad \text{and} \quad \text{Var} P(A) = \frac{P_0(A)\{1 - P_0(A)\}}{a + 1}.$$

The Dirichlet process has various uses as a probabilistic model for a random distribution. It is also well-suited for **inference after observations from unknown distribution**. **Master lemma** says that if (i) $P \sim \text{Dir}(aP_0)$ and (ii) $y_1, \dots, y_n \mid P$ are i.i.d. from P , then

$$P \mid \text{data} \sim \text{Dir}(aP_0 + \delta(y_1) + \dots + \delta(y_n)),$$

i.e. with posterior measure $aP_0 + nP_n$, with $P_n = n^{-1} \sum_{i=1}^n \delta(y_i)$, the empirical measure with point mass $1/n$ in each data point.

In particular:

$$\begin{aligned}\hat{P}(A) &= \text{E} \{P(A) \mid \text{data}\} = \frac{a}{a+n} P_0(A) + \frac{n}{a+n} P_n(A), \\ \text{Var} \{P(A) \mid \text{data}\} &= \frac{1}{a+n+1} \hat{P}(A) \{1 - \hat{P}(A)\}.\end{aligned}$$

May also form **confidence bands** for $P(A)$, and may e.g. **simulate 1000 realisations** from $P \mid \text{data}$, from which we can **read off posterior** for any $\theta(P)$.

How to simulate a $\text{Dir}(aP_0)$, over the sample space S ?

(i) May discretise, **sample space cut into tiny pieces**, and use a finite-dimensional Dirichlet ($aP_0(A_1), \dots, aP_0(A_{1000})$).

(ii) May write $P(A) = G(A)/G(S)$, with G a **Gamma process** with independent pieces over disjoint sets, $G(A) \sim \text{Gamma}(aP_0(A))$.

(iii) Via the Tiwari–Sethuraman **representation theorem**:

$$P = \sum_{j=1}^{\infty} p_j \delta(\theta_j),$$

where the random locations $\theta_1, \theta_2, \dots$ are i.i.d. from P_0 , and where the random **stick-breaking probabilities** are

$$p_1 = B_1, p_2 = (1 - B_1)B_2, p_3 = (1 - B_1)(1 - B_2)B_3, \dots,$$

with B_1, B_2, B_3, \dots i.i.d. $\text{Beta}(1, a)$.

Note that the random distribution P is **discrete**.

C: The Beta process

For survival analysis, consider life-times T , with distribution F , and cumulative hazard rate function A :

$$dA(s) = \frac{dF(s)}{F[s, \infty)} = \Pr\{\text{die in } [s, s + ds] \mid \text{still alive at } s\}.$$

The **survival curve** is

$$S(t) = \Pr\{T \geq t\} = \prod_{[0, t]} \{1 - dA(s)\}.$$

Hjort (1985, 1990) constructs a **Beta process**, with **independent increments**: With $A_0(\cdot)$ the prior mean function, and $c(\cdot)$ a prior strength function,

$$dA(s) \approx_d \text{Beta}(c(s) dA_0(s), c(s)\{1 - dA_0(s)\}).$$

The **existence is non-trivial**, since a sum of Betas is not a Beta; a fine-limit-argument is needed (cf. **Lévy processes**).

As for the Dirichlet process, also the Beta process has various probabilistical uses in studies of [random transitions phenomena](#), and as the de Finetti measure of the [Indian Buffet Processes](#).

They are particularly well-suited for [survival analysis](#). Survival data (t_i, δ_i) , with $t_i = \min(t_i^0, z_i)$ and $\delta_i = I\{t_i^0 \leq z_i\}$ the indicator for non-censoring: The classical nonparametric estimators for cumulative hazard and survival are

$$\tilde{A}(t) = \int_0^t \frac{dN(s)}{Y(s)} \quad \text{and} \quad \tilde{S}(t) = \prod_{[0,t]} \left\{ 1 - \frac{dN(s)}{Y(s)} \right\},$$

the [Nelson–Aalen](#) and [Kaplan–Meier](#) estimators. Here $Y(s)$ is the number at risk at time s and $dN(s)$ the number of those dying in $[s, s + ds]$.

With the Beta process, we reach [Bayesian nonparametric extensions of these](#).

If $A \sim \text{Beta}(c, A_0)$, then

$$A \mid \text{data} \sim \text{Beta}(c + Y, \hat{A}),$$

with

$$\hat{A}(t) = E\{A(t) \mid \text{data}\} = \int_0^t \frac{c(s) dA_0(s) + dN(s)}{c(s) + Y(s)}.$$

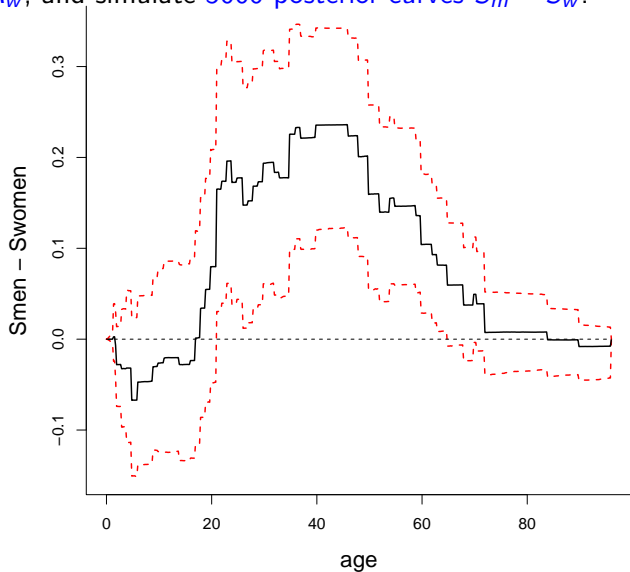
The Bayes estimator for survival is

$$\hat{S}(t) = E\{S(t) \mid \text{data}\} = \prod_{[0,t]} \left\{ 1 - \frac{c(s) dA_0(s) + dN(s)}{c(s) + Y(s)} \right\}.$$

With $c(s) \rightarrow 0$ we do not trust the prior, and we get the Nelson–Aalen and Kaplan–Meier estimators.

May also **simulate 1000 posterior realisations** from the distributions $A \mid \text{data}$ and $S \mid \text{data}$, and read off relevant features and probabilities.

Example: Analyse life-lengths from ancient Egypt, for 82 men and 59 women, via posterior distribution for $\Pr(T_m \geq t) - \Pr(T_w \geq t)$. I start with Beta process priors for A_m and A_w , and simulate 5000 posterior curves $S_m - S_w$.



Sir David Cox (b. 1924) is an Eternal Guru of Statistics (the first ever winner of the [International Prize in Statistics](#), 2017). His [most important invention](#) (from 1972) is the hazard rate regression model

$$\alpha_i(s) = \alpha(s) \exp(x_i^t \beta)$$

along with deep methodology for handling such and similar models, starting with [survival data](#) (t_i, δ_i, x_i) .

The [canonical semiparametric Bayesian extension](#) of this method (Hjort, 1990) starts with

$$1 - dA_i(s) = \{1 - dA(s)\}^{\exp(x_i^t \beta)},$$

a prior for β , and $A \sim \text{Beta}(c, A_0)$.

There is a long list of further generalisations and uses of the Beta process in [models with transitions over time](#) (medicine, demography, biology, [event history analysis](#), etc.).

D: Bernshteĭn–von Mises theorems (do not always hold)

For ordinary parametric inference, there's a **comforting general theorem** saying frequentist and Bayesian inferences 'agree in the end', with enough data.

Suppose n data points or vectors have been observed, from a model $f(y, \theta)$, with $\hat{\theta}$ the maximum likelihood estimator. **First**,

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N_p(0, J^{-1}).$$

Second, having started with **any prior**, the posterior $\pi(\theta | \text{data})$ is such that

$$\sqrt{n}(\theta - \hat{\theta}) | \text{data} \rightarrow_d N(0, J^{-1}).$$

So frequentist and (every) Bayesian inference tend to agree, with $\hat{\mu} \pm 1.96 \hat{\kappa} / \sqrt{n}$ as the 95% confidence or credibility interval, etc.

The **nonparametric world** is bigger and scarier (?), however. Various reasonable-looking nonparametric Bayesian schemes don't work – **lack of consistency**, **wrong coverage**, etc.

E: Stationary time series

Bayesian nonparametrics for **covariance functions** with application to **time series** – the following reports briefly on joint work with **Gudmund Hermansen**.

- ▶ Covariance and correlation functions: via spectral measure F
- ▶ Prior on $F \Rightarrow$ **prior on covariances and correlations**
- ▶ F a Dirichlet: $C(h) = \int_0^\pi \cos(h\omega) dF(\omega)$ is a valid correlation sequence
- ▶ Stationary time series: **full nonparametric Bayes inference**
- ▶ Other spatial and spatial-temporal models: prior, posterior, Bayesian inference ok; but fewer hard results

Plan:

- ▶ 1. Priors for stationary Gaussian time series – Spectral representation: F first, then C
- ▶ 2. Frequentist analysis – Periodogramme, cumulative, Brownian motion
- ▶ 3. Exact and approximate Bayesian updating – Whittle approximation, MCMC
- ▶ 4. Limit theorems and Bernshteĭn–von Mises
- ▶ 5. [Illustration: sun spots, etc.; not here (!)]

E1. Priors for stationary Gaussian time series

Let Y_1, Y_2, \dots be a zero-mean stationary Gaussian time series with unknown **covariance function**

$$C(h) = C(|i - j|) = \text{cov}(Y_i, Y_j) \quad \text{for } |j - i| = h.$$

Wish to use **Bayesian nonparametrics** for modelling $C(\cdot)$ – and hence any function of $C(\cdot)$, i.e. any function of the $n \times n$ covariance matrix.

If we manage, this leads to full Bayesian inference for a time series with ‘**uncertain covariance function**’. Can then also answer **predictive questions**, like the Nordmarkaesque

$$\alpha = \Pr\{Y_{n+1} \geq y_0, Y_{n+2} \geq y_0, Y_{n+3} \geq y_0 \mid \text{data}\}.$$

Would also wish to **centre** the random $C(\cdot)$ as some given $C_0(\cdot)$, say $C_0(h) = \sigma^2 \rho^h$ for AR(1).

Not easy to do it 'directly' – placing a random band around ρ^h quickly produces outcomes that are non-valid, i.e. the associated covariance matrices may be negative definite. We need

- ▶ the random $C(\cdot)$ is positive definite;
- ▶ clear interpretation of prior;
- ▶ big (or full) prior support;
- ▶ simulations (or approximations) to the posterior;
- ▶ posterior consistency;
- ▶ perhaps more, e.g. Bernshteĭn–von Mises.

General approach (but not the only one): modelling $C(\cdot)$ via spectral measure $F(\cdot)$ on $[0, \pi]$. Wold's theorem:

$$C(h) = 2 \int_0^\pi \cos(hu) dF(u) \quad \text{for } h = 0, 1, 2, \dots,$$

with F nondecreasing and finite; in particular

$C(0) = 2F(\pi) = \sigma^2 = \text{Var } Y_i$. Hence also for correlation function:

$$\text{corr}(h) = \int_0^\pi \cos(hu) \frac{dF(u)}{F(\pi)} \quad \text{with } \frac{F}{F(\pi)} \text{ random cdf.}$$

Main idea & programme:

- ▶ model for $F(\cdot)$
- ▶ \Rightarrow model for $C(\cdot)$
- ▶ \Rightarrow model for full covariance matrix and all related quantities
- ▶ \Rightarrow full posterior distribution of 'everything', when coupled with **data likelihood**, which is

$$L_n \propto \exp\left(-\frac{1}{2} \log |\Sigma_n| - \frac{1}{2} y^\top \Sigma_n^{-1} y\right).$$

To understand F models from C properties (and vice versa, they are a 'Fourier couple'):

$$C(h) = 2 \int_0^\pi \cos(hu) dF(u),$$
$$f(u) = \frac{1}{2\pi} \sum_{h=-\infty}^{\infty} \exp(-ihu) C(h) = \frac{\sigma^2}{2\pi} + \frac{1}{\pi} \sum_{h=1}^{\infty} \cos(hu) C(h).$$

Prior for F should match prior knowledge for $C(\cdot)$. May **centre** F at F_0 that matches some C_0 .

Example: AR(1). Here $C_0(h) = \sigma^2 \rho^h$ and spectral density becomes

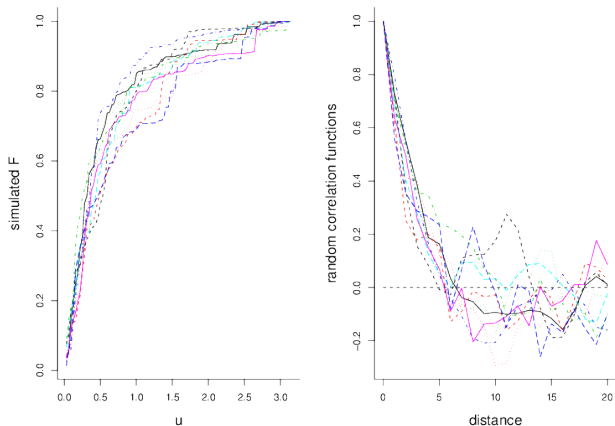
$$f_0(u) = \frac{\sigma^2}{2\pi} \frac{1 - \rho^2}{1 - 2\rho \cos u + \rho^2} \quad \text{for } u \in [0, \pi].$$

May e.g. choose a prior for F with prior mean matching $F_0(u) = \int_0^u f_0(v) dv$, and with uncertainty band matching prior uncertainty about $C(h)$ around $C_0(h)$.

Many possibilities: any random finite measure F gives rise to a random C etc.

Take J Bayesian nonparametrics papers for random finite measures \Rightarrow new papers [paper₁](#), ..., [paper_J](#) for Bayesian nonparametric analysis of stationary time series.

Good tool: F a Gamma process $(aF_0, 1) \Rightarrow F/F(\pi)$ a Dirichlet (cF_0) on $[0, \pi]$, etc.



Left: simulated Dirichlet processes, $F/F(\pi) \sim \text{Dir}(cF_0)$ on $[0, \pi]$;
 right: the accompanying *random correlation functions*
 $C(h) = \int_0^\pi \cos(hu) dF(u)/F(\pi)$.

E2. Frequentist analysis

Time series analysis is nearly always **parametric** (typically also with model selection issues etc.), though nonparametric analysis is also possible – for spectral density f , its cdf F , and hence the covariance function. Consider the **periodogramme** (Schuster, 1898)

$$I_n(u) = \frac{1}{n} \frac{1}{2\pi} \left| \sum_{k=1}^n \exp(-iku) y_k \right|^2 \quad \text{for } u \in [0, \pi].$$

One has $I_n(u_{n,j}) \approx f_{\text{true}}(u) \text{Expo}(1)$ when $u_{n,j} \rightarrow u$, and asymptotic independence between these limits when $u_{n,j} \doteq \pi j/n$.

There's a big literature on **smoothed periodogrammes** etc., for estimating f_{true} , but here we are more interested in F than f .

We may use either of

$$F_n(u) = \int_0^u I_n(v) dv \quad \text{and} \quad \widehat{F}_n(u) = \frac{\pi}{n} \sum_{\pi j/n \leq u} I_n(\pi j/n),$$

and have **process convergence**:

$$Z_n(u) = \sqrt{n} \{ \widehat{F}_n(u) - F_{\text{true}}(u) \} \rightarrow_d W \left(2\pi \int_0^u f_{\text{true}}(v)^2 dv \right).$$

The associated estimators of covariances $C(h)$ are

$\widehat{C}_n(h) = 2 \int_0^\pi \cos(uh) d\widehat{F}_n(u)$. From $Z_n = \sqrt{n}(\widehat{F}_n - F_{\text{true}}) \rightarrow_d Z$, a **time-transformed Brownian motion**, follows

$$\sqrt{n} \{ \widehat{C}_n(h) - C_{\text{true}}(h) \} \rightarrow_d A_h = 2 \int_0^\pi \cos(uh) dZ(u),$$

and variances and covariances may be written down and estimated consistently. We have also good **large-sample nonparametric control** over all other smooth functions of the Σ_n matrix, and can put down normal approximations and confidence intervals etc.

E3. Exact and approximate Bayesian updating

Attractive class of priors: let F have **independent increments**, and split the spectral domain into windows,

$$[0, \pi] = W_1 \cup \dots \cup W_m,$$

perhaps of equal width, $W_j = (\pi(j-1)/m, \pi j/m]$. The log-likelihood also almost splits into m components, across windows:

$$\begin{aligned}\ell_n &= -\frac{1}{2} \log |\Sigma_n| - \frac{1}{2} y^t \Sigma_n^{-1} y + \text{const}, \\ \tilde{\ell}_n &= -\frac{1}{2} n \frac{1}{\pi} \int_0^\pi \left\{ \log f(u) + \frac{I_n(u)}{f(u)} \right\} du + \text{const} = \sum_{j=1}^m \tilde{\ell}_{n,j}.\end{aligned}$$

Here $\tilde{\ell}_n$ is the **Whittle approximation** to the exact ℓ_n , and $I_n(u)$ the periodogramme. We need F to have **a density** f .

Hence **Bayesian updating** can be undertaken 'window by window'.

Special prior: locally constant spectral density,

$$f(u) = f_j \quad \text{for } u \in \text{window } W_j, \quad j = 1, \dots, m,$$

with priors $\pi_1(f_1), \dots, \pi_m(f_m)$ for these constants. The random F is continuous and **piecewise linear**.

Exact posterior distribution

$$\pi(f_1, \dots, f_m \mid \text{data}) \propto \pi_1(f_1) \cdots \pi_m(f_m) \exp\{\ell_n(f_1, \dots, f_m)\},$$

can be worked with, both practically (MCMC) and theoretically. For **growing** n (and **windows not too small**) it is close enough to its easier **Whittle approximation**:

$$\propto \pi_1(f_1) \exp\{\tilde{\ell}_{n,1}(f_1)\} \cdots \pi_m(f_m) \exp\{\tilde{\ell}_{n,m}(f_m)\},$$

where

$$\tilde{\ell}_{n,j} = -\frac{1}{2}n \frac{1}{\pi} \int_{W_j} \left\{ \log f_j + \frac{I_n(u)}{f_j} \right\} du = -\frac{1}{2}n \frac{1}{\pi} \left(w_j \log f_j + \frac{v_{n,j}}{f_j} \right)$$

for window W_j , with w_j length of W_j and $v_{n,j} = \int_{W_j} I_n(u) du$.

Full Bayesian analysis may now be carried out, from a given number of windows and given priors for the spectral heights f_1, \dots, f_m .

We may use exact posterior via MCMC, or approximate posterior via Whittle and independence,

$$\pi(f_j | \text{data}) \propto \pi_j(f_j) \exp\left\{-\frac{1}{2}n(1/\pi)(w_j \log f_j + v_{n,j}/f_j)\right\}$$

for $j = 1, \dots, m$, with $v_{n,j} = \int_{W_j} I_n(u) du$ and w_j the width of W_j .

Can use inverse gamma priors for the local constants (convenient updating), but there are reasons for preferring gamma priors, say $f_j \sim \text{Gam}(af_{0,j}, a)$.

We may compute posterior mean and variance directly (involving Bessel functions etc.), and also draw samples from $\pi(f_j | \text{data})$.

E4. Large-sample results

We have proven nice large-sample theorems that in a **Bernshteĭn–von Mises** fashion mirror the frequentist results. The essential conditions are $m \rightarrow \infty$ and $m/\sqrt{n} \rightarrow 0$, ‘more and more windows, but not too many’. Then:

- ▶ the Whittle approximation becomes good enough (same limit with exact and with Whittle);
- ▶ the parametric BvM theorem has time to kick in, for each window:

$$f_j | \text{data} \approx N\left(w_j^{-1} \int_{W_j} d\hat{F}_n(u), 2\pi \int_{W_j} f_{\text{true}}(v)^2 dv/n\right),$$

with x_j midpoint of W_j ;

- ▶ nonparametric **BvM process convergence**:

$$\sqrt{n}(F - \hat{F}_n) | \text{data} \rightarrow_d W\left(2\pi \int_0^u f_{\text{true}}(v)^2 dv\right);$$

with ‘invariance theorem’ consequences for $C(h)$ etc.

F: Nonparametric quantile inference

Suppose x_1, \dots, x_n are i.i.d. from a distribution F , with **quantile function**

$$Q(y) = F^{-1}(y) = \inf\{t: F(t) \geq y\}.$$

So $Q(\frac{1}{2})$ is the median; $Q(\frac{1}{4})$ and $Q(\frac{3}{4})$ the two quartiles, etc.

A **quantile pyramid** is constructed in this fashion:

- 1 give a prior for $Q(\frac{1}{2})$;
- 2 give priors for $Q(\frac{1}{4})$ and $Q(\frac{3}{4})$, given $Q(\frac{1}{2})$;
- 3 give priors for $Q(\frac{1}{8})$, $Q(\frac{3}{8})$, $Q(\frac{5}{8})$, $Q(\frac{7}{8})$, given $Q(\frac{1}{4})$, $Q(\frac{2}{4})$, $Q(\frac{3}{4})$;
- &cetera, &cetera.

Under **some conditions**, this pans out well (Hjort and Walker, Annals, 2009) – the full $Q = \{Q(y): y \in (0, 1)\}$ exists; there is a characterisation of $Q \mid \text{data}$; this may be computed and simulated from.

The class of Quantile Pyramids is **very large**. It may be used for **purely probabilistic analyses** of different types of phenomena, and for **statistical quantile inference**. A broad model is

$$z_i = m(x_i) + \varepsilon_i, \quad \text{where } \varepsilon_i \text{ has quantile process } Q,$$

with a prior process for $m(x)$. May then reach inference for

$$(Q(0.05 | x), Q(0.50 | x), Q(0.95 | x),$$

presented as bands in x , etc.

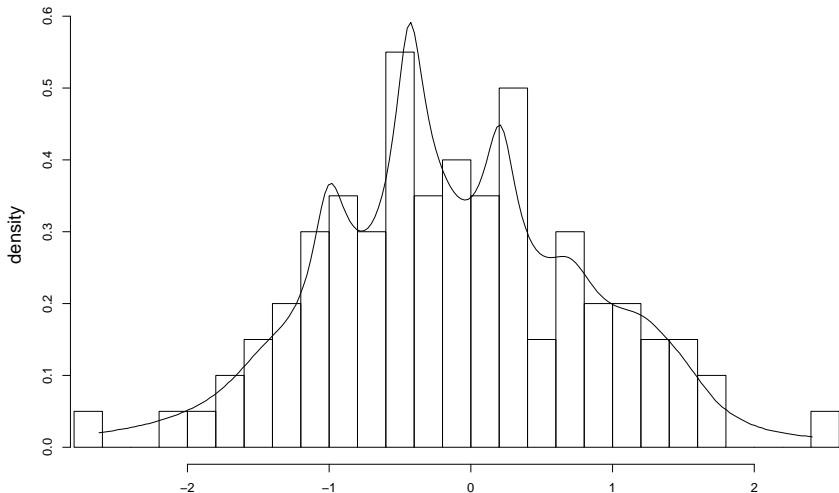
A **special case** of the Quantile Pyramid $Q = \{Q(y) : y \in (0, 1)\}$ corresponds to $F = \{F(x) : x \in R\}$ being a $\text{Dir}(aF_0)$. Cute quantile estimator:

$$\hat{Q}(y) = \sum_{i=1}^n \binom{n-1}{i-1} y^{i-1} (1-y)^{n-i} x_{(i)} \quad \text{for } 0 \leq y \leq 1.$$

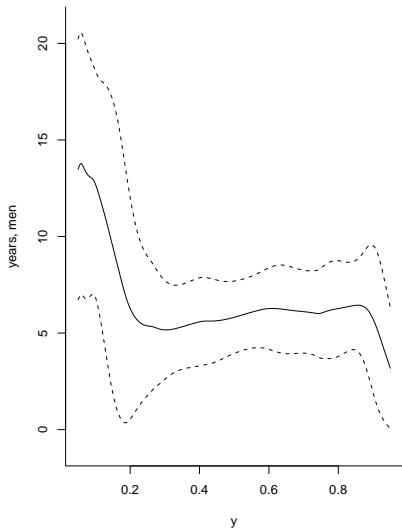
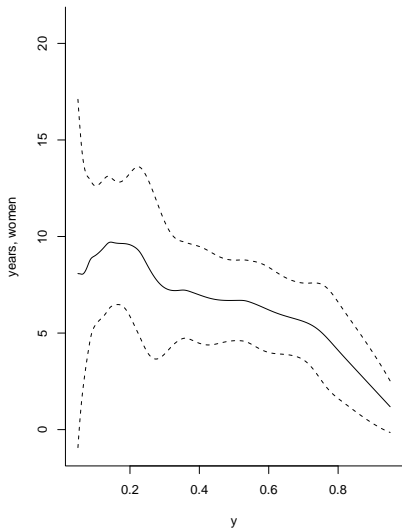
It has $\hat{Q}(0) = x_{(1)}$ and $\hat{Q}(1) = x_{(n)}$.

A **fully automatic density estimator** (Hjort and Petrone, 2007):
solve $\hat{Q}(y) = x$ to identify $y = \hat{F}(x)$, and then

$$\hat{f}(x) = \left[\sum_{i=1}^{n-1} (x_{(i+1)} - x_{(i)}) \text{beta}(\hat{F}(x); i, n - i) \right]^{-1}.$$



Quantile difference function $G^{-1}(y) - F^{-1}(y)$, with bands, for F age at hospitalisation and G age at death, for women and for men. (I've re-analysed data from Laake, Laake, Aaberge, 1985.)



G: Models and methods for clusters

A **simple prototype setup**: Data points y_1, \dots, y_n are to be clustered, say as belonging to $N(\xi_j, 1)$, with cluster centres ξ_j , and we do not know the number of clusters in advance (so this is not **k-means** or similar).

- ▶ 1 Let $P \sim \text{Dir}(aP_0)$.
- ▶ 2 Let μ_1, \dots, μ_n i.i.d. P – **but only D_n of these n will be distinct.**
- ▶ 3 Let $y_i \sim N(\mu_i, 1)$ for $i = 1, \dots, n$.

May then set up a posterior scheme simulating from (μ_1, \dots, μ_n) . From this one reads off both $\pi(D_n | \text{data})$ and the **positions** of cluster centres.

\exists hundreds of variations – many in heavy use. Note that a influences size of D_n : $D_n \approx a \log n$. (Can also have a prior on a .)

Bigger Things

Hierarchies of hierarchies, Dirichlet process of Dirichlet processes;
Beta process of Beta processes; ...

$$G_0 \mid \gamma, H \sim \text{Dir}(\gamma, H),$$
$$G_j \mid \alpha, G_0 \sim \text{Dir}(\alpha, G_0),$$

with $G_0 = \sum_{j=1}^{\infty} \delta(\theta_j) p_j$ etc. There's a **sharing of atoms** involved.

- ▶ Information retrieval: 'term frequency / inverse document frequency' for ranking of documents.
- ▶ Multipopulation haplotype phasing.
- ▶ Topic modelling.
- ▶ HMMs with infinite state spaces.
- ▶ Hierarchical clustering.
- ▶ Automatic translation.

H: Concluding remarks

Bayesian nonparametrics has **grown drastically**, from c. 1973 to now – in horizon size, ambition level, flexibility, convenience, popularity (!), computational feasibility, applicability, **maturity**. It's close friends with branches of **probability theory and applications** and with **machine learners** and with all uses of Big Hierarchical Constructions.

Its uses include **more flexibility around stricter models** (nonparametric envelopes around parametric models).

It links with machine learning for **nonparametric regression and classification**; for **hierarchical structures** ('Dirichlet process of Dirichlet processes'); and for clustering and allocation processes (Chinese Restaurant Process, **Indian Buffet Process**).

FocuStat Workshop May 2015: [CDs and Related Fields](#)

FocuStat Workshop May 2016: [FICology](#)

FocuStat Workshop May 23–25 2017: [Building Bridges](#), 'from parametrics to nonparametrics', including Bayesian nonparametrics.

