

Exercises and Lecture Notes

STK 9190, Autumn 2019

Version 0.81, 3-x-2019

Nils Lid Hjort

Department of Mathematics, University of Oslo

Abstract

Exercises and Lecture Notes collected here are indeed for the Bayesian Nonparametrics course STK 9190, given for the first time in the spring semester 2018 and then for the second time in the autumn semester of 2019. It is still useful to go through some prototype lower-dimensional Bayesian work, so a few exercises of that type are also included. This relates to clarifying concepts and principles, and also to Bayesian Nonparametrics constructions that use lower-dimensional pieces – as the famous interlocking versatile Lego bricks pieces.

1. Prior to posterior updating with Poisson data

This exercise illustrates the basic prior to posterior updating mechanism for Poisson data.

- (a) First make sure that you are reasonably acquainted with the Gamma distribution. We say that $Z \sim \text{Gamma}(a, b)$ if its density is

$$g(z) = \frac{b^a}{\Gamma(a)} z^{a-1} \exp(-bz) \quad \text{on } (0, \infty).$$

Here a and b are positive parameters. Show that

$$\text{E} Z = \frac{a}{b} \quad \text{and} \quad \text{Var} Z = \frac{a}{b^2} = \frac{\text{E} Z}{b}.$$

In particular, low and high values of b signify high and low variability, respectively.

- (b) Now suppose $y|\theta$ is a Poisson with parameter θ , and that θ has the prior distribution $\text{Gamma}(a, b)$. Show that $\theta|y \sim \text{Gamma}(a + y, b + 1)$.
- (c) Then suppose there are repeated Poisson observations y_1, \dots, y_n , being i.i.d. $\sim \text{Pois}(\theta)$ for given θ . Use the above result repeatedly, e.g. interpreting $p(\theta|y_1)$ as the new prior before observing y_2 , etc., to show that

$$\theta|y_1, \dots, y_n \sim \text{Gamma}(a + y_1 + \dots + y_n, b + n).$$

Also derive this result directly, i.e. without necessarily thinking about the data having emerged sequentially.

- (d) Suppose the prior used is a rather flat $\text{Gamma}(0.1, 0.1)$ and that the Poisson data are 6, 8, 7, 6, 7, 4, 11, 8, 6, 3. Reconstruct a version of Figure 0.1 in your computer, plotting the ten curves $p(\theta | \text{data}_j)$, where data_j is y_1, \dots, y_j , along with the prior density. Also compute the ten Bayes estimates $\hat{\theta}_j = E(\theta | \text{data}_j)$ and the posterior standard deviations, for $j = 0, \dots, 10$.
- (e) The mathematics turned out to be rather uncomplicated in this situation, since the Gamma continuous density matches the Poisson discrete density so nicely. Suppose instead that the initial prior for θ is a uniform over $[0.5, 50]$. Try to compute posterior distributions, Bayes estimates and posterior standard deviations also in this case, and compare with what you found above.

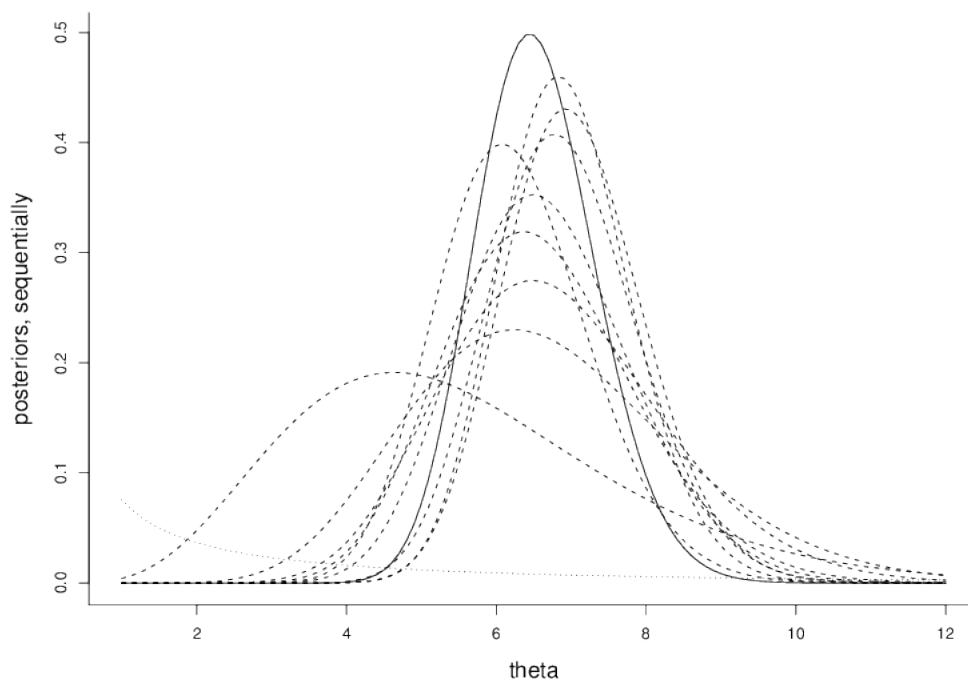


Figure 0.1: Eleven curves are displayed, corresponding to the $\text{Gamma}(0.1, 0.1)$ initial prior density for the Poisson parameter θ along with the ten updates following each of the observations 6, 8, 7, 6, 7, 4, 11, 8, 6, 3.

2. The Master Recipe for finding the Bayes solution

I decide to copy in this particular exercise from the lower-dimensional lower-ambition Bayes course, without changing the terms or the notation. The meta-exercise, however, is to understand that all of this still applies in the higher-level world of Bayesian Nonparametrics, partly at the price of the required higher-level mathematical abstraction level. Basically, where one for Bayesian Parametrics writes model likelihoods in terms of the famous generic θ , below, one needs for Bayesian Nonparametrics to think and write and work in terms of a very-high-dimensional or even infinite-dimensional parameter vector. This could be an unknown cumulative distribution function F , an unknown median regression function $m(x)$, an intensity function $\lambda(t)$, etc., rather than the prototypical θ . Often enough there are no clear-and-simple likelihood functions coming out of such

constructions, however, as we shall see during the course. This does not stop us from trying to crunch our way from priors to posteriors.

Crucially and amazingly, the basic concepts of decision functions, prior and posterior, loss functions and risk functions, and the optimal Bayesian strategy, carry over. As long as the statistician has data y , a model in terms of some distribution P (i.e. rather than the ubiquitous θ), a clear (nonparametric) prior for this P , and a loss function $L(P, a)$ encountered for decision a if the truth is P – then there will be (a) a posterior $\pi(P | \text{data})$; (b) a clear strategy for reaching the Bayes solution \hat{a}_B ; and (c) this strategy is unbeatable, the sole gold medal winner, in the Olympic competition against other strategies.

Consider a general framework with data y , in a suitable sample space \mathcal{Y} ; having likelihood $p(y | \theta)$ for given parameter θ (stemming from an appropriate parametric model), with θ being inside a parameter space Ω ; and with loss function $L(\theta, a)$ associate with decision or action a if the true parameter value is θ , with a belonging to a suitable action space \mathcal{A} . This could be the real line, if a parameter space is called for; or a two-valued set $\{\text{reject}, \text{accept}\}$ if a hypothesis test is being carried out; or the set of all intervals, if the statistician needs a confidence interval.

A statistical *decision function*, or procedure, is a function $\hat{a}: \mathcal{Y} \rightarrow \mathcal{A}$, getting from data y the decision $\hat{a}(y)$. Its *risk function* is the expected loss, as a function of the parameter:

$$R(\hat{a}, \theta) = E_{\theta} L(\theta, \hat{a}) = \int L(\theta, \hat{a}(y)) p(y | \theta) dy.$$

(In particular, in this expectation operation the random element is y , having its $p(y | \theta)$ distribution for given parameter, and the integration range is that of the sample space \mathcal{Y} .)

So far the framework does not include Bayesian components per se, and is indeed a useful one for frequentist statistics, where risk functions for different decision functions (be they estimators, or tests, or confidence intervals, depending on the action space and the loss function) may be compared.

We are now adding one more component to the framework, however, which is that of a *prior distribution* $p(\theta)$ for the parameter. The overall risk, or *Bayes risk*, associated with a decision function \hat{a} , is then the overall expected loss, i.e.

$$\text{BR}(\hat{a}, p) = E R(\hat{a}, \theta) = \int R(\hat{a}, \theta) p(\theta) d\theta.$$

(Here θ is the random quantity, having its prior distribution, making also the risk function $R(\hat{a}, \theta)$ random.) The *minimum Bayes risk* is the smallest possible Bayes risk, i.e.

$$\text{MBR}(p) = \min\{\text{BR}(\hat{a}, p) : \text{all decision functions } \hat{a}\}.$$

The *Bayes solution* for the problem is the strategy or decision function \hat{a}_B that succeeds in minimising the Bayes risk, with the given prior, i.e.

$$\text{MBR}(p) = \text{BR}(\hat{a}_B, p).$$

The *Master Theorem* about Bayes procedures is that there is actually a recipe for finding the optimal Bayes solution $\hat{a}_B(y)$, for the given data y (even without taking into account other values y' that could have been observed).

- (a) Show that the *posterior density* of θ , i.e. the distribution of the parameter given the data, takes the form

$$p(\theta | y) = k(y)^{-1} p(\theta) p(y | \theta),$$

where $k(y)$ is the required integration constant $\int p(\theta)p(y|\theta) d\theta$. This is the *Bayes theorem*.

(b) Show also that the *marginal distribution* of y becomes

$$p(y) = \int p(y|\theta)p(\theta) d\theta.$$

(I follow a certain semi-classical convention here, regarding using the ‘ p ’ multipurposedly, and with each ‘ p ’ to be understood by the reader from the context.)

(c) Show that the overall risk may be expressed as

$$\begin{aligned} \text{BR}(\hat{a}, p) &= \text{E} L(\theta, \hat{a}(Y)) \\ &= \text{E} \text{E} \{L(\theta, \hat{a}(Y)) | Y\} \\ &= \int \left\{ \int L(\theta, \hat{a}(y)) p(\theta | y) d\theta \right\} p(y) dy. \end{aligned}$$

The inner integral, or ‘inner expectation’, is $\text{E}\{L(\theta, \hat{a}(y)) | y\}$, the expected loss given data.

(d) Show then that the optimal Bayes strategy, i.e. minimising the Bayes risk, is achieved by using

$$\hat{a}_B(y) = \text{argmin } g = \text{the value } a_0 \text{ minimising the function } g,$$

where $g = g(a)$ is the expected posterior loss,

$$g(a) = \text{E}\{L(\theta, a) | y\}.$$

The g function is evaluated and minimised over all a , for the given data y . This is the Bayes recipe. – For examples and illustrations, with different loss functions, see the Nils 2008 Exercises.

2B. Two binomials: illustrating the Bayes rule

Consider a very simple setup, with two independent binomials, $y_0 \sim \text{Bin}(n_0, p_0)$ and $y_1 \sim \text{Bin}(n_1, p_1)$. The point with the exercise is to illustrate the use of the Master Recipe (see Exercise 2) to find the Bayes solution, i.e. the Optimal Decision!, with two different loss functions.

Feel free to create your own mini-narrative around this, but a simple story line, not too artificial, is as follows. The $p_0 = \text{Pr}(E)$ probability is ‘status quo’, the chance of success with some established procedure, and that p_1 is the possibly higher probability of the same good event E with a new procedure. Here one would be interested both in the gap or relative change $\delta = p_1 - p_0$ and in whether such a change is sufficiently big to warrant implementation of a new and perhaps costly system.

How should one teach Master level courses in statistics at the Department of Mathematics, University of Oslo? Nils teaches three hours a week, for a given course, resulting in $y_0 = 33$ achieving A or B for the last $n_0 = 100$ students. Imagine that Nils forced another $n_1 = 100$ through a more laborious six hours a week, with the number of successes a higher $y_1 = 44$.

(a) With independent uniform priors on p_0 and p_1 , show that what we now know, after counting $y_0 = 33$ and $y_1 = 44$, corresponds to $p_0 \sim \text{Beta}(34, 68)$ and $p_1 \sim \text{Beta}(45, 57)$. Simulate 10^5 realisations of $\delta = p_1 - p_0$ given data, and check `plot(density(deltasim))` in R (or, more simply, check the histogram).

- (b) For the loss function $L_0((p_0, p_1), \hat{\delta}) = |\hat{\delta} - \delta|$, with $\delta = p_1 - p_0$, compute the Bayes estimate. Give also the 0.05 and 0.95 points in the posterior distribution for δ . Would you ‘bet’, with or without quotation marks, that p_1 is higher than p_0 ?
- (c) The department can now *either* go for status quo (no change), *or* drastically decide that over the coming years, all courses should be taught using the Nils Method (six hours a week per course). Consider the loss function

$$L((p_0, p_1), \text{do nothing}) = \begin{cases} 0 & \text{if } p_1 \leq p_0, \\ 1 & \text{if } p_1 > p_0, \end{cases}$$

$$L((p_0, p_1), \text{change system}) = \begin{cases} 0 & \text{if } p_1 > p_0, \\ k & \text{if } p_1 \leq p_0, \end{cases}$$

where the cost is measured in millions of future student happiness quality units, and where $k = 25$. Comment briefly on such a loss function. Compute the posterior expected losses,

$$r_0 = E\{L((p_0, p_1), \text{do nothing}) \mid \text{data}\} \quad \text{and} \quad r_1 = E\{L((p_0, p_1), \text{change system}) \mid \text{data}\}.$$

What is your proposal, for the Department of Mathematics?

- (d) The loss function used above is arguably a bit too simple, since increasing p_0 from say 0.33 to 0.75 is rather more important and impressive (for the future lives of the students) than if p_0 is merely pushed from 0.33 to 0.36. Construct a suitably modified loss function, say $L^*((p_0, p_1), A)$, where the action A now can be either ‘do nothing’, or ‘let Nils teach 100 more students the new way, before we make a final decision’, or ‘we change the system right away’. With your L^* loss function, what is your proposal to the Department of Mathematics?
- (e) This exercise has so far been utterly low-parametric. Importantly, however, the same type of reasoning, and the same way of finding the optimal decision, apply also for the much more complicated models, setups, priors, posteriors, for Bayesian nonparametrics, i.e. for this course. To make it into an exercise, or half of an exercise, suppose an Old System produces outcomes y following a cumulative distribution function (cdf) F_0 , but a possible New System leads to outcomes that follow a potentially more benevolent cdf F_1 . We shall soon learn ways of putting nonparametric priors for F_0 and F_1 , and to compute their posteriors. Invent a loss function and think through how optimal decisions can be reached.

3. The Dirichlet-multinomial model

The Beta-binomial model, with a Beta distribution for the binomial probability parameter, is on the ‘Nice List’ where the Bayesian machinery works particularly well: Prior elicitation is easy, as is the updating mechanism. This exercise concerns the generalisation to the Dirichlet-multinomial model, which is certainly also on the Nice List and indeed in broad and frequent use for a number of statistical analyses.

- (a) Let (y_1, \dots, y_m) be the count vector associated with n independent experiments having m different outcomes A_1, \dots, A_m . In other words, y_j is the number of events of type A_j , for

$j = 1, \dots, m$. Show that if the vector of $\Pr(A_j) = p_j$ is constant across the n independent experiments, then the probability distribution governing the count data is

$$f(y_1, \dots, y_m) = \frac{n!}{y_1! \cdots y_m!} p_1^{y_1} \cdots p_m^{y_m}$$

for $y_1 \geq 0, \dots, y_m \geq 0, y_1 + \cdots + y_m = n$. This is the multinomial model. Explain how it generalises the binomial model.

(b) Show that

$$E Y_j = np_j, \quad \text{Var } Y_j = np_j(1 - p_j), \quad \text{cov}(Y_j, Y_k) = -np_j p_k \text{ for } j \neq k.$$

(c) Now define the Dirichlet distribution over m cells with parameters (a_1, \dots, a_m) as having probability density

$$\pi(p_1, \dots, p_{m-1}) = \frac{\Gamma(a_1 + \cdots + a_m)}{\Gamma(a_1) \cdots \Gamma(a_m)} p_1^{a_1-1} \cdots p_{m-1}^{a_{m-1}-1} (1 - p_1 - \cdots - p_{m-1})^{a_m-1},$$

over the simplex where each $p_j \geq 0$ and $p_1 + \cdots + p_{m-1} \leq 1$. Of course we may choose to write this as

$$\pi(p_1, \dots, p_{m-1}) \propto p_1^{a_1-1} \cdots p_{m-1}^{a_{m-1}-1} p_m^{a_m-1},$$

with $p_m = 1 - p_1 - \cdots - p_{m-1}$; the point is however that there are only $m - 1$ unknown parameters in the model as one knows the m th once one learns the values of the other $m - 1$. Show that the marginals are Beta distributed,

$$p_j \sim \text{Beta}(a_j, a - a_j) \quad \text{where } a = a_1 + \cdots + a_m.$$

(d) Infer from this that

$$E p_j = p_{0,j} \quad \text{and} \quad \text{Var } p_j = \frac{1}{a+1} p_{0,j}(1 - p_{0,j}),$$

in terms of $a_j = ap_{0,j}$. Show also that

$$\text{cov}(p_j, p_k) = -\frac{1}{a+1} p_{0,j} p_{0,k} \quad \text{for } j \neq k.$$

For the ‘flat Dirichlet’, with parameters $(1, \dots, 1)$ and prior density $(m - 1)!$ over the simplex, find the means, variances, covariances.

(e) Now for the basic Bayesian updating result. When (p_1, \dots, p_m) has a $\text{Dir}(a_1, \dots, a_m)$ prior, then, given the multinomial data, show that

$$(p_1, \dots, p_m) \mid \text{data} \sim \text{Dir}(a_1 + y_1, \dots, a_m + y_m).$$

Give formulae for the posterior means, variances, and covariances. In particular, explain why

$$\hat{p}_j = \frac{a_j + y_j}{a + n}$$

is a natural Bayes estimate of the unknown p_j . Also find an expression for the posterior standard deviation of the p_j .

- (f) In order to carry out easy and flexible Bayesian inference for p_1, \dots, p_m given observed counts y_1, \dots, y_m , one needs a recipe for simulating from the Dirichlet distribution. One such is as follows: Let X_1, \dots, X_m be independent with $X_j \sim \text{Gamma}(a_j, 1)$ for $j = 1, \dots, m$. Then the ratios

$$Z_1 = \frac{X_1}{X_1 + \dots + X_m}, \dots, Z_m = \frac{X_m}{X_1 + \dots + X_m}$$

are in fact $\text{Dir}(a_1, \dots, a_m)$. Try to show this from the transformation law for probability distributions: If X has density $f(x)$, and $Z = h(X)$ is a one-to-one transformation with inverse $X = h^{-1}(Z)$, then the density of Z is

$$g(z) = f(h^{-1}(z)) \left| \frac{\partial h^{-1}(z)}{\partial z} \right|$$

(featuring the determinant of the Jacobian of the transformation). Use in fact this theorem to find the joint distribution of (Z_1, \dots, Z_{m-1}, S) , where $S = X_1 + \dots + X_m$ (one discovers that the Dirichlet vector of Z_j is independent of their sum S).

- (g) The Dirichlet distribution has a nice ‘collapsibility’ property: If say (p_1, \dots, p_8) is $\text{Dir}(a_1, \dots, a_8)$, show that then the collapsed vector $(p_1 + p_2, p_3 + p_4 + p_5, p_6, p_7 + p_8)$ is $\text{Dir}(a_1 + a_2, a_3 + a_4 + a_5, a_6, a_7 + a_8)$.

4. Gott würfelt nicht

... but I do so, on demand. I throw a certain moderately strange-looking die 30 times and have counts $(2, 5, 3, 7, 5, 8)$ of outcomes 1, 2, 3, 4, 5, 6.

- (a) Use either of the priors (i) ‘flat’, $\text{Dir}(1, 1, 1, 1, 1, 1)$; (ii) ‘symmetric but more confident’, $\text{Dir}(3, 3, 3, 3, 3, 3)$; (iii) ‘unwilling to guess’, $\text{Dir}(0.1, 0.1, 0.1, 0.1, 0.1, 0.1)$

for the probabilities (p_1, \dots, p_6) to assess the posterior distribution of each of the following quantities:

$$\begin{aligned} \rho &= p_6/p_1, \\ \alpha &= (1/6) \sum_{j=1}^6 (p_j - 1/6)^2, \\ \beta &= (1/6) \sum_{j=1}^6 |p_j - 1/6|, \\ \gamma &= (p_4 p_5 p_6)^{1/3} / (p_1 p_2 p_3)^{1/3}. \end{aligned}$$

- (b) The above priors are slightly artificial in this context, since they do not allow the explicit possibility that the die in question is plain boring utterly simply a correct one, i.e. that $p = p_0 = (1/6, \dots, 1/6)$. The priors used hence do not give us the possibility to admit that ok, then, perhaps $\rho = 1, \alpha = 0, \beta = 0, \gamma = 1$, after all. This motivates using a mixture prior which allows a positive chance for $p = p_0$. Please therefore redo the Bayesian analysis above, with the same $(2, 5, 3, 7, 5, 8)$ data, for the prior $\frac{1}{2} \delta(p_0) + \frac{1}{2} \text{Dir}(1, 1, 1, 1, 1, 1)$. Here $\delta(p_0)$ is the ‘degenerate prior’ that puts unit point mass at position p_0 . Compute in particular the posterior probability that $p = p_0$, and display the posterior distributions of $\rho, \alpha, \beta, \gamma$.

- (c) [xx something re the same ideas being generalisable to the fuller Dirichlet process case. xx]

5. The Dirichlet Process: definition, existence, constructions

Let \mathcal{X} be some sample space, like the real line, with subsets A belonging to an appropriate sigma-algebra \mathcal{A} . Let P_0 be a fixed probability distribution on \mathcal{X} , and a a positive scalar. We say that P is a Dirichlet process on \mathcal{X} , with parameter aP_0 , and write $P \sim \text{Dir}(aP_0)$ to indicate this, if it is the case for each partition (A_1, \dots, A_m) , we have

$$(p_1, \dots, p_m) = (P(A_1), \dots, P(A_m)) \sim \text{Dir}(aP_0(A_1), \dots, aP_0(A_m)).$$

This is required for any number m of elements in the partition.

- (a) Show that the basic ‘logic coherence’ property is satisfied, that we may put some of the A_j together where the resulting distribution does not clash with the start definition. For example, with sets A_1, \dots, A_8 in such a partition, deduce the distribution for

$$(P(A_1) + P(A_2), P(A_3) + P(A_4) + P(A_5), P(A_6), P(A_7) + P(A_8)),$$

and verify that this is as it should be (i.e. the same distribution as dictated from the start definition). This is the ‘collapsibility property’ for the Dirichlet distribution, cf. Exercise 3(g). Without this property, the start definition would not make sense, and there would be no Dirichlet process.

- (b) The full existence of the $\text{Dir}(aP_0)$ is not a trivial matter, however. There are several routes to proving that yes, lo \mathcal{E} behold, it exists. Think a bit about the paths of proofs brief indicated below. If sufficiently curious (now or later), with enough time, go ad fontem and check the arguments.
- (i) Check the original argument used by Ferguson (1973, *Annals*), appealing to Kolmogorov’s consistency (or ‘inherent coherence’) theorem. Under a few natural and clearly necessary conditions, Kolmogorov proved that these are also sufficient; there will be no cognitive dissonance. Ferguson then verified the Kolmogorov dictated conditions. It is worth noting that in this fashion he ‘only’ got a random $P = \{P(A) : A \in \mathcal{A}\}$, with a certain well-defined probability distribution \mathcal{P} , in the enormous space $[0, 1]^{\mathcal{A}}$ of all function P on the enormous space \mathcal{A} , with values $P(A)$ in $[0, 1]$ for every A . He could then could go on to prove that $\mathcal{P}(\mathcal{M}) = 1$, where \mathcal{M} is the space of all probability measures on \mathcal{X} . This is still not the same as having created a \mathcal{P} working directly on \mathcal{M} . Several of the other Dirichlet process constructions are more direct than this, however.
- (ii) Check also Ferguson (1974, *Annals*), where a representation in the form of $P = Z/Z(\mathcal{X})$ is worked through, with $Z(\cdot)$ a Gamma process.
- (iii) Hjort (1976, last chapter) showed that the distribution \mathcal{P} of a $P \sim \text{Dir}(aP_0)$ can be reached as the well-defined limit in distribution of say \mathcal{P}_m , where \mathcal{P}_m is an easier finite-dimensional construction, basically a Dirichlet process $aP_{0,m}$ for a simpler discrete $P_{0,m}$ concentrated in only finitely many positions (for which the Dirichlet process existence is immediate). With the $P_{0,m}$ sequence constructed to tend in distribution to the perhaps continuous P_0 , Hjort showed that \mathcal{P}_m is tight; that its finite-dimensional distributions converge; that it must have a unique limit; and this limit is identical to Ferguson’s $\text{Dir}(aP_0)$. Care needs to be exercised regarding the convergence of probability measures on a space of probability measures (yes,

you heard that right). In other words, the complicatedness of the statement $\mathcal{P}_m \rightarrow_a \mathcal{P}$ needs to be examined carefully, as part of the construction.

‘Det er å håpe at denne alternative konstruksjonen av en Dirichlet-prosess ikke bare er av teoretisk verdi. Konstruksjonen gir informasjon utover det tre år gamle faktum at Dirichlet-prosessen eksisterer.’ (Hjort, 1976, last chapter.) Hjort’s 1976 construction takes place directly on the subspace \mathcal{M}_0 of all *discrete* probability measures on $(\mathcal{X}, \mathcal{A})$, so Ferguson’s non-trivial 1973 theorem that \mathcal{P} with probability 1 selects a discrete probability measure is here automatic.

- (iv) Tiwari and Sethuramam (1982, Purdue Symposium), and later Sethuraman (1994, Statistica Sinica), have given an intriguing explicit representation of a Dirichlet process, in the form of

$$P = \sum_{h=1}^{\infty} w_h \delta(\xi_h),$$

where the ξ_h are i.i.d. from P_0 , and the random probability weights w_h constructed in a certain way, discussed in Exercise [xx ... xx] below. Here, $\delta(\xi_h)$ means the degenerate point-mass measure with value 1 at position ξ_h .

- (v) Hjort (1990, Annals). [xx via the Beta process. xx]

- (vi) Hjort (2003, HSSS book). [xx via the symmetric representation and then the limit. xx]

6. Some properties for the Dirichlet process

Let $P \sim \text{Dir}(aP_0)$ on some space \mathcal{X} . Here are a few properties to go through, shedding light on the behaviour of the random P . Note that the Dirichlet process provides a model for random probability measures (hence also for random distribution functions, etc.), with independent or separate interest. The broader appeal lies however in its use as a prior for an unknown distribution, from which we then have observations, say X_1, \dots, X_n . See exercises and notes below.

- (a) With A a given set, show that

$$P(A) \sim \text{Beta}(aP_0(A), aP_0(A^c)),$$

with mean and variance

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

Thus P_0 is the mean of P , hence often called simply the prior mean. The a parameter indicates strength of belief in the prior guess; a large a means a tight distribution around P_0 , and vice versa for a smaller a .

- (b) Find the covariance and then correlation between $P(A)$ and $P(B)$, first for A and B disjoint, then with potential overlap.
- (c) With $g: \mathcal{X} \rightarrow \mathcal{R}$ a function, consider the random mean

$$\theta = \int g \, dP = \int g(x) \, dP(x).$$

Show that

$$\mathbb{E} \theta = \theta_0 = \int g \, dP_0,$$

so the mean of the random mean is the prior mean. Show also that

$$\text{Var} \theta = \frac{\sigma_0^2}{a+1},$$

with $\sigma_0^2 = \int (g - \theta_0)^2 \, dP_0$ the prior variance.

- (d) For two functions g_1, g_2 , consider the two random means $\theta_1 = \int g_1 \, dP$ and $\theta_2 = \int g_2 \, dP$. Find expressions for the covariance and correlation between these two random means.

7. The basic updating theorem for the Dirichlet process

Suppose $P \sim \text{Dir}(aP_0)$, and that $X | P$ follows the P distribution:

$$\mathcal{P}\{X \in A | P\} = P(A) \quad \text{for all } A.$$

In yet other words, X is a sample of size $n = 1$ from the given P , where P is selected randomly from the $\text{Dir}(aP_0)$ machine first.

- (a) Show that X has distribution P_0 . Start from

$$\mathbb{E} \{I(X \in A) | P\} = P(A)$$

and use double expectation.

- (b) The task is then to deduce the distribution of P given $X = x$. Attempt to show that if A_1, \dots, A_m is a partition, where x happens to lie in say the first of these, then

$$(P(A_1), \dots, P(A_m)) | (X = x) \sim \text{Dir}(aP_0(A_1) + 1, aP_0(A_2), \dots, aP_0(A_m)).$$

- (c) This is an indication that P given x is actually itself a Dirichlet process, with updated parameter $aP_0 + \delta(x)$. This also fits nicely with the finite-dimensional situation, see Exercise 3(f). You may attempt to give a formal proof of this basic updating statement for the Dirichlet process. See Ferguson (1973, Annals) or Ghosal and van der Vaart (2017, CUP book, Ch. 4).
- (d) Then consider a random sample X_1, \dots, X_n from the randomly selected P , with the defining property that

$$\mathcal{P}\{X_1 \in A_1, \dots, X_n \in A_n | P\} = P(A_1) \cdots P(A_n)$$

for all A_1, \dots, A_n . With P from the Dirichlet aP_0 , this defines a joint probability measure for (P, X_1, \dots, X_n) . Show, perhaps by induction, that

$$P | x_1, \dots, x_n \sim \text{Dir}\left(aP_0 + \sum_{i=1}^n \delta(x_i)\right).$$

This is really a wondrously and convenient convincing result, which matches the classical Dirichlet-multinomial situation examined in Exercise 3. Note that the parameter of the posterior Dirichlet process can be written

$$aP_0 + \sum_{i=1}^n \delta(x_i) = aP_0 + nP_n,$$

with $P_n = \sum_{i=1}^n (1/n)\delta(x_i)$ the empirical distribution for the n data points.

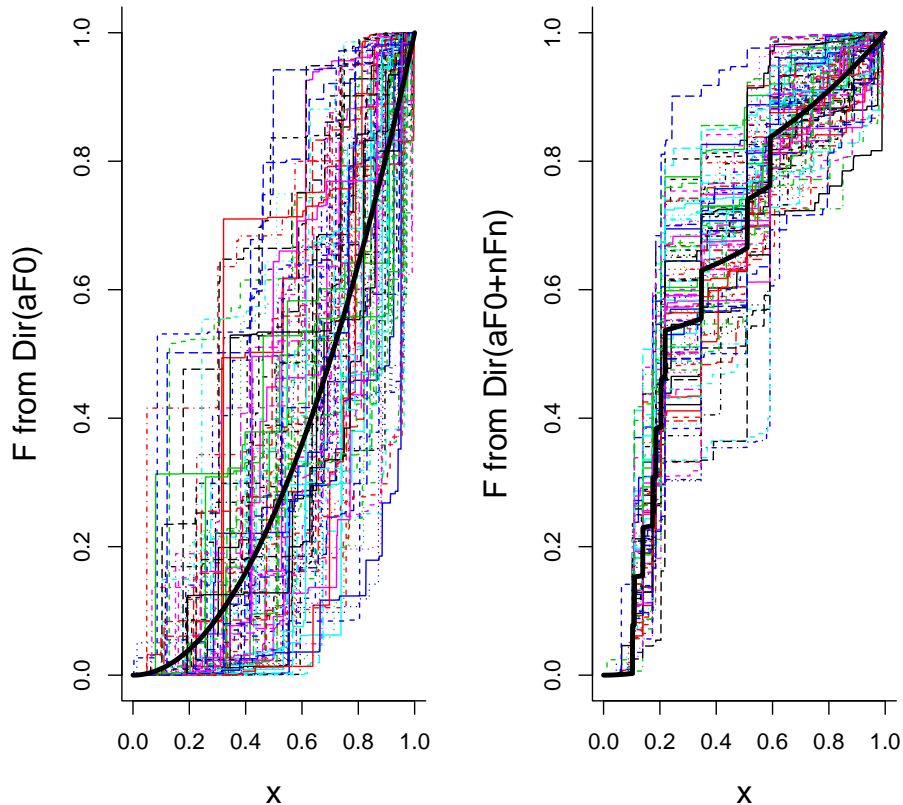


Figure 0.2: 100 simulations of F from the $\text{Dir}(aF_0)$ prior (left); then 100 simulations of F from the $\text{Dir}(aF_0 + nF_n)$ posterior (right), with the $n = 10$ data points of Exercise 8. The fat black curves are the prior mean and posterior mean, respectively.

8. Simulating from the prior and posterior, for a Dirichlet process

We need to be able to simulate realisations from the prior and the posterior, and here, specifically, from a given Dirichlet process. There are indeed several recipes for accomplishing this, but the simplest and most direct is to cut the space into a high number of smaller boxes, and then use the ensuing finite-dimensional Dirichlet as a fully adequate approximation. To carry out such finite-dimensional simulation we may use the recipe implicit in Exercise 3(g), which here means simulating a long list of small Gamma pieces and then normalising in the end.

Suppose you observe the following data points on the unit interval:

$$0.103, 0.110, 0.140, 0.175, 0.186, 0.205, 0.219, 0.348, 0.511, 0.592.$$

I have actually generated these from another distribution, namely the $\text{Beta}(1, 2)$, but the statistician seeing and about to analyse the data does not know this. For the prior for the unknown cumulative distribution function (cdf) F , take $F \sim \text{Dir}(aF_0)$, with F_0 the $\text{Beta}(2, 1)$.

- (a) Simulate say 100 realisations $F = \{F(x) : x \in [0, 1]\}$ from the prior, using the ‘lots of tiny boxes’ scheme of things. See the left panel of Figure 0.2, where I’ve used $a = 3.333$.

(b) Then simulate say 100 realisations F from the posterior, where

$$F \mid \text{data} \sim \text{Dir}(aF_0 + nF_n),$$

with $nF_n = \sum_{i=1}^n \delta(x_i)$. See the right panel of Figure 0.2.

(c) Show that the Bayes estimator, under quadratic loss, is

$$\widehat{F}_B(x) = \mathbb{E}\{F(x) \mid \text{data}\} = \frac{aF_0(x) + nF_n(x)}{a+n} = \frac{a}{a+n}F_0(x) + \frac{n}{a+n}F_n(x),$$

with F_n the empirical distribution function, i.e. the one having point-mass $1/n$ at each data point. Show furthermore that the posterior variance is

$$\widehat{\tau}^2(x) = \text{Var}\{F(x) \mid \text{data}\} = \frac{1}{n+a+1} \widehat{F}_B(x) \{1 - \widehat{F}_B(x)\}.$$

(d) Given realisations from F , these may be used to read off outcomes for parameters of interest, like $F(0.70) - F(0.60)$, the mean $\int_0^1 x dF(x)$, or the median

$$\mu = \min\{x: F(x) \geq \frac{1}{2}\}.$$

Carry out analysis for this random median, by computing the $\mu = \mu(F)$ for each realisation of F , for the prior and the posterior. This leads to Figure 0.3, where I used 10^4 simulations.

(e) Play with your code a bit, to see the influence of a small a or a large a , and of the choice of the prior mean cdf F_0 . You should also monitor what happens if you have say $n = 40$ data points from the underlying data generating mechanism, not only $n = 10$. You should get something similar to the right panel of Figure 0.2, but now with a slimmer and tighter spread around the Bayes estimator \widehat{F}_B .

(f) Then try $a = 0.0001$, a very tiny value, to see that happens with the posterior distribution of the median μ . You should learn that it has a distribution concentrated in the n data points. Try to find explicit formulae for these point masses,

$$\mathcal{P}(\mu = x_i \mid \text{data}), \quad \text{for } i = 1, \dots, 10.$$

9. War and peace, before and after Vietnam

Access the Tolstoyean `krigogfred-data` dataset on the course website and download it to your computer. It provides

$$(x_i, z_i) \quad \text{for } i = 1, \dots, 95,$$

the 95 inter-state wars from 1823 to 2003 with at least 1000 battle deaths; here x_i is time of onset and z_i the number of battle deaths, for war i . Look through Hjort's FocuStat Blog Post (which apparently impressed Steven Pinker enough to cause an admiring tweet about it, to his 368,001 followers), and also the Cunen, Hjort, Nygård (2019) paper, to get a sense of the themes, the questions, the predictions for our common future, and the controversies.

From these data, carry out the following two follow-up operations. First, limit attention to the 51 wars where $z_i \geq z_0$, with $z_0 = 7061$, a certain threshold value selected by A. Clauset, with the statistical intention that above this threshold, the density is proportional to $1/z^\alpha$, for an appropriate α . This is related to power laws and fat tails etc.; see again the Hjort blog post.

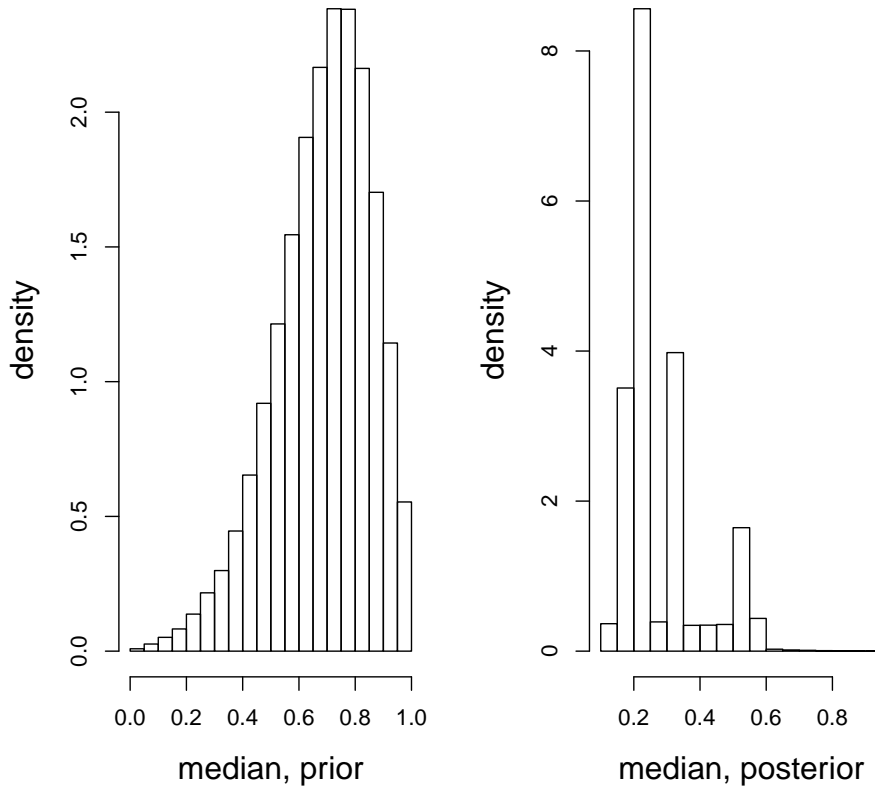


Figure 0.3: For the random median $\mu = \min\{x: F(x) \geq \frac{1}{2}\}$, I give histograms of its distribution, for the prior (left) and the posterior (right), based on 10^4 simulations, for each case.

Second, divide the remaining 51 value of (x_i, z_i) into a Left part, those 37 wars where $x_i \leq 1965.103$ (the onset-time for the Vietnam War), and a Right part, those 14 wars where $x_i > 1965.103$.

The statistical task is now to model and analyse the distributions of

$$y_i = \log(z_i/z_0) = \log z_i - \log 7061, \quad \text{for } i = 1, \dots, 51,$$

divided into

$$\begin{aligned} & y_1, \dots, y_{37}, \quad \text{with } x_i \text{ before and up to Vietnam,} \\ & y_{38}, \dots, y_{51}, \quad \text{with } x_i \text{ after Vietnam.} \end{aligned}$$

Specifically, we take the 37 before and including Vietnam to be i.i.d. from some F_L , and the 14 after Vietnam to be i.i.d. from some F_R .

- (a) Suppose Z has a density with the power law tail property that $f(z)$ is proportional to $1/z^\alpha$ for all z above some threshold z_0 . Show that this is equivalent to saying that $Y = \log(Z/z_0)$ has an exponential tail, specifically that $\Pr\{Y \geq y \mid Y \geq y_0\} = \exp(-\theta(y - y_0))$ for $y \geq y_0$, with $\theta = \alpha - 1$. Power law tail behaviour for the z_i , the battle deaths counts, can therefore be examined and in terms of exponential tails for the $\log z_i$.

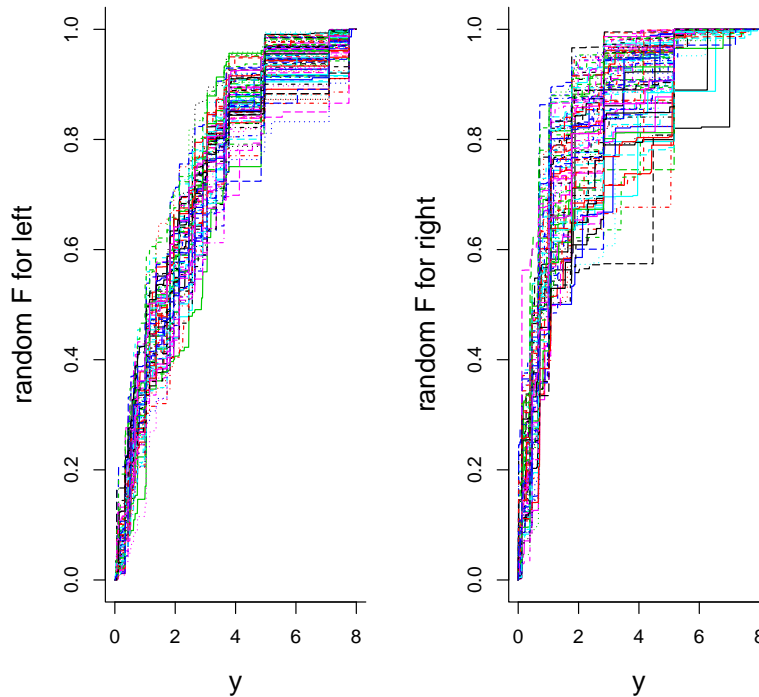


Figure 0.4: 100 simulated realisations of F_L , representing the past up to Vietnam (left), and 100 realisations of F_R , representing post-Vietnam period (right). The scale here is that of $y = \log(z/7061)$, for all wars with battle death counts at least 7061.

- (b) It makes sense to take the same prior $\text{Dir}(aF_0)$ for both F_L and F_R , since there is controversy in claiming that there is a difference between them at all; see Clauset's papers (2017, 2018). Take indeed $F_0(y) = 1 - \exp(-0.5y)$, and exponential, and $a = 3.333$ (later on you may tinker with that strength parameter). Work out the posterior distributions, and simulate say 100 realisations from each of them, as I have done to create Figure 0.4.
- (c) Carry out the consequent Bayesian nonparametric inference for the difference function $\delta(y) = F_L(y) - F_R(y)$. Plot the Bayes estimate $\hat{\delta}(y) = \text{E}\{\delta(y) | \text{data}\}$, along with a pointwise 90% credibility interval. The latter can be constructed accurately, via simulations, or via $\pm 1.645 \hat{\kappa}(y)$, where $\hat{\kappa}(y)$ is the posterior standard deviation. Attempt both methods.
- (d) [xx something more. inference for median of F_L minus median of F_R . xx]

10. The marginal distribution of a Dirichlet process sample

Suppose that $P \sim \text{Dir}(aP_0)$, and that data points are subsequently drawn independently from that P . The defining property for a sample of size n , is again that

$$\mathcal{P}\{X_1 \in A_1, \dots, X_n \in A_n | P\} = P(A_1) \cdots P(A_n),$$

for all sets A_1, \dots, A_n . Here we look at a few properties.

- (a) Let X be one of these points, say the first point. Show that its distribution is P_0 ; see also Exercise 7.

- (b) Consider next (X_1, X_2) , the two first data points. Show that their distribution can be expressed as

$$Q_2(A \times B) = \mathcal{P}\{X_1 \in A, X_2 \in B\} = \mathbb{E} P(A)P(B).$$

Then give formulae for this expression, (i) when A and B are disjoint; (ii) when they are identical; (iii) in the general case.

- (c) Show that

$$Q_2 = \frac{a}{a+1} P_0 \times P_0 + \frac{1}{a+1} P_{0,12},$$

where $P_{0,12}(A \times B) = P_0(A \cap B)$. We may think about this latter probability component $P_{0,12}$ as a mechanism that first picks $X_1 \sim P_0$ and then automatically takes the X_2 equal to the first.

- (d) Next study the joint distribution of three observations from a Dirichlet process. Note that X_1, X_2, X_3 are indeed i.i.d. given P , but the randomness in P makes the three dependent. Start from

$$Q_3(A \times B \times C) = \mathcal{P}_3\{X_1 \in A, X_2 \in B, X_3 \in C\} = \mathbb{E} P(A)P(B)P(C),$$

and give a formula for the case where A, B, C are disjoint.

- (e) [xx then finish this, give clear representation of Q_3 , find Hjort (1976). xx]

11. Ties and the steadily slowing stream of new guys from a Dirichlet

[xx about the process of seeing some old guys, occasionally a new guy, for Dirichlet samples. xx]

12. The number of discrete values in a Dirichlet sample

[xx to be written and polished. xx] we have $D_n = R_1 + \dots + R_n$ representation. we find $D_n / \log n \rightarrow_{pr} a$, and limiting normality from Nils 1976,

$$(\log n)^{1/2}(D_n / \log n - a) \rightarrow_d N(0, a).$$

Also, the simple $D_n / \log n$ is large-sample equivalent to the maximum likelihood estimator.

13. A simple models for clusters in data

[xx to be written out and polished. xx] We consider a simple hierarchical model which in a natural fashion leads to clusters, or groups, in the data, and where the number of such clusters is not specified in advance. The setup can be described as a three-step machinery, as follows:

- (i) A distribution P is taken from $\text{Dir}(aP_0)$;
- (ii) model parameters $\theta_1, \dots, \theta_n$ are sampled from P (which in particular means various ties);
- (iii) observations y_1, \dots, y_n are independent, given the $\theta_1, \dots, \theta_n$, and $y_i | \theta_i \sim f(y_i | \theta_i)$.

The Bayesian task is to understand the posterior distribution of $P, \theta_1, \dots, \theta_n$ given the observations y_1, \dots, y_n .

To make this clear and understandable in a simple prototype setup, consider a case where the parameters θ_i form a sample from P , where $P \sim \text{Dir}(aP_0)$, with $P_0 = N(0, \sigma_0^2)$. We also take $y_i \sim N(\theta_i, \sigma^2)$, with known σ . [xx more to come here. xx]

14. A clustering illustration

[xx a simple illustration here, with Dirichlet producing the model parameters, with lots of ties, etc. xx]

15. The Sethuraman stick-breaking representation

A somewhat surprising representation of the Dirichlet process, stemming from Sethuraman and Tiwari (1982, Purdue Symposium) and written out more fully in Sethuraman (1994, Sinica), is described here. With P_0 a probability measure, and a positive, we start with B_1, B_2, B_3, \dots being i.i.d. from $\text{Beta}(1, a)$. From these we form weights w_1, w_2, w_3, \dots , from

$$w_1 = B_1, \quad w_2 = (1 - B_1)B_2, \quad w_3 = (1 - B_1)(1 - B_2)B_3, \quad w_h = (1 - B_1) \cdots (1 - B_{h-1})B_h.$$

In addition, we draw an infinite i.i.d. sequence ξ_1, ξ_2, \dots from P_0 . The stick-breaking representation is

$$P = \sum_{h=1}^{\infty} w_h \delta(\xi_h),$$

with $\delta(\xi_h)$ the unit point-mass in position ξ_h .

(a) Show that

$$1 - w_1 - w_2 - w_3 = (1 - B_1)(1 - B_2)(1 - B_3),$$

with the immediate generalisation to $1 - w_1 - \dots - w_n$. Show from this that $\sum_{h=1}^{\infty} w_h = 1$, with probability 1.

(b) Try to comprehend & then sell the stick-breaking metaphor.

(c) For a fixed set A , consider the random probability $p = P(A)$, using the representation above. Show that p has mean $p_0 = P_0(A)$, and that

$$\text{Var } p = \text{E}(p - p_0)^2 = p_0(1 - p_0)/(a + 1).$$

(d) Attempt to prove that $p = P(A)$ is a $\text{Beta}(ap_0, a(1 - p_0))$.

(e) For a given bounded function g , consider the random mean

$$\theta = \int g \, dP = \sum_{h=1}^{\infty} w_h g(\xi_h).$$

Show that it has mean $\theta_0 = \int g \, dP_0$ and variance $\sigma_0^2/(a + 1)$, with $\sigma_0^2 = \int (g - \theta_0)^2 \, dP_0$.

(f) Consider disjoint sets A and B , and work with

$$p = P(A) = \sum_{\xi_h \in A} w_h \quad \text{and} \quad q = P(B) = \sum_{\xi_h \in B} w_h.$$

Calculate the covariance between p and q , from the stick-breaking representation above, and verify that you get the correct answer, i.e. what we should have if P indeed is a $\text{Dir}(aP_0)$.

(g) Then attempt to prove that the Sethuraman–Tiwari representation is correct, i.e. that P above with probability 1 becomes a $\text{Dir}(aP_0)$. – You may check the Sethuraman (1994, Sinica) paper. A more concise proof is given in Ghosal and van der Vaart (2017, Ch. 4), but this needs certain other properties which must be established separately, including a distributional equation property that uniquely characterises the Dirichlet process. See also Hjort and Ongaro (2005, SISP).

16. Dependent Dirichlet processes, using stick-breaking representations

[xx something here. check nils's discussion contribution to Gelfand and Petrone. xx] the basic idea, for two Dirichlet processes, which now become dependent: suppose (ξ_h, ξ'_h) are i.i.d. pairs, from some joint distribution, like the standardised binormal with correlation ρ . let P_0 and P'_0 be the marginals of this joint distribution for pairs. then construct

$$P = \sum_{h=1}^{\infty} w_h \delta(\xi_h) \quad \text{and} \quad P' = \sum_{h=1}^{\infty} w_h \delta(\xi'_h),$$

with the same stick-breaking sequence of probabilities w_h as in Exercise [xx 15 xx]. by construction, $P \sim \text{Dir}(aP_0)$ and $P' \sim \text{Dir}(aP'_0)$, and they are dependent. a quick illustration.

17. Quantile inference for the Dirichlet process

[xx something here. nils and sonia. cute formula

$$\widehat{Q}_0(y) = \sum_{i=1}^n \binom{n-1}{i-1} y^{i-1} (1-y)^{n-i} x_{(i)}.$$

start from $F \sim \text{Dir}(aF_0)$, and consider the random quantile function

$$Q(y) = \min\{x : F(x) \geq y\}.$$

find a clear expression for the distribution of $Q(y)$. check case of $a \rightarrow 0$ separately. also the resulting cute enough nonparametric automatic bandwidth-free density estimator $\widehat{f}_0(x)$. xx]

18. Quantile pyramids

[xx something here. Hjort and Walker (2009, Annals) and their quantile pyramids. first construction, then MCMC for posterior. xx]

19. Brownian motion via convergence of a partial-sum process

Here I briefly describe the construction of Brownian motion as the proper limit in distribution of an empirical partial-sum process. This is of interest in its own right, as it also gives a proof of the existence of the relevant Gaussian process. The point is also that similar constructions (where 'similar' could mean 'very similar' or 'somewhat similar' or 'long-distance part-time similar') in different straits are useful for Bayesian Nonparametrics, such as the Gamma and Beta processes down the road.

The Brownian motion, or Wiener process, say $W = \{W(t) : t \geq 0\}$, is a Gaussian process (all finite-dimensional distributions are Gaussian), with mean zero, independent increments, with $W(t) - W(s) \sim N(0, t - s)$. The existence of such a process is a non-trivial delicate matter, but the construction I give below has 'yes, the Brownian motion process exists' as a by-product.

We start from an i.i.d. sequence $\varepsilon_1, \varepsilon_2, \dots$, with mean zero and variances one, and then build the empirical process

$$Z_n(t) = (1/\sqrt{n}) \sum_{i/n \leq t} \varepsilon_i.$$

- (a) Above we define Brownian motion via the property that the independent increments have $N(0, t - s)$ distributions. Prove that if we somehow had started with $N(0, |t - s|^\gamma)$ distributions instead, for some $\gamma \neq 1$ for the variances, then things would quickly backfire, turning

the universe into massive cognitive dissonance. The Kolmogorov coherence theorem way of proving existence of Brownian motion indeed starts with checking that coherence matters are in order.

- (b) Verify that $Z_n = \{Z_n(t) : t \geq 0\}$ has independent increments, mean zero, and variance $[nt]/n$, where $[nt]$ is the integer part of nt (so $[17] = 17$, $[17.01] = 17$, $[17.99] = 17$, etc.). Show also that

$$\text{Var} \{Z_n(t) - Z_n(s)\} = (1/n) \text{Var} \sum_{s < i/n \leq t} \varepsilon_i = [nt]/n - [ns]/n \rightarrow t - s,$$

for each $s < t$.

- (c) Show that, for each t , we have $Z_n(t) \rightarrow_d N(0, t)$. This is essentially the central limit theorem at work.
- (d) For $t_1 < \dots < t_k$, show that the vector of random differences

$$(Z_n(t_1), Z_n(t_2) - Z_n(t_1), \dots, Z_n(t_k) - Z_n(t_{k-1}))$$

tends to the distribution of (D_1, \dots, D_k) , where these are independent, with $D_j \sim N(0, t_j - t_{j-1})$ (writing also $t_0 = 0$).

- (e) Use this to verify that indeed

$$(Z_n(t_1), \dots, Z_n(t_k)) \rightarrow_d (W(t_1), \dots, W(t_k)),$$

for any $t_1 < \dots < t_k$.

- (f) The theory of convergence of probability measures, see e.g. the classic Billingsley (1968, Wiley), tells us that (d) is necessary, but sufficient, for properly proving that $Z_n \rightarrow_d W$, in the space $D[0, \tau]$ of all right-continuous functions $x : [0, \tau] \rightarrow \mathcal{R}$ with left-hand limits, and equipped with the Skorokhod topology. We do not go into the details here, but the added necessity factor is that of *tightness*, a condition that secures that the Z_n does not have any mass escaping away, or turning itself into too high oscillation with too high probability. Mathematically, tightness of the Z_n sequence means that for each $\varepsilon > 0$, there should exist a compact set A such that $\Pr\{Z_n \in A\} \geq 1 - \varepsilon$ for all large n . A condition securing such tightness, which again secures what we're after, namely full process convergence $Z_n \rightarrow_d W$, is

$$\text{E} \{Z_n(t) - Z_n(s)\}^2 \{Z_n(u) - Z_n(t)\}^2 \leq \{K(u) - K(s)\}^2, \quad (0.1)$$

for all triples $s < t < u$, for a suitable monotone, continuous function K . Verify this condition here.

- (g) For other application I note the following variations of this sufficient condition for tightness of a sequence Z_n of processes in $D[0, 1]$. First, it sufficient that

$$\text{E} |Z_n(t) - Z_n(s)|^a |Z_n(u) - Z_n(t)|^a \leq \{K(u) - K(s)\}^b \quad (0.2)$$

holds, for some $a > 0$ and $b > 1$, for all triples $s < t < u$, again with a monotone continuous function K . Often this is easiest to work with for $a = 2$, as above. Second, the following variation is *not* written out in the course of the classic text Billingsley (1968, Chs. 3-4), but

I've found it useful in several applications, and it follows essentially from the proof he gives for his Theorem 15.6. The variation is that it suffices to have

$$E |Z_n(t) - Z_n(s)|^a |Z_n(u) - Z_n(t)|^a \leq \{K_n(u) - K_n(s)\}^b \quad (0.3)$$

for all triples $s < t < u$, again with $a > 0$ and $b > 1$, where K_n is a monotone function converging pointwise to a continuous K .

- (h) Note that the above construction and reasoning hold, regardless of the actual distribution of the building blocks $\varepsilon_1, \varepsilon_2, \dots$. In particular, we may take the two very different start distributions $\varepsilon_i \sim N(0, 1)$, or ε_i equal to 1 or -1 with equal probability $\frac{1}{2}$, and have the same Brownian limit. Simulate some paths, from these two partial-sum processes, for say $n = 1000$, and check if you can tell the difference.
- (i) When this cornerstone theorem is in place, there is a long list of implications and corollaries and new constructions. Let me mention the Brownian bridge, $W^0 = \{W^0(t) : 0 \leq t \leq 1\}$. It is Gaussian, zero mean, and covariance function $\text{cov}\{W^0(s), W^0(t)\} = s(1-t)$ for $s \leq t$. It emerges in several ways, including starting from the Wiener process and then forming

$$W^0(t) = W(t) - tW(1) \quad \text{for } 0 \leq t \leq 1.$$

It is also the limit of the bridged version of the above empirical process,

$$Z_n^0(t) = Z_n(t) - tZ_n(1) \quad \text{for } 0 \leq t \leq 1.$$

- (j) [xx perhaps a few illustrations of the 'invariance theorem' aspect of this Donsker theorem. xx]

20. A little lemma

We shall encounter situations involving long products of the type $a_n = \prod_{i \leq n} (1 + z_{n,i})$, where there for each n is a well-defined sequence of $z_{n,i}$ for $i = 1, \dots, n$. If these are small and their sum converges, the sequence of products will converge. Specifically, assume

- (i) that $\sum_{i \leq n} z_{n,i} \rightarrow z$;
- (ii) that $\delta_n = \max_{i \leq n} |z_{n,i}| \rightarrow 0$;
- (iii) that $\sum_{i \leq n} |z_{n,i}|$ remains bounded.

Show that then $a_n = \prod_{i \leq n} (1 + z_{n,i}) \rightarrow a = \exp(z)$. It is helpful here to write

$$\log(1 + z) = z - \frac{1}{2}z^2 + z^2K(z),$$

where $|K(z)| \leq \frac{1}{2}$ for all $|z| \leq \frac{1}{2}$.

Similar results also hold when the product is taken over suitable subsets of i/n , like

$$\prod_{s < i/n \leq t} (1 + z_{n,i}) \rightarrow \exp(z_{s,t}),$$

if $\sum_{s < i/n \leq t} z_{n,i} \rightarrow z_{s,t}$, etc.

21A. Time-discrete Gamma and Poisson processes

As we've seen in Exercise 1, if $y | \theta \sim \text{Pois}(\theta)$, and $\theta \sim \text{Gamma}(a, b)$, then very nicely and conveniently $\theta | y \sim \text{Gamma}(a + y, b + 1)$. This may be lifted without too many problems to *time-discrete Gamma processes* with accompanying *time-discrete Poissons processes*. Key points are that both the Poisson and the Gamma have additive properties – a sum of four independent Poissons is a Poisson, and a sum of four independent $\text{Gamma}(a_j, b_j)$ is a Gamma, but only provided the b_j parameters are identical. Later on we will delve into continuous limits of these constructions.

- (a) Suppose we have a sequence of time-points, perhaps a long string $t_1 < \dots < t_m$, where count variables v_1, \dots, v_m are observed. First consider a vector $G = (G_1, \dots, G_m)$ of independent Gamma contributions, associated with these time-points, with $G_j \sim \text{Gamma}(a_j, b)$. Form from these the process of cumulative sums,

$$Z = (Z(t_1), \dots, Z(t_m)) \quad \text{with} \quad Z(t_j) = G_1 + \dots + G_j \quad \text{for } j = 1, \dots, m.$$

Show that $Z(t_j)$ is a Gamma, with parameters $(a_1 + \dots + a_m, b)$. Put up expressions for the mean and variance of $Z(t_j)$.

- (b) Then, given such an underlying process of cumulative intensities, suppose $V = (V_1, \dots, V_m)$ is a vector of independent Poisson variables, with parameters G_1, \dots, G_m . Form from these the increasing process of cumulative sums,

$$Y = (Y(t_1), \dots, Y(t_m)) \quad \text{with} \quad Y(t_j) = V_1 + \dots + V_j \quad \text{for } j = 1, \dots, m.$$

Show that $Y(t_j)$, conditional on the rate parameters, is Poisson with parameter $Z(t_j)$.

- (c) Show also that

$$Z(t_j) | \text{data} \sim \text{Gamma}(a_1 + \dots + a_j + Y(t_j), b + 1) \quad \text{for } j = 1, \dots, m.$$

Put up expressions for the posterior mean and variance for the cumulative intensity Z process.

- (d) Set up a simulation experiment, with say $m = 100$, and some suitable choice for a_1, \dots, a_m and b . Simulate (i) the cumulative intensity process and (ii) Poisson outcomes, and then give a plot for $\hat{A}(t_j) = E\{Z(t_j) | \text{data}\}$, along with a $[0.05, 0.95]$ credibility band.
- (e) How can the setup of this exercise be extended to (i) a continuous time cumulative intensity Gamma process; (ii) a continuous time Poisson process; (iii) the posterior process for the cumulative intensity function?

21B. The Gamma process

For a given monotone function $M(t)$, starting at $M(0) = 0$, we may define a Gamma process $Z = \{Z(t) : t \geq 0\}$ with the property that it has independent increments with $Z(t) - Z(s) \sim \text{Gamma}(M(t) - M(s), 1)$. Existence of such a process is not entirely obvious, but one is of course helped by the fact that

$$\text{Gamma}(M(t) - M(s), 1) + \text{Gamma}(M(u) - M(t), 1) \sim \text{Gamma}(M(u) - M(s), 1)$$

for $s < t < u$, with the two components on the left hand side being independent.

The purpose of this exercise is to work through some of the crucial details for the Gamma process, which also opens the door for more general constructions later on, like the extended Gamma process in the next exercise.

- (a) Let $G \sim \text{Gamma}(a, b)$, with density proportional to $x^{a-1} \exp(-bx)$. Show that its Laplace transform may be written as

$$E \exp(-uG) = \frac{b^a}{\Gamma(a)} \frac{\Gamma(a)}{(b+u)^a} = \frac{1}{(1+u/b)^a} = \exp\{-a \log(1+u/b)\}.$$

- (b) Use this to show that if G_1, \dots, G_m are independent Gamma distributed variables, with parameters $(a_1, b), \dots, (a_m, b)$, then their sum is also Gamma distributed, with parameters $(\sum_{i=1}^m a_i, b)$.

- (c) Show that the negative exponent in the Laplace transform can be expressed as

$$a \log(1+u/b) = \int_0^\infty \{1 - \exp(-us)\} dL(s),$$

with

$$dL(s) = as^{-1} \exp(-bs) ds.$$

- (d) Suppose as above that $M(t)$ is monotone, with $M(0) = 0$; in various applications, it will be a cumulative intensity function and of the form $M(t) = \int_0^t a(s) ds$, with an underlying nonnegative intensity function $a(s)$. Consider the process

$$Z_m(t) = \sum_{i/m \leq t} G_{m,i} \quad \text{for } t \geq 0,$$

where the $G_{m,i}$ are independent, and $G_{m,i} \sim \text{Gamma}(a_{m,i}, b)$, with $a_{m,i} = M(i/m) - M((i-1)/m)$ for $i \geq 1$. For the case of M being the integral of a , it is useful to think of $a_{m,i}$ as $a(i/m)(1/m)$. Show that the mean and variance converge properly,

$$E Z_m(t) \rightarrow M(t)/b \quad \text{and} \quad \text{Var } Z_m(t) \rightarrow M(t)/b^2.$$

- (e) Show that the Laplace transform converges,

$$E \exp\{-uZ_m(t)\} = \prod_{i/m \leq t} \exp\{-a_{m,i} \log(1+u/b)\} \rightarrow \exp\{-M(t) \log(1+u/b)\}.$$

This establishes existence of the Gamma process with parameter $(M(\cdot), b)$, via process convergence arguments as in Exercise [xx 19 and more xx]. Show in fact that the Z_m sequence is tight, in the appropriate function space $D[0, \tau]$ of all right-continuous functions $x: [0, \tau] \rightarrow R$ with left-hand limits, equipped with the Skorohod topology, using conditions of the type (0.3).

- (f) Show that the arguments above also work in the case where the underlying $M(\cdot)$ function is replaced by a $M_m(\cdot)$ function, which converges to a limit $M(\cdot)$. In particular, things go through for the case of $M_m(t) = \sum_{i/m \leq t} a(i/m)(1/m)$, tending to $\int_0^t a(s) ds$.

- (g) [xx a bit more xx]

22. The Extended Gamma process

In the course of this exercise I build a more general process, which I term an Extended Gamma process. [xx which has been worked with earlier, actually; find one or two references, from the two Walker students. xx] We start with independent and inherently small gammas,

$$G_{m,i} \sim \text{Gamma}(a(i/m)(1/m), b(i/m)) \quad \text{for } i = 1, 2, \dots,$$

where $a(s)$ and $b(s)$ are functions, taken positive and continuous, or at least piecewise continuous, and with $b(s)$ bounded above zero. From these we form the partial sum process

$$Z_m(t) = \sum_{i/m \leq t} G_{m,i} \quad \text{for } t \geq 0.$$

Below we demonstrate that it has a proper limit in distribution, say $Z = \{Z(t) : t \geq 0\}$, which we term the Extended Gamma process with parameter functions $(a(t), b(t))$. A special case is of course that of $b(t) = b$ constant, in which case we have a Gamma process with $Z(t) \sim \text{Gamma}(A(t), b)$, where $A(t) = \int_0^t a(s) ds$, i.e. as in Exercise 21.

(a) Show that the mean and variance converge,

$$\begin{aligned} \mathbb{E} Z_m(t) &= \sum_{i/m \leq t} \frac{a(i/m)(1/m)}{b(i/m)} \rightarrow \int_0^t \frac{a(s)}{b(s)} ds, \\ \text{Var } Z_m(t) &= \sum_{i/m \leq t} \frac{a(i/m)(1/m)}{b(i/m)^2} \rightarrow V(t) = \int_0^t \frac{a(s)}{b(s)^2} ds. \end{aligned}$$

(b) Show that the Laplace transform converges properly:

$$\mathbb{E} \exp\{-uZ_m(t)\} = \prod_{i/m \leq t} \mathbb{E} \exp(-uG_{m,i}) = \exp\left[-\sum_{i/m \leq t} a(i/m)(1/m) \log\{1 + u/b(i/m)\}\right],$$

which indeed tends to

$$\exp\left[-\int_0^t a(s) \log\{1 + u/b(s)\} ds\right].$$

(c) With $V(t)$ as above, work with bounds for both

$$\mathbb{E}\{Z_m(t) - Z_m(s)\}\{Z_m(u) - Z_m(t)\}$$

and

$$\mathbb{E}\{Z_m(t) - Z_m(s)\}^2\{Z_m(u) - Z_m(t)\}^2,$$

for triples $s < t < u$, and use tightness criteria of Exercise 19 to establish full process convergence $Z_m \rightarrow_d Z$.

(d) [xx a bit more. the Lévy representation of things. xx]

$$\mathbb{E} \exp\{-uZ(t)\} = \exp\left[-\int_0^\infty \{1 - \exp(-us)\} dL_t(s)\right].$$

23. The Extended Gamma process with a Poisson process

[xx part of the nils-emil story. pointer to later exercise with covariates. pointer also to Beta process version of things. xx] I start with the time-discrete version of things. Consider a sequence of independent pairs $(\theta_{m,i}, z_{m,i})$, to be thought of as evolving over time points i/m , with

$$\theta_{m,i} \sim \text{Gamma}(a(i/m)(1/m), b(i/m)) \quad \text{and} \quad z_{m,i} | \theta_{m,i} \sim \text{Pois}(\theta_{m,i}).$$

In particular, $G_m(t) = \sum_{i/m \leq t} \theta_{m,i}$ is the cumulative intensity process, and $Z_m(t) = \sum_{i/m \leq t} z_{m,i}$ the cumulative Poisson count of events.

(a) Show that

$$\theta_{m,i} | \text{data} \sim \text{Gamma}(a(i/m)(1/m) + z_{m,i}, b(i/m) + 1).$$

(b) Then consider the time-continuous version of this story, corresponding to letting $m \rightarrow \infty$ above. This leads to $G_m \rightarrow_d G$, an Extended Gamma process, with parameter functions $a(s), b(s)$. Then, given G , there is a limit $Z_m \rightarrow Z$, a nonhomogeneous observed Poisson process $Z = \{Z(t) : t \geq 0\}$, with cumulative intensity function G . Show that G given data is another Extended Gamma process, with $b_{\text{new}}(s) = b(s) + 1$, and $A_{\text{new}}(t) = \int_0^t a(s) ds + Z(t)$. This translates to

$$dG(s) | \text{data} \sim \text{Gamma}(a(s) ds + dZ(s), b(s) + 1),$$

with $dZ(s) = Z[s, s + ds]$ the number of Poisson events observed in the small time window $[s, s + ds]$.

(c) In particular, writing event times as $T_1 < T_2 < \dots$, show that the posterior mean becomes

$$\hat{G}(t) = \int_0^t \frac{a(s) ds + dZ(s)}{b(s) + 1} = \int_{\text{no jumps}} \frac{a(s)}{b(s) + 1} ds + \sum_{\text{jumps} \leq t} \frac{1}{b(T_j) + 1}.$$

Find also an expression for the posterior variance.

(d) Suppose next that there are several observed nonhomogeneous Poisson processes, say Z_1, \dots, Z_k , with the same underlying G . Show that G given the data is again an Extended Gamma process, with

$$dG(s) | \text{data} \sim \text{Gamma}\left(a(s) ds + \sum_{j=1}^k dZ_j(s), b(s) + k\right).$$

(e) With such observed processes, show that the Bayes estimator for the cumulative intensity process, i.e. the posterior mean, becomes

$$\hat{G}(t) = \int_0^t \frac{a(s) ds + \sum_{j=1}^k dZ_j(s)}{b(s) + k} = \int_{\text{between jumps}} \frac{a(s)}{b(s) + k} ds + \sum_{\text{jumps} \leq t} \frac{1}{b(T_j) + k},$$

now with ‘jumps’ referring to jumps in any of the k observed nonhomogeneous Poisson processes. Also, give a formula for the posterior variance.

23B. Bayesian inference for a Poisson intensity process

Here you are to illustrate the setup with a Gamma intensity process and Poisson data.

- (a) Consider a cumulative intensity Gamma process $Z(t)$ on the time interval $[0, 1]$, where $Z(t) \sim \text{Gamma}(cA(t), c)$, with $A(t)$ the integral of the intensity function $a(s) = 10 + 0.75 \cos(2\pi s)$. What is the mean and variance of $Z(t)$? Simulate and display say 25 $Z(t)$ processes, using a small, a moderate, and a large value of c .
- (b) Then fake a data Poisson process $Y = \{Y(t) : t \in [0, 1]\}$ in your computer, taking $Y(t) \sim \text{Pois}(Z_{\text{true}}(t))$, with $Z_{\text{true}}(t) = \int_0^t z_{\text{true}}(s) ds$, where $z_{\text{true}}(s) = 10 + 0.05 \cos(2\pi s)$.
- (c) Find the posterior process $Z(t)$ given the observed Y process. Compute and display the Bayes estimator $\hat{Z}(t) = E\{Z(t) | \text{data}\}$, along with a credibility band with width proportional to the posterior standard deviation.

- (d) Draw say 25 simulations of $Z(t)$ given your data.
- (e) Generalise the setup, in a couple of directions, where one such direction is that of having more data. Assume, for example, that the same Poisson intensity process $Z(t)$ is at work, for each of $k = 10$ occasions, leading to observed Poisson counting processes Y_1, \dots, Y_k . Find $\widehat{Z}(t) = E\{Z(t) | \text{data}\}$ and the posterior standard deviation function, etc.

24. The biggest jumps of a Gamma process

[xx to be polished. xx] Consider a Gamma process $Z = \{Z(t) : t \geq 0\}$, at first with plain linear mean function, so that $Z(t) \sim \text{Gamma}(at, 1)$. It has jumps, in fact infinitely many jumps on each time interval. Consider $J(t)$, the biggest of these jumps in the course of the time interval $[0, t]$. We shall find its distribution.

- (a) We start with a little investigation of the size of $X_\varepsilon \sim \text{Gamma}(\varepsilon, 1)$, with ε small. Writing Γ_ε for its cdf, i.e. $\Gamma_\varepsilon(v) = G(v, \varepsilon, 1)$ in terms of the Gamma distribution cdf with parameters $(\varepsilon, 1)$, show that

$$\begin{aligned} \Gamma_\varepsilon(v) &= 1 - \int_v^\infty \frac{1}{\Gamma(\varepsilon)} x^{\varepsilon-1} \exp(-x) dx \\ &= 1 - \varepsilon \int_v^\infty \frac{1}{\Gamma(1+\varepsilon)} x^\varepsilon x^{-1} \exp(-x) dx = 1 - \varepsilon E_1(v) \{1 + O(\varepsilon)\}, \end{aligned}$$

in terms of the exponential integral function

$$E_1(v) = \int_v^\infty (1/x) \exp(-x) dx.$$

- (b) We know that $Z(\cdot)$ can be seen as the limit in distribution of the process $Z_m(t) = \sum_{i/m \leq t} G_{m,i}$, with independent components $G_{m,i} \sim \text{Gamma}(a/m, 1)$. Work with the biggest of these, say $J_m(t)$, and show

$$\Pr\{J_m(t) \leq v\} = \Pr\{G_{m,i} \leq v \text{ for all } i/m \leq t\} = \prod_{i/m \leq t} \Gamma_{a/m}(v).$$

- (c) That $J_m(t) \rightarrow_d J(t)$ is intuitively correct; some finer details regarding this are in Hjort and Ongaro (2006). Hence the distribution of $J(t)$ can be found by taking the limit of the expression above. Use indeed the above to prove

$$\Pr\{J_m(t) \leq v\} \rightarrow \exp\{-at E_1(v)\} \quad \text{for } v > 0.$$

This is the sought-for result for the distribution of $J(t)$.

- (d) Next consider a Gamma process with $Z(t) \sim \text{Gamma}(aM(t), 1)$, for a monotone $M(\cdot)$ function. With $J(t)$ the biggest jump during $[0, t]$, show that

$$\Pr\{J(t) \leq v\} = \exp\{-aM(t)E_1(v)\} \quad \text{for } v > 0.$$

More generally, with Z and Extended Gamma process with parameter functions $(a(t), b(t))$, as in Exercise 22, show that the biggest jump over the time window $[t_1, t_2]$ has distribution

$$\Pr\{J(t_1, t_2) \leq v\} = \exp\left\{-\int_{t_1}^{t_2} \frac{a(s)}{b(s)} ds E_1(v)\right\}.$$

(e) A further variation of interest, also for model building aspects (which I intend to come back to in a later exercise), is as follows. For a Gamma process with $Z(t) \sim \text{Gamma}(A(t), 1)$, consider sizes of jumps with respect to a boundary, i.e. $\Delta Z(t)/v(t)$, with $\Delta Z(t)$ the jump size and $v(t)$ a function. I formalise this via $J_m(t) = \max_{i/m \leq t} G_{m,i}/v(i/m)$, and can then work with

$$\begin{aligned} \Pr\{J_m(t) \leq 1\} &= \prod_{i/m \leq t} \Pr\{G_{m,i} \leq v(i/m)\} \\ &= \prod_{i/m \leq t} [1 - a(i/m)(1/m)E_1(v(i/m))\{1 + O(1/m)\}] \end{aligned}$$

which is then seen to converge to

$$\Pr\{J(t) \leq 1\} = \exp\left\{-\int_0^t a(s)E_1(v(s)) ds\right\}.$$

25. The Beta process

Hjort (1985, SJS, invited discussion contribution to the SJS paper by P.K. Andersen and Ø. Borgan on counting process models) introduced the Beta process, used as a prior process for cumulative hazard functions, and gave the crucial conjugacy property when used for survival data. A fuller account was then given in Hjort (1990, Annals). The present exercise indicates how the Beta process can be constructed from a limit operation for a partial-sum process involving small Beta components.

We start with a function $a_0(s)$, intended to be like a prior guess hazard function, with cumulative $A_0(t) = \int_0^t a_0(s) ds$. For given m , let $B_{m,1}, B_{m,2}, \dots$ be independent Beta random variables, with

$$B_{m,i} \sim \text{Beta}\left(c\left(\frac{i}{m}\right)a_0\left(\frac{i}{m}\right)\frac{1}{m}, c\left(\frac{i}{m}\right) - c\left(\frac{i}{m}\right)a_0\left(\frac{i}{m}\right)\frac{1}{m}\right).$$

Here $c(s)$ is a positive function, with at most finitely many discontinuities; it may e.g. be a constant. Our process is

$$A_m(t) = \sum_{i/m \leq t} B_{m,i} \quad \text{for } t \geq 0.$$

(a) Show that

$$E Z_m(t) = \sum_{i/m \leq t} a_0(i/m)(1/m) \rightarrow A_0(t).$$

Show also that

$$\text{Var } A_m(t) = \sum_{i/m \leq t} \frac{a_0(i/m)(1/m)\{1 - a_0(i/m)(1/m)\}}{c(i/m) + 1} \rightarrow \int_0^t \frac{a_0(s) ds}{c(s) + 1}.$$

(b) Hjort (1985, 1990) proves that A_m really converges to a well-defined limit process $A = \{A(t) : t \geq 0\}$, with independent increments all inside $[0, 1]$, and calls this the Beta process, with parameters (c, A_0) . Proving convergence and existence of this limit process takes some care and tools from empirical processes. The crucial point here is that the Laplace transform has a well-defined limit, so let us work with

$$E \exp\{-uA_m(t)\} = \prod_{i/m \leq t} E \exp(-uB_{m,i}) = \prod_{i/m \leq t} (1 + z_{m,i}),$$

say. We must then work hard enough with the $z_{m,i}$ to be able to apply the Little Lemma of Exercise XX. Show via Beta moments that

$$\begin{aligned} \mathbb{E} \exp(-uB_{m,i}) &= 1 + z_{m,i} \\ &= 1 + \sum_{j=1}^{\infty} (-1)^j \frac{u^j}{j!} \frac{\Gamma(c(i/m))}{\Gamma(c(i/m)a_0(i/m)(1/m))} \frac{\Gamma(c(i/m)a_0(i/m)(1/m) + j)}{\Gamma(c(i/m) + j)}. \end{aligned}$$

(c) Then use

$$\Gamma(\varepsilon + j)/\Gamma(\varepsilon) = \varepsilon(\varepsilon + 1) \cdots (\varepsilon + j - 1) = (j - 1)! \varepsilon + O(\varepsilon^2)$$

for small ε to deduce

$$\mathbb{E} \exp(-uB_{m,i}) = 1 + \sum_{j=1}^{\infty} (-1)^j \frac{u^j}{j} \frac{\Gamma(c(i/m))}{\Gamma(c(i/m) + j)} c(i/m)a_0(i/m)(1/m) + O(1/m^2).$$

(d) Show that this leads to

$$\mathbb{E} \exp\{-uA_m(t)\} \rightarrow \exp\left\{-\int_0^t \sum_{j=1}^{\infty} (-1)^j \frac{u^j}{j!} \frac{\Gamma(c(z))\Gamma(j)}{\Gamma(c(z) + j)} a_0(z) dz\right\}.$$

(e) [xx more xx] Link to Lévy representation

$$\int_0^1 \{1 - \exp(-us)\} dL_t(s).$$

– The above establishes the existence of a Beta process, with parameters (c, A_0) ; for a fuller discussion, see Hjort (1990, Annals). It is a independent and nonnegative increments, and these are all in $[0, 1]$. The intuitive interpretation for a Beta process is that

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

These tiny increments are not exactly Beta distributed, though; that distribution does not have any easy convolution properties, unlike e.g. the Gamma.

26. A time-discrete framework for survival analysis

Consider the following framework for life-times, now with time-discrete outcomes in $\{0, 1, 2, \dots\}$, rather than the usual time-continuous setup of $[0, \infty)$. A random variable T then has probability masses

$$f_j = \Pr\{T = j\} \quad \text{for } j = 0, 1, 2, \dots,$$

with cumulative $F_j = \Pr\{T \leq j\} = f_0 + f_1 + \dots + f_j$. It is also very fruitful to work with the hazards

$$\alpha_j = \Pr\{T = j | T \geq j\} = f_j / (f_j + f_{j+1} + \dots),$$

along with the cumulative hazards $A_j = \alpha_0 + \alpha_1 + \dots + \alpha_j$.

Part of what I present in this exercise was also included in Hjort (1990, Annals, along with further results, extensions, and discussion). The framework and methods for the time-discrete setup have separate interest, and it inspired the invention of the Beta process, as a fine limit of the time-discrete grid.

(a) Show that

$$F_j = 1 - \prod_{k=0}^j (1 - \alpha_k), \quad \text{for } j \geq 0.$$

(b) Show then that

$$f_j = (1 - \alpha_0)(1 - \alpha_1) \cdots (1 - \alpha_{j-1})\alpha_j,$$

and give an interpretation of this identity.

- Assume now that we have observations (t_i, δ_i) for $i = 1, \dots, n$, for different individuals, with $\delta_i = 1$ if the life-time is observed and $\delta_i = 0$ if there merely is censored information that the real life-time is larger than t_i . From these, define

$$Y_j = \sum_{i=1}^n I\{t_i \geq j\} \quad \text{and} \quad N_j = \sum_{i=1}^n I\{t_i = j, \delta_i = 1\},$$

the at-risk counter and the counting process of observed life-times. In particular, let ΔN_j be the jump of N at time point j , the number of the Y_j at risk who experience a transition at time j .

(d) Show, with the necessary efforts of bureaucratic book-keeping, that the likelihood information from such a dataset can be expressed as

$$L = \prod_{i: \delta_i=1} f_i \prod_{i: \delta_i=0} (1 - F_i) = \prod_{j=0}^{\infty} (1 - \alpha_j)^{Y_j - \Delta N_j} \alpha_j^{\Delta N_j}.$$

Discuss how or to what extent this can be interpreted as a succession of binomial trials, with $\Delta N_j | Y_j$ a binomial (Y_j, α_j) .

(e) The representation above invites the idea of independent Beta priors for the hazards. Let in fact $\alpha_j \sim \text{Beta}(c_j \alpha_{j,0}, c_j - c_j \alpha_{j,0})$, for $j = 0, 1, 2, \dots$, and deduce that these are independent and Beta distributed also given data, with updated parameters

$$\alpha_j | \text{data} \sim \text{Beta}(c_j \alpha_{j,0} + \Delta N_j, c_j - c_j \alpha_{j,0} + Y_j - \Delta N_j).$$

(f) Show that the Bayes estimator for the cumulative hazard function is

$$\hat{A}_j = E(A_j | \text{data}) = \sum_{k=0}^j \frac{c_k \alpha_{k,0} + \Delta N_k}{c_k + Y_k}.$$

The noninformative case of the c_j becoming small leads to $\sum_{k=0}^j \Delta N_k / Y_k$, a time-discrete version of the Nelson–Aalen estimator (see Exercise [xx ... xx]).

(g) Then show that the Bayes estimator for the survival function $\Pr\{T > j\} = \prod_{k=0}^j (1 - \alpha_k)$ is

$$\hat{S}_j = \prod_{k=0}^j \left(1 - \frac{c_k \alpha_{k,0} + \Delta N_k}{c_k + Y_k}\right).$$

For the noninformative case of the $c_j \rightarrow 0$, we find $\prod_{k=0}^j (1 - \Delta N_k / Y_k)$, a time-discrete version of the Kaplan–Meier estimator (see again Exercise [xx ... xx]).

27. The Beta process for survival data

[xx need polish. xx] Suppose a survival dataset of the usual form (t_i, δ_i) is available, to inform us about an underlying survival distribution F on $[0, \infty)$. As per tradition, $\delta_i = 1$ is an indicator for non-censoring, and means a fully observed life-time, whereas $\delta_i = 0$ means that the life-time involved is censored, but one knows that it is larger than t_i . The survival distribution is $S(t) = 1 - F(t) = \Pr\{T > t\}$, and the cumulative hazard rate is

$$A(t) = \int_0^t dA(s), \quad \text{where} \quad dA(s) = \frac{dF(s)}{F[s, \infty)} = \Pr\{T \in [s, s + ds] | T \geq s\}.$$

In Aalen–Borgan notation, consider the at-risk counter and the counting process of observed life-times,

$$Y(s) = \sum_{i=1}^n I\{t_i \geq s\} \quad \text{and} \quad N(t) = \sum_{i=1}^n I\{t_i \leq t, \delta_i = 1\}.$$

In particular, $dN(s)$ is 1 if a life-time has been observed in $[s, s + ds]$, and 0 if not. The famous Nelson–Aalen and perhaps even more famous Kaplan–Meier estimator, for the cumulative hazard rate and the survival curve, are

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

- (a) Let $A \sim \text{Beta}(c, A_0)$, with A_0 the prior mean and c the strength function. Try to show, perhaps using some intuitive arguments, based on the approximate prior distribution of $dA(s)$, that

$$dA(s) | \text{data} \approx_d \text{Beta}\{c(s) dA_0(s) + dN(s), c(s) - c(s) dA_0(s) + Y(s) - dN(s)\},$$

and that these increments must be independent.

- (b) Try, again perhaps using heuristic arguments, to show that this means that the posterior distribution of A is an updated Beta process,

$$A | \text{data} \sim \text{Beta}(c + Y, \widehat{A}),$$

with posterior mean function

$$\widehat{A}(t) = \int_0^t \frac{c dA_0 + dN}{c + Y}.$$

This is the basic conjugacy property for the Beta process with survival data, proven in Hjort (1990, Annals) – involving, he says, ‘heroic integrations’.

- (c) Use the product integral representation

$$F(t) = 1 - \prod_{[0, t]} \{1 - dA(s)\}$$

to find the posterior mean of the survival function,

$$\widehat{S}(t) = \prod_{[0, t]} \left\{1 - \frac{dN(s)}{Y(s)}\right\}.$$

- (d) Show that when the $c(s)$ function tends to zero, or if the data volume is relatively large compared to the $c(s)$, then we’re back to the Nelson–Aalen and Kaplan–Meier estimators.

- (e) Explain how one may simulate realisations of A and then S from the posterior distribution. This may then be used to read off what we might wish for from these, like the posterior median

$$\mu = \min\{t: F(t) \geq \frac{1}{2}\}.$$

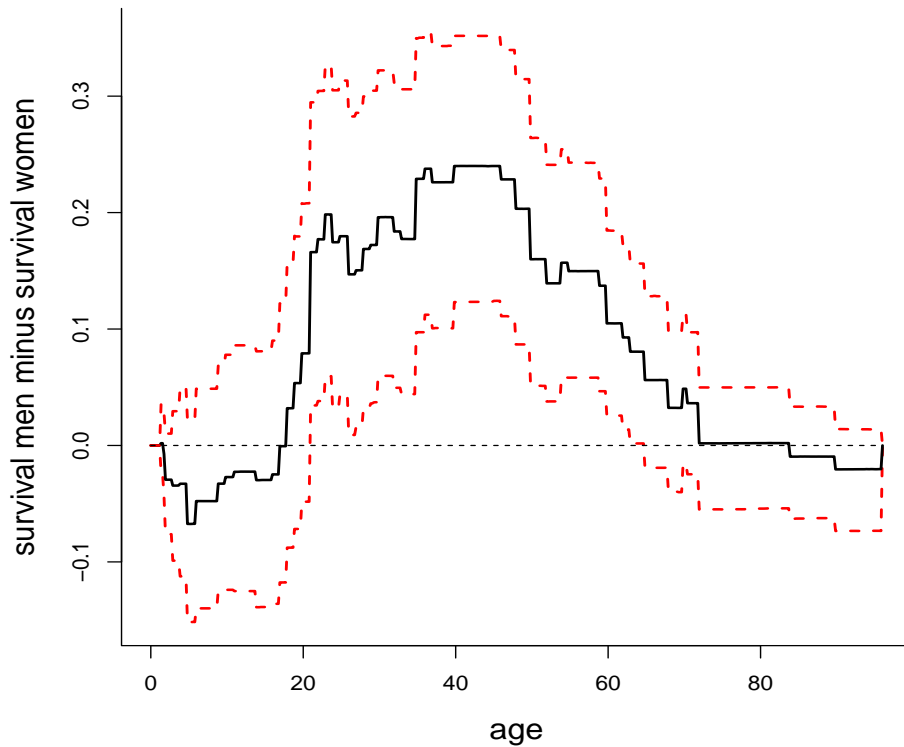


Figure 0.5: Roman Era Egypt lifetimes: difference $S_m(t) - S_w(t)$ between men's and women's survival curves, with 0.05, 0.50, 0.95 pointwise quantiles from the posterior distribution.

28. Lifelengths in Roman Era Egypt

[xx this to be polished. the exercise is too long and will be broken into two separate ones. xx]
 Access the `egypt-data` dataset from the course website, pertaining to the life-lengths of 82 men and 59 women from Roman Era Egypt, the 1st century b.C. This was a relatively peaceful society, without major wars, etc., and the life-lengths can be seen as having been sampled from the upper classes of that society. I've taken the data from the very first issue of *Biometrika* (1901), where Karl Pearson briefly discussed aspects of the life-lengths distribution, comparing this to Britain 1900. I've analysed aspects of these data both in Claeskens and Hjort (2008, Ch. 2) and in Schweder and Hjort (2016, Ch. xx).

Here we are interested in aspects of the underlying distributions F_w and F_m , for women and men, respectively, and, in particular, aspects where we might identify differences between the two. Let A_w and A_m be the cumulative hazard rate functions, along with survival curves

$$S_w(t) = \prod_{[0,t]} \{1 - dA_w(s)\} \quad \text{and} \quad S_m(t) = \prod_{[0,t]} \{1 - dA_m(s)\}. \quad (\text{eg1})$$

We use Beta process priors for the cumulative hazard rates, $A_w \sim \text{Beta}(c_w, A_{0,w})$ and $A_m \sim \text{Beta}(c_m, A_{0,m})$.

- (a) Assume for about two minutes that A_w and A_m are continuous functions. Then show from the product integrals that the familiar formulae

$$S_w(t) = \exp\{-A_w(t)\} \quad \text{and} \quad S_m(t) = \exp\{-A_m(t)\} \quad (\text{eg2})$$

emerge. With the Beta process priors to be used, however, there are discrete components, and we prefer (eg1) over (eg2), in terms of setup, modelling, prior to posterior, analysis, and interpretation. See also the general discussion regarding this point in Hjort (1990, Annals).

- (b) To make this concrete, choose the same Beta process prior for men and for women, with prior guess $A_0(t) = \int_0^t \alpha_0(s) ds$ corresponding to a Gamma with mean 30.00 and standard deviation 20.00, and then your own $c(s)$ strength function. Simulate realisations from A_w, A_m , and by implication S_w, S_m , on your screen.

- (c) Then update the Beta processes, given the data from the heroic Egyptian women and men, to say

$$A_w | \text{data} \sim \text{Beta}(c_w + Y_w, \hat{A}_w) \quad \text{and} \quad A_m | \text{data} \sim \text{Beta}(c_m + Y_m, \hat{A}_m).$$

In particular, compute and display both

$$\hat{A}_w(t) = \int_0^t \frac{c_w dA_0(s) + dN_w(s)}{c_w(s) + Y_w(s)} \quad \text{and} \quad \hat{A}_m(t) = \int_0^t \frac{c_m dA_0(s) + dN_m(s)}{c_m(s) + Y_m(s)},$$

and the survival curves

$$\hat{S}_w(t) = \prod_{[0,t]} \left\{ 1 - \frac{c_w(s) dA_0(s) + dN_w(s)}{c_w(s) + Y_w(s)} \right\} \quad \text{and} \quad \hat{S}_m(t) = \prod_{[0,t]} \left\{ 1 - \frac{c_m(s) dA_0(s) + dN_m(s)}{c_m(s) + Y_m(s)} \right\}.$$

- (d) Compute and display also the standard deviation curves, say $\hat{\kappa}_w(t)$ and $\hat{\kappa}_m(t)$ for A_w and A_m , and $\hat{\tau}_w(t)$ and $\hat{\tau}_m(t)$ for S_w and S_m .

- (e) Display the easy and simulation free approximate pointwise 90% confidence bands, of the type

$$\hat{A}_w(t) \pm 1.645 \hat{\kappa}_w(t) \quad \text{and} \quad \hat{A}_m(t) \pm 1.645 \hat{\kappa}_m(t),$$

and similarly for the survival curves. Crucially, in order to check the differences between the female and male populations, do this also for $A_w - A_m$ and $S_w - S_m$.

- (f) Then re-do the above point, without formulae, but via simulations from the posterior Beta processes.

- (g) This thing looks cool and relevant: Consider the survival curve ratio

$$\rho(t) = \frac{S_m(t)}{S_w(t)} = \prod_{[0,t]} \frac{1 - dA_m(s)}{1 - dA_w(s)}.$$

Find formulae for the prior and posterior mean of $\rho(t)$, and display the resulting $\hat{\rho}(t)$. Supplement this with a pointwise 90% credibility band, from simulations, or from conditional variances.

- (h) Summarise your findings properly. Yes, the women and the men of Roman Era Egypt had different life-length distributions. For which age interval is this most clear? And what could be the underlying mechanism or explanations?
- (i) [xx nils then includes a couple of Old Egyptian plots here. prior used is a $\text{Gamma}(a, b)$ with (a, b) chosen so that the mean is 30.0 years and the standard deviation is 20.0 years. investigate this prior mean choice. choose your own $c(s)$, perhaps reflecting less certainty for higher age values. check Figures 0.5 and 0.7, pertaining to the plain difference $S_m(t) - S_w(t)$ and the ratio $\rho(t) = S_m(t)/S_w(t)$ over time, with posterior median along with a 90% pointwise credibility band. the men were better off than women, for the age span 20 to 60. xx]

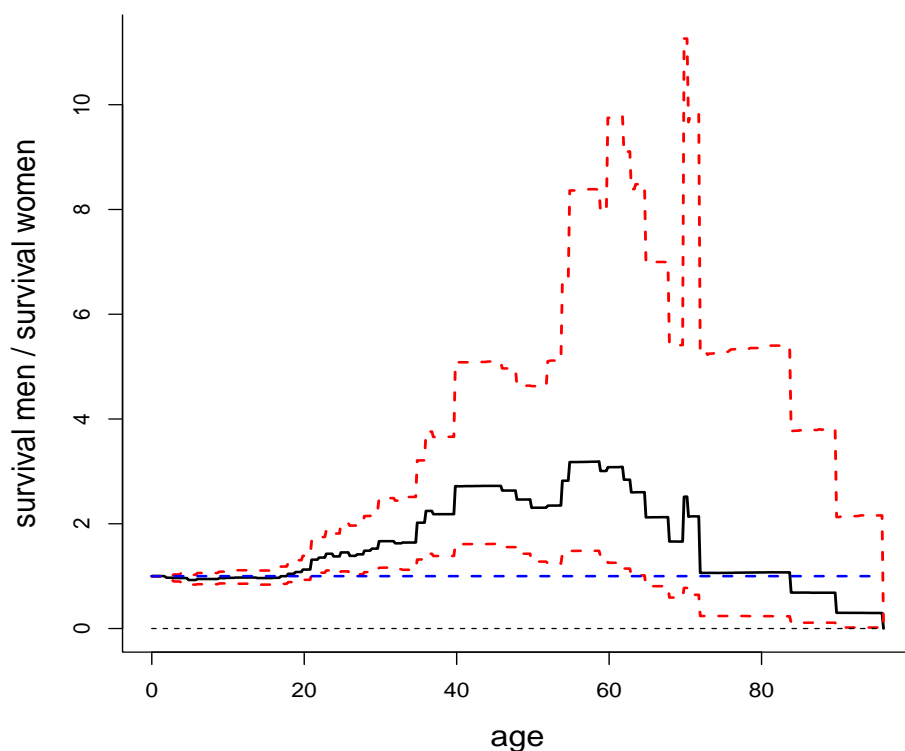


Figure 0.6: Roman Era Egypt lifetimes: ratio $S_m(t)/S_w(t)$ between men's and women's survival curves, with 0.05, 0.50, 0.95 pointwise quantiles from the posterior distribution.

29. The Bernoulli process and the Poisson process

We learn if not in kindergarten then perhaps in high school that a binomial (n, p) is close to a Poisson if n is big and p is small. This exercise exhibits generalisations of this basic result, leading also to a nonhomogeneous Poisson process limit of a suitably defined Bernoulli events process.

- (a) For $y \sim \text{Bin}(n, p)$, show that its Laplace transform is

$$L_n(u) = \text{E} \exp(-uy) = \{\exp(-u)p + 1 - p\}^n.$$

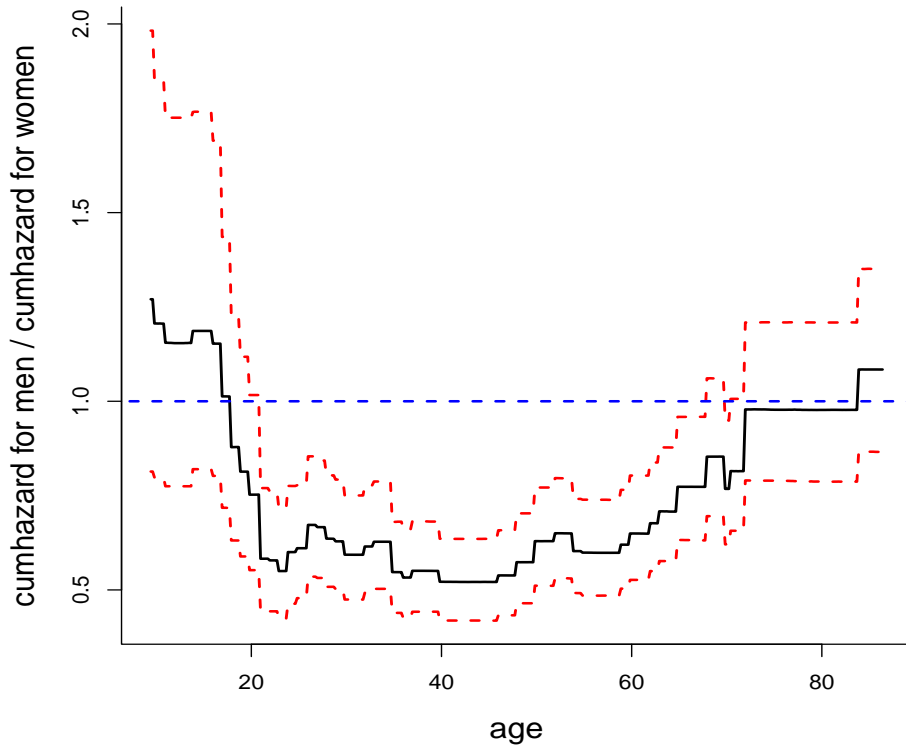


Figure 0.7: Roman Era Egypt lifetimes: ratio $A_m(t)/A_w(t)$ between men's and women's cumulative hazard functions, with 0.05, 0.50, 0.95 pointwise quantiles from the posterior distribution. We learn that the hazard rates are not proportional.

(b) Show also that when n increases and p decreases in such a way that $np \rightarrow \theta$, then

$$L_n(u) = [1 - p\{1 - \exp(-u)\}]^n \rightarrow \exp[-\theta\{1 - \exp(-u)\}].$$

Verify that this limit is the Laplace transform of a $\text{Pois}(\theta)$.

(c) Now study a Bernoulli event process, of the form

$$Z_m(t) = \sum_{i/m \leq t} B_{m,i} \quad \text{for } t \geq 0,$$

with independent Bernoulli components $B_{m,i} \sim \text{Bin}(1, a(i/m)(1/m))$. Here $a(s)$ is some nonnegative function, perhaps constant, perhaps evolving over time, and with cumulative intensity function $A(t) = \int_0^t a(s) ds$. Show that Z_m has independent increments, and that its Laplace transform converges,

$$\mathbb{E} \exp\{-uZ_m(t)\} = \prod_{i/n \leq t} [1 - a(i/m)(1/m)\{1 - \exp(-u)\}] \rightarrow \exp[-A(t)\{1 - \exp(-u)\}].$$

This means that the fine-grid Bernoulli process has properly converged, in the time-continuous limit, to a nonhomogeneous Poisson process, with $A(t) = \int_0^t a(s) ds$ as cumulative intensity process.

30. The jumps of a Gamma process via a Poisson process

In Exercise 24 we worked with the biggest jumps of a Gamma process over a given time window. Specifically, when $Z(t) \sim \text{Gamma}(A(t), 1)$ is such a process, with $J(t)$ the biggest jump over $[0, t]$, then $\Pr\{J(t) \leq v\} = \exp\{-A(t)E_1(v)\}$ for $v > 0$, with $E_1(v) = \int_v^\infty (1/u) \exp(-u) du$ the exponential integral function. The present exercise goes into a general Poisson process explanation for this result, leading also to more general results concerning the distribution of the 2nd biggest jump, the 3rd biggest jump, etc. These results again can be used for building survival analysis models, as in Exercise XX.

- (a) Let as in Exercise 24 $Z_m(t) = \sum_{i/m \leq t} G_{m,i}$, with the $G_{m,i} \sim \text{Gamma}(a(i/m)(1/m), 1)$ and independent. Consider the process counting occasions where the components are above threshold v , say

$$N_m(t) = \sum_{i/m \leq t} I\{G_{m,i} > v\}.$$

Show that

$$E N_m(t) = \sum_{i/m \leq t} \Pr\{G_{m,i} > v\} = \sum_{i/m \leq t} \{a(i/m)(1/m)E_1(v) + O(1/m^2)\} \rightarrow A(t)E_1(v),$$

with $A(t) = \int_0^t a(s) ds$.

- (b) Show also that the variance of $N_m(t)$ has the same limit $A(t)E_1(v)$.
- (c) Show indeed that $N_m(t)$ converges in distribution to a Poisson process $N(t)$ with independent increments and $N(t) \sim \text{Pois}(B_v(t))$, where $B_v(t) = A(t)E_1(v)$.
- (d) Give another proof, based on this, for the biggest jump $J(t)$ over $[0, t]$ having the distribution $\Pr\{N(t) = 0\} = \exp\{-B_v(t)\}$.
- (e) Let $J_{m,3}(t)$ be the 3rd biggest jump experienced by the Z_m process over the $[0, t]$ time window. Show that

$$\Pr\{J_{m,3}(t) \leq v\} \rightarrow \Pr\{N(t) \leq 2\} = \exp\{-B_v(t)\} \{1 + B_v(t) + \frac{1}{2}B_v(t)^2\}.$$

Generalise properly.

- (f) Let T_3 be the first time the 3rd biggest jump of the Gamma process Z exceeds threshold v . Show that the survival distribution becomes

$$S(t) = \Pr\{T_3 > t\} = \Pr\{J_3(t) \leq v\} = \Pr\{N(t) \leq 2\},$$

and that this corresponds to a cumulative hazard rate function

$$H_3(t) = B_v(t) - \log\{1 + B_v(t) + \frac{1}{2}B_v(t)^2\}.$$

- (g) Show that the hazard rate function can be written

$$h_3(t) = b_v(t) \frac{\frac{1}{2}B_v(t)^2}{1 + B_v(t) + \frac{1}{2}B_v(t)^2} = b_v(t)Q_v(t),$$

say, where $b_v(t) = a(t)E_1(v)$ is the basis hazard rate and $Q_v(t)$ a frailty correction function, climbing from zero to one as time increases.

31. The Beta process with a Bernoulli process

[xx to be written down. part of nils-emil story. xx] prior $A \sim \text{Beta}(c, A_0)$ for the cumulative intensity of a Bernoulli process Z . then

$$A \mid \text{data} \sim \text{Beta}(c + 1, \widehat{A}),$$

where

$$\widehat{A}(t) = \int_0^t \frac{c(s) dA_0(s) + dZ(s)}{c(s) + 1}.$$

with variation: extended Gamma. perhaps with marks or covariates. cross-ref to other exercises here.

32. The Gamma process, with a Poisson process, with covariates

[xx to be polished. part of the nils-emil story with police tweets, now with covariates. xx] In Exercise [xx 23 xx] we worked with the extended Gamma process as a prior G for the cumulative intensity function of nonhomogeneous Poisson processes, say Z_1, \dots, Z_k . The present exercise takes us through a certain statistically important generalisation, where covariate information is available for the Z_j event counting processes.

- (a) As a start generalisation of the framework of Exercise [xx 23 xx], suppose that there is a sequence of independent pairs $(\theta_{m,i}, z_{m,i})$, where

$$\theta_{m,i} \sim \text{Gamma}(a(i/m)(1/m), b(i/m)) \quad \text{and} \quad z_{m,i} \mid \theta_{m,i} \sim \text{Pois}(w(i/m)\theta_{m,i}).$$

At the moment, the $w(s)$ is to be thought of as a given function, producing multiplicative Poisson intensity factors $w(i/m)$. Show that this leads to

$$\theta_{m,i} \mid \text{data} \sim \text{Gamma}(a(i/m)(1/m) + z_{m,i}, b(i/m) + w(i/m)).$$

- (b) In the time-continuous limit, with $G_m(t) = \sum_{i/m \leq t} \theta_{m,i}$ tending to a cumulative intensity process $G(t)$, and the nonhomogeneous Poisson counting process $Z_m(t) = \sum_{i/m \leq t} z_{m,i}$ to a proper $Z(t)$, show that

$$dG(s) \mid \text{data} \sim \text{Gamma}(a(s) ds + dZ(s), b(s) + w(s)).$$

- (c) Suppose there are several Poisson event processes being observed, say Z_1, \dots, Z_k , which are conditionally independent given G , and with

$$dZ_j(s) \mid G \sim \text{Pois}(w_j(s) dG(s)) \quad \text{for } j = 1, \dots, k,$$

where the $w_j(s)$ are multiplicative Poisson factor functions. Show that G given data again becomes an Extended Gamma process, with

$$dG(s) \mid \text{data} \sim \text{Gamma}\left(a(s) ds + \sum_{j=1}^k dZ_j(s), b(s) + \sum_{j=1}^k w_j(s)\right).$$

- (d) Assume there are covariates x_1, \dots, x_k at work for the Poisson event counting processes Z_1, \dots, Z_k ; these may also depend on time, say with $x_j(s)$ related to the outcome $dZ_j(s)$. Let $w_j = \exp(x_j^t \beta)$, with a prior $\pi(\beta)$ for this regression parameter. The model at work then says

(i) that β is drawn from the prior $\pi(\beta)$; (ii) that $G(\cdot)$ is an Extended Gamma process, with parameters $(a(s), b(s))$; (iii) that the Poisson processes Z_1, \dots, Z_k have intensity functions $\exp(x_j^t \beta) dG(s)$, hence cumulative intensity functions $\int_0^t \exp(x_j^t \beta) dG(s)$. Show that G , given both the data and β , is another Extended Gamma process, with parameters

$$\left(a(s) + \sum_{j=1}^k dZ_j(s), b(s) + \sum_{j=1}^k \exp(x_j^t \beta) \right).$$

(e) Show that

$$\widehat{G}(t | \beta) = \mathbb{E} \{ G(t) | \text{data}, \beta \} = \int_0^t \frac{a(s) ds + \sum_{j=1}^k dZ_j(s)}{b(s) + \sum_{j=1}^k \exp(x_j^t \beta)},$$

which may also be written out as an integral over intervals with zero jumps plus the component summed over the precise jump times.

(f) Then work out an expression for the posterior density of β . It may be required to set up an MCMC scheme for simulation from this $\pi(\beta | \text{data})$. This leads in particular to the Bayes estimator

$$\widehat{G}(t) = \int \widehat{G}(t | \beta) \pi(\beta | \text{data}) d\beta.$$

(g) [xx just a little more here. point to similar story for Beta process. also link to Nils Beta process with Cox regression type data. and then to Nils-Emil Beta- and Gamma-Process Police Department's Tweetery. xx]

33. The Gamma process, with a Poisson process, with a marks process

[xx something here. actuarial statistics is fond of compound Poisson processes, with total claim size

$$W(t) = \sum_{T_j \leq t} \xi_j,$$

summed over claim times $T_1 < T_2 < \dots$. may here build a Bayesian nonparametrics story, with a prior for the nonhomogeneous Poisson process etc. xx]

34. The biggest jumps of a Beta process

Consider a Beta process A with parameters (c, A_0) . It has infinitely many jumps on each interval. Let $J(t)$ be the biggest jump experienced during the time interval $[0, t]$. To find its distribution, it is convenient to go via the time-discrete version, and then take the limit. Thus let $A_m(t) = \sum_{i/m \leq t} B_{m,i}$, with independent small Beta components

$$B_{m,i} \sim \text{Beta} \left(c \left(\frac{i}{m} \right) a_0 \left(\frac{i}{m} \right) \frac{1}{m}, c \left(\frac{i}{m} \right) - c \left(\frac{i}{m} \right) a_0 \left(\frac{i}{m} \right) \frac{1}{m} \right),$$

as in Exercise 25, and let $J_m(t) = \max\{B_{m,i} : i/m \leq t\}$.

(a) Let $X_\varepsilon \sim \text{Beta}(c\varepsilon, c(1-\varepsilon))$. Show that its cumulative distribution function can be approximated as

$$\begin{aligned} \Pr\{X_\varepsilon \leq v\} &= 1 - \int_v^1 \frac{\Gamma(c)}{\Gamma(c\varepsilon)\Gamma(c(1-\varepsilon))} u^{c\varepsilon-1} (1-u)^{c-c\varepsilon-1} ds \\ &= 1 - \varepsilon F_c(v) \{1 + O(\varepsilon)\}, \end{aligned}$$

with

$$F_c(v) = \Gamma(c) \int_v^1 u^{-1} c(1-u)^{c-1} du.$$

(b) Show from this that

$$\Pr\{J_m(t) \leq v\} = \prod_{i/m \leq t} [1 - \Gamma(c(i/m)) a_0(i/m) (1/m) F_{c(i/m)}(v) \{1 + O(1/m)\}]$$

and that this converges to

$$\Pr\{J(t) \leq v\} = \exp\left\{-\int_0^t \Gamma(c(s)) a_0(s) F_{c(s)}(v) ds\right\}$$

(c) For the case of $c(s) = c$, constant over the time window of interest, show that

$$\Pr\{J(t) \leq v\} = \exp\{-A_0(t) F_c(v)\}.$$

For the special cases of $A \sim \text{Beta}(1, A_0)$, i.e. with concentration function $c(s) = 1$, show that the biggest jump over $[0, t]$ has the simpler distribution

$$\Pr\{J(t) \leq v\} = v^{A_0(t)},$$

which means $J(t) \sim \text{Beta}(A_0(t), 1)$.

35. The jumps of a Beta process via a Poisson process

[xx to be written down cleanly. xx] Consider $A_m(t) = \sum_{i/m \leq t} B_{m,i}$, with independent small Beta contributions

$$B_{m,i} \sim \text{Beta}(c(i/m) a_0(i/m) (1/m), c(i/m) (1 - a_0(i/m) (1/m))),$$

as in Exercise 34. Then study

$$N_m(t) = \sum_{i/m \leq t} I\{B_{m,i} > v\},$$

the number of jumps above level v .

(a) Show that

$$E N_m(t) = \sum_{i/m \leq t} a_0(i/m) (1/m) F_{c(i/m)}(v) \{1 + O(1/m)\} \rightarrow K(t) = \int_0^t a_0(s) F_{c(s)}(v) ds.$$

(b) Show that $N_m(t)$ converges to a Poisson process with cumulative intensity function $K(t)$. In particular, if $c(s) = c$ is constant, the cumulative intensity is $K(t) = A_0(t) F_c(v)$.

(c) Suppose an item is born with a Beta process in its rucksack and that it is alive as long as the jumps do not exceed threshold level v . Derive again the result about its survival distribution from Exercise 34.

(d) Suppose next that this item is alive until the 4th biggest jump exceeds v . Show that the survival function becomes

$$S_4(t) = \Pr\{N(t) \leq 3\} = \exp\{-K(t)\} \{1 + K(t) + \frac{1}{2} K(t)^2 + \frac{1}{6} K(t)^3\}.$$

[xx then something more here. xx]

36. Bernshtein–von Mises theorems

[xx to be written down xx] first for Dirichlet, with fairly clear details. but it takes the Donsker and Kolmogorov thing. then for Beta processes.

37. The Bayesian bootstrap

[xx to be briefly written down. describe Rubin’s (1981) Bayesoian bootstrap. explain that this actually corresponds to sampling from $\text{Dir}(nF_n)$, though this is not in Rubin’s paper. indicate the basics for why this goes well. xx]

38. Hjort’s informative Bayesian bootstrap

[xx to be briefly written down, from Nils Stanford 1985. translate sampling from $F | \text{data} \sim \text{Dir}(aF_0 + nF_n)$ to informative Bayesian bootstrapping. a couple of lemmas spelling out the Bernshtein–von Mises things, with a Brownian bridge. also large-sample equivalent to Efron’s classic nonparametric bootstrap. xx]

39A. The multinormal distribution

‘Multivariate statistics’ is broadly speaking the area of statistical modelling and analysis where data exhibit dependencies. The most important multivariate distribution is the multinormal one. We say that $X = (X_1, \dots, X_k)^t$ is multinormal with mean vector ξ (a k -vector) and variance matrix Σ (a positive definite $k \times k$ matrix) if its density has the form

$$f(x) = (2\pi)^{-k/2} |\Sigma|^{-1/2} \exp\{-\frac{1}{2}(x - \xi)^t \Sigma^{-1} (x - \xi)\} \quad \text{for } x \in \mathcal{R}^k.$$

We write $X \sim N_k(\xi, \Sigma)$ to indicate this. For dimension $k = 1$ this corresponds to the traditional Gaußian $N(\xi, \sigma^2)$.

- (a) Show that if $X \sim N_k(\xi, \Sigma)$ and A is $k \times k$ of full rank, and b a k -vector, then

$$Y = AX + b \sim N_k(A\xi + b, A\Sigma A^t).$$

Generalise to the situation where A is of dimension $m \times k$ (rather than merely $k \times k$).

- (b) Show that if $X \sim N_k(\xi, \Sigma)$, then indeed

$$E X = \xi \quad \text{and} \quad \text{Var } X = \Sigma,$$

justifying the semantic terms used above.

- (c) Show that X is multinormal if and only if all linear combinations are normal. In particular, if $X \sim N_k(\xi, \Sigma)$, then $a^t X = a_1 X_1 + \dots + a_k X_k$ is $N(a^t \xi, a^t \Sigma a)$. – We will also allow saying ‘ $X \sim N_k(\xi, \Sigma)$ ’ in cases where Σ has less than full rank. in particular, a constant may be seen as a normal distribution with zero variance.

- (d) An important property of the multinormal is that a subset of components, conditional on another subset of components, remains multinormal. Show in fact that if

$$X = \begin{pmatrix} X^{(1)} \\ X^{(2)} \end{pmatrix} \sim N_{k_1+k_2} \left(\begin{pmatrix} \xi^{(1)} \\ \xi^{(2)} \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right),$$

then

$$X^{(1)} | \{X^{(2)} = x^{(2)}\} \sim N_{k_1}(\xi^{(1)} + \Sigma_{12} \Sigma_{22}^{-1} (x^{(2)} - \xi^{(2)}), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}).$$

- (e) How tall is Professor Hjort? Assume that the heights of Norwegian men above the age of twenty follows the normal distribution $N(\xi, \sigma^2)$, with $\xi = 180$ cm and $\sigma = 9$ cm. Thus, if you have *not yet seen* or bothered to notice this particular aspect of Professor Hjort and his lectures, your point estimate of his height ought to be $\xi = 180$ and a 95% prediction interval for his height would be $\xi \pm 1.96\sigma$, or $[162.4, 197.6]$. – Assume now that you learn that his four brothers are actually 195 cm, 207 cm, 196 cm, 200 cm tall, and furthermore that correlations between brothers' heights in the population of Norwegian men is equal to $\rho = 0.80$. Use this information about his four brothers (still assuming that you have not noticed Professor Hjort's height) to revise your initial point estimate of Professor Hjort's height. Is he a five-percent statistical outlier in his family (i.e. outside the 95% prediction interval)?
- (f) Assume Professor Hjort has n brothers (rather than merely four) and that you're learning their heights X_1, \dots, X_n . What is the conditional distribution of Professor Hjort's height X_0 , given this information? Represent this as a $N(\xi_n, \sigma_n^2)$ distribution, with proper formulae for its parameters. How small is σ_n for a large number of brothers? (The point here is partly that even if you observe and measure my 99 brothers, there's still a limit to how much you can infer about me.)

39B. Simulating realisations of a Gaussian process

[xx to be written down and polished. xx] We say that $Z = \{Z(x) : x \in [a, b]\}$ is a Gaussian process if all its finite-dimensional distributions are Gaussian. In particular, $Z(x)$ is normal, say $N(m(x), \sigma^2(x))$, and $(Z(x), Z(x'))$ is binormal, with correlation say $\rho(x, x')$.

- (a) Explain why giving the mean function $m(x)$, the standard deviation function $\sigma(x)$, and the correlation function $\rho(x, x')$, is actually sufficient to determine the full distribution of Z .
- For some Gaussian processes there are specialised techniques making it easier-than-brute-force to simulate realisations. In general, however, we can't do much better than brute-force, which means simulating $Z^* = (Z(x_1), \dots, Z(x_n))$, for a fine enough grid x_1, \dots, x_n . The implied distribution is multinormal,

$$Z^* \sim N_n(\xi, \Sigma),$$

with ξ having components $m(x_i)$ and Σ of size $n \times n$ and with components $\sigma(x_i)\sigma(x_j)\rho(x_i, x_j)$. Thus simulating from Z becomes practically the same as being able to simulate from a general multinormal $N_n(0, \Sigma)$.

- (c) The R algorithm `rmvnorm` may be used, for simulating from a given multinormal, but my impression is that it might not work well for higher n . A general technique that can be used here is as follows. First, find a unitary matrix P such that

$$P\Sigma P^t = D = \text{diag}(\lambda_1, \dots, \lambda_n).$$

A unitary or orthonormal matrix Q is one having the property that $QQ^t = I = Q^tQ$. Finding such a P , for given Σ , can be achieved via the `eigen` algorithm in R. Then define, compute, and store the root-matrix

$$\Sigma^{1/2} = PD^{1/2}P^t, \quad \text{with } D^{1/2} = \text{diag}(\lambda_1^{1/2}, \dots, \lambda_n^{1/2}).$$

Verify that $\Sigma^{1/2}\Sigma^{1/2} = \Sigma$. Then use

$$z = \Sigma^{1/2}\varepsilon, \quad \text{where } \varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^t \sim N_n(0, I_n),$$

i.e. these are independent standard normals. Verify that z then has the desired multinormal distribution. Check that the following R code works:

```
squareroot <- function(Sigma)
{rootL <- 0*Sigma
  diag(rootL) <- sqrt(eigen(Sigma, symmetric = T)$values)
  P <- eigen(Sigma, symmetric = T)$vectors
  P %*% rootL %*% t(P)}
```

- (d) Consider an Ornstein–Uhlenbeck process Z on $[0, 10]$, with mean zero and covariance function $\text{cov}\{Z(x), Z(x')\} = \exp(-a|x - x'|)$, say with $a = 1.3579$. Simulate and plot 50 realisations of the Z process.

40. Bayesian Kriging

[xx to be written out and polished. xx] Suppose there is a continuous process $Z(x)$ on $[0, 1]$, which we have observed only in a small number of locations. How can we estimate $Z(x)$ where we have not seen it, along with a measure of precision? This translates to ‘spatial interpolation’ and so on, and with Kriging one of its names (from the Master Thesis of Daniel Gerhardus Krige, 1919–2013, a South African geostatistician).

Suppose $Z(x)$ is Gaussian, with constant mean function a , and covariance function

$$\text{cov}\{Z(x), Z(x')\} = \sigma^2 K_0(|x - x'|),$$

where $K_0(r)$ is the correlation function. This means a stationary setup, where $Z(x)$ and $Z(x + r)$ have a correlation independent of position x .

- (a) Use $a = 1.3579$ and $K_0(r) = \exp(-\lambda r)$, with $\lambda = 2.222$. Simulate realisations of $Z(x)$, for $x \in [0, 1]$. Take $\sigma = 1$ here (but later on we may tinker with this precision parameter).
- (b) Assume now that the scientific team has come back from their expedition and report that for positions 0.11, 0.22, 0.33, 0.77, 0.88, they found that $Z(x)$ is equal to 0.99, 1.33, 1.66, 1.22, 1.11 (yes, I’m inventing this, and will search for a real application later on). Find expressions giving the posterior distribution of $Z = \{Z(x) : x \in [0, 1]\}$.
- (c) Find in particular an expression for $\widehat{Z}(x) = E\{Z(x) \mid \text{data}\}$, and plot that curve. Show in fact, using general formulae from Exercise 39A, that with observed data $z_{\text{data}} = (Z(x_1), \dots, Z(x_n))^t$ in locations x_1, \dots, x_n , we have

$$\widehat{Z}(x) = a + k(x)^t \Sigma^{-1}(z_d - a\mathbf{1}),$$

with Σ the variance matrix of all $K_0(|x_i - x_j|)$, $k(x)$ the vector $(K_0(|x - x_1|), \dots, K_0(|x - x_n|))^t$, and $\mathbf{1}$ the vector $(1, \dots, 1)^t$.

- (d) Find also a formula for $\widehat{\kappa}(x)^2 = \text{Var}\{Z(x) \mid \text{data}\}$, and plot the 90% prediction confidence band $\widehat{Z}(x) \pm 1.645 \widehat{\kappa}(x)$. In fact, show

$$\text{Var}\{Z(x) \mid \text{data}\} = \sigma^2 \{1 - k(x)^\top \Sigma^{-1} k(x)\}.$$

Show that this conditional variance is zero at locations x_1, \dots, x_n where $Z(x)$ is actually observed.

- (e) Try to produce versions of Figure 0.8, where you should also play with different values of the smoothness parameter q , in the correlation function $K_0(r) = \exp(-\lambda r^q)$. The geostatistical default value is apparently $q = 1$, with bigger values leading to more smoothness and smaller values to more roughness. The mathematically allowed span of such smoothness parameter values is $q \in (0, 2)$.
- (f) Above we have reached clear, understandable, implementable, plottable formulae for the mean and standard deviation for each not-observed $Z(x)$ given the data $Z(x_1) = z_1, \dots, Z(x_n) = z_n$ (then to be used for the tiny set of $n = 5$ data locations). To properly understand the full conditional process, and to be able to simulate realisations from this Bayesian nonparametrics posterior distribution, we also need the conditional covariance structure. Show indeed that

$$\begin{aligned} K^*(u_1, u_2) &= \text{cov}\{Z(u_1), Z(u_2) \mid \text{data}\} \\ &= \sigma^2 \{K_0(|u_1 - u_2|) - k(u_1)^\top \Sigma^{-1} k(u_2)\}, \end{aligned}$$

which properly generalises the formula from (d).

- (g) Simulate say 50 realisations from the distribution of $Z = \{Z(x) : x \in [0, 1]\}$ given the observed data, and plot these.

41. Skiing days at ♡ Bjørnholt

[xx to be polished with decent care. xx] Figure 0.9 shows the number of skiing days per season at Bjørnholt, 1897 to 2012. This historical time series process,

$$Z(t) = \text{number of skiing days for winter season } t,$$

has however a gap, with no data for the seasons 1938 to 1954. Use a couple of Bayesian Kriging models and methods to ‘fill in’ the skiing days for these years, with a pointwise 90% band.

- (a) Use first a simple model with $E Z(t) = a$, $\text{Var} Z(t) = \sigma^2$,

$$\text{corr}\{Z(s), Z(t)\} = \exp(-\lambda |s - t|),$$

and with values for a, σ, λ that you pretend are known, and $\lambda = 1.11$.

- (b) Then the same with $E Z(t) = a + b(t - 1900)$ with values for (a, b) that you pretend are known.
- (c) Then a somewhat bigger Bayesian Kriging exercise, where you start with priors for (a, b, σ) , but with $\lambda = 1.11$ still taken known.
- (d) Finally a setup with independent priors for a, b, σ, λ .

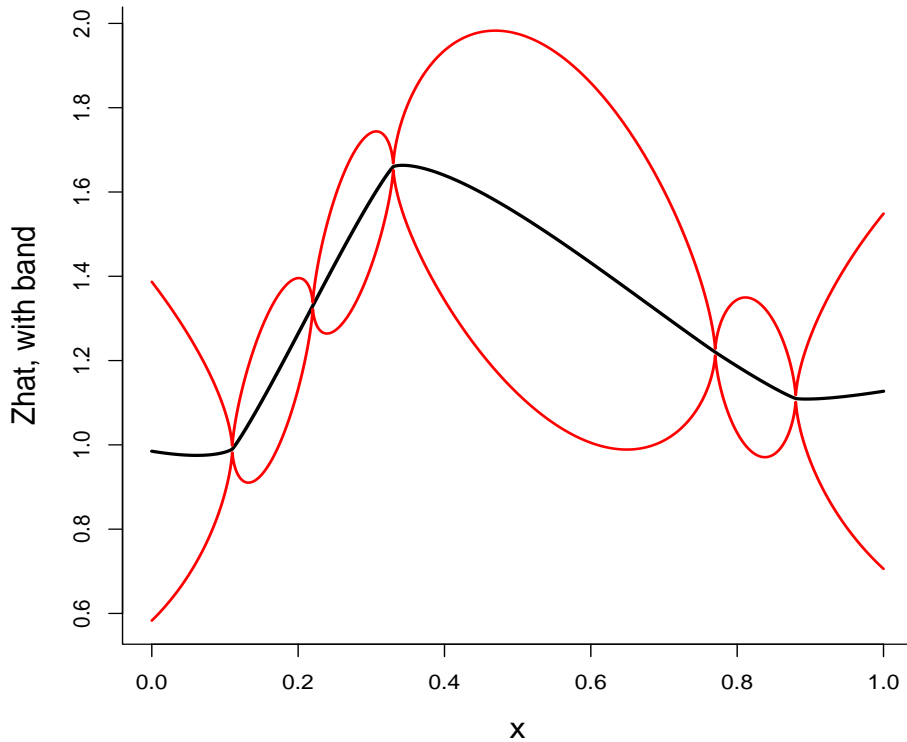


Figure 0.8: Bayesian nonparametric interpolation, after having observed the Gaussian process $Z(x)$ in $n = 5$ locations, as per Exercise 40. The black curve is the posterior mean and the red curves give a pointwise 90% credibility band. Here I've used the correlation function $K_0(r) = \exp(-\lambda r^q)$, with $q = 1.25$, and otherwise taken $\sigma = 0.50$, $a = 1.3579$, $\lambda = 2.222$.

(e) Then extrapolation, from 2013 go 2028.

42. Bayesian nonparametric regression

[xx to be written out and polished. xx] model is

$$y_i = m(x_i) + \varepsilon_i \quad \text{for } i = 1, \dots, n,$$

where the ε_i are i.i.d. from $N(0, \sigma^2)$. Suppose $m(x)$ is Gaussian, with mean function $m_0(x)$ and covariance function for the form $\sigma_0^2 K_0(|x - x'|)$, with a given correlation function $K_0(r)$.

Then find expressions for the conditional mean, the conditional variance, and conditional covariance, of the process $m(x)$, given the data (x_i, y_i) .

43. Bayes and minimax

[xx to be polished. xx] Consider the general framework with data y from model $f(y, \theta)$, with loss function $L(\theta, a)$ associated with action a . An action algorithm $\hat{a} = \hat{a}(y)$ then has risk function

$$R(\hat{a}, \theta) = E_\theta L(\theta, \hat{a}(Y)) = \int L(\theta, \hat{a}(y)) f(y, \theta) dy.$$

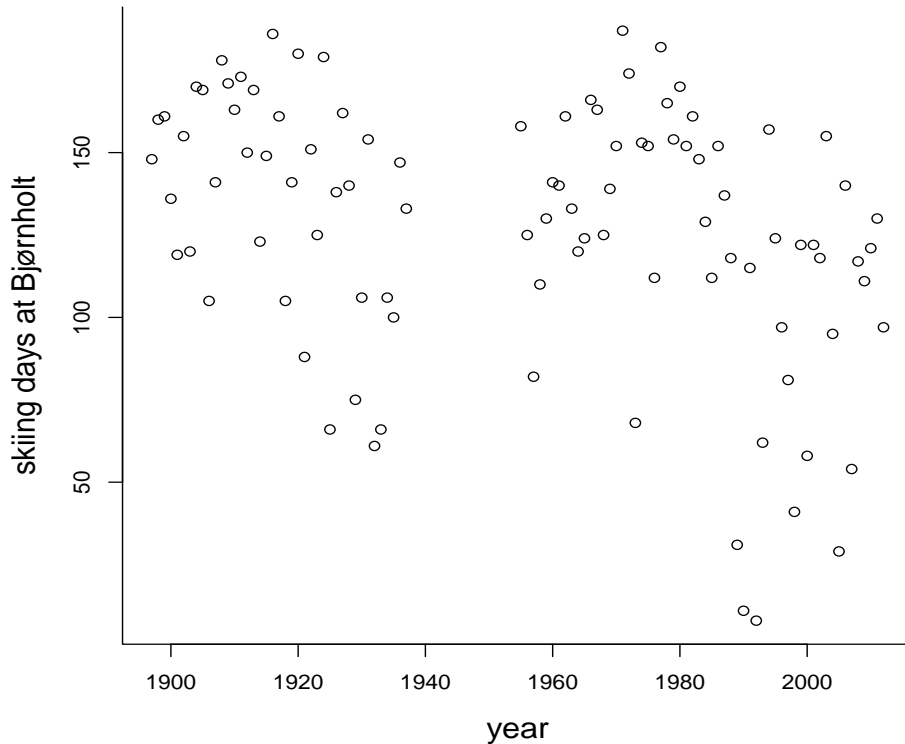


Figure 0.9: The wondrously falling snow flakes [you may sing Sibelius' *Demanten på marssnön* here] represent skiing days per winter season at Bjørnholt, from 1897 to 2012, but with no data for the seasons 1938 to 1954.

Its maximum risk value is

$$R_{\max}(\hat{a}) = \sup R(\hat{a}, \theta).$$

For an arbitrary prior $\pi(\theta)$, the Bayes solution is \hat{a}_B , minimising the posterior expected loss

$$E\{L(\theta, a) | y\} = \int L(\theta, a)\pi(\theta | y) d\theta.$$

It minimises the Bayes risk $BR(\hat{a}, \pi)$ over all procedures \hat{a} , and achieves the minimum Bayes risk $MBR(\pi)$; cf. the general setup described in Exercise 2.

(a) For any prior π and action plan \hat{a} , show that

$$MBR(\pi) \leq R_{\max}(\hat{a}).$$

(b) Suppose the action method $a^* = a^*(y)$ is such that there is a prior π for which $MBR(\pi) = R_{\max}(a^*)$. Show that a^* is then a minimax strategy, i.e. $R_{\max}(a^*) \leq R_{\max}(\hat{a})$ for any competitor \hat{a} .

(c) Suppose somewhat more generally that a^* is such that there is a sequence of priors π_j with $MBR(\pi_j) \rightarrow R_{\max}(a^*)$ (with the convergence in question taking place upwards, not downwards, in view of point (a)). Show that a^* is minimax.

44. Some minimax estimators

Here we through a little list of situations where minimax estimators can be identified.

- (a) Suppose $y|\theta$ is a $N(\theta, 1)$, with θ to be estimated with squared error loss. Show that $\theta^* = y$ itself, the maximum likelihood estimator, has risk function equal to the constant 1. With a prior $N(0, \tau^2)$ for θ , find the posterior distribution, the Bayes estimator $\hat{\theta}_\tau$, its risk function, and finally the minimum Bayes risk $\text{MBR} = \tau^2/(\tau^2 + 1)$. Hence show that y is minimax.
- (b) Suppose now that $y|\theta$ is a $N_p(\theta, \Sigma)$, with a known variance matrix Σ , and with θ to be estimated with the loss function $(\hat{\theta} - \theta)^t \Sigma^{-1} (\hat{\theta} - \theta)$. Show that the ML estimator is y itself; that its risk function is the constant p .
- (c) Show that $\hat{\theta} = y$ is minimax, by working with the prior that takes $\theta \sim N_p(0, \tau^2 \Sigma)$.
- (d) Generalise the above to the case where y_1, \dots, y_n are i.i.d. from $N_p(\theta, \Sigma)$, and still with loss function as in (b). Show that the sample average $\bar{y} = (1/n) \sum_{i=1}^n y_i$ is minimax for θ .
- (e) Consider a nonparametric regression setup with $y_i = m(x_i) + \varepsilon_i$ for $i = 1, \dots, n$, where the ε_i are i.i.d. error terms from the $N(0, \sigma^2)$. The task is estimation of m at positions x_1, \dots, x_n , with loss function $\sum_{i=1}^n \{\hat{m}(x_i) - m(x_i)\}^2$. Show that the very simple (and not particularly clever) estimator $m^*(x_i) = y_i$ is minimax. It is not admissible, however, cf. the big literature on Stein estimation. There are competing estimators, which are also minimax and with the same maximum risk $p\sigma^2$, but with considerably lower risk in important parts of the parameter space.

45. Minimax estimators for multiple Poisson means

Here we work with minimax estimators for situations involving multiple Poisson parameters.

- (a) Let $y|\theta$ come from the $\text{Pois}(\theta)$ distribution, with θ to be estimated via the weighted quadratic loss function

$$L(\theta, \hat{\theta}) = \frac{(\hat{\theta} - \theta)^2}{\theta}.$$

Show that the ML estimator is y and that its risk function is the constant 1. With the $\text{Gamma}(a, b)$ prior for θ , find the Bayes estimator, its risk function, and the minimum Bayes risk $\text{MBR}(a, b)$. When is this maximal? Show that $\theta^* = y$ indeed is minimax.

- (b) Now consider an infinite (or finite) sequence of independent Poisson variables, say $y_i|\theta_i \sim \text{Pois}(\theta_i)$, with the long vector of means to be estimated with the loss function

$$L(\theta, \hat{\theta}) = \sum_{i=1}^{\infty} w_i \frac{(\hat{\theta}_i - \theta_i)^2}{\theta_i},$$

where the weights w_i are given numbers with a convergent series. Show that the estimator which uses $\theta_i^* = y_i$ for all i has a constant risk function, and that it is minimax. – It is however not admissible; see Stoltenberg and Hjort (2019a).

46. Minimax estimators for binomial parameters

Here we work with minimax estimators for situations involving the binomial distribution.

- (a) Let $y|\theta$ be a simple binomial (n, θ) , with θ to be estimated with squared error loss. For the ML estimator y/n , find the risk function and its maximum value. Then let θ have the Beta($a\theta_0, a(1 - \theta_0)$) prior. Show that the Bayes estimator becomes

$$\hat{\theta} = \frac{a\theta_0 + y}{a + n} = c_n\theta_0 + (1 - c_n)y/n,$$

with $c_n = a/(a + n)$. Find its risk function and the associated minimum Bayes risk

$$\text{MBR}(a, \theta_0) = \theta_0(1 - \theta_0) \frac{a}{(a + 1)(a + n)}.$$

When is this maximised, giving the Nature's maximin prior?

- (b) Show that the estimator

$$\theta^* = \frac{\sqrt{n}}{\sqrt{n} + 1} \frac{y}{n} + \frac{1}{\sqrt{n} + 1} \frac{1}{2}$$

has constant risk function, and that it is minimax. Draw the risk functions for θ^* and the ML estimator y/n in the same diagram.

- (c) [xx To generalise from two to more boxes, consider (x, y, z) having the trinomial (n, p, q, r) distribution, with $p + q + r = 1$ and hence also $x + y + z = n$. Suppose the problem is joint estimation of (p, q) , with loss function

$$L((p, q), (\hat{p}, \hat{q})) = (\hat{p} - p)^2 + (\hat{q} - q)^2.$$

Find the risk function, along with its maximum, for the ML estimators $(x/n, y/n)$.

- (d) Work with the Dirichlet prior (ap_0, aq_0, ar_0) for (p, q, r) . Find the Bayes estimator and its risk function. [xx then attempt to find a minimax estimator. xx]
- (e) A sometimes more natural loss function in this trinomial situation is that of

$$L((p, q, r), (\hat{p}, \hat{q}, \hat{r})) = \frac{(\hat{p} - p)^2}{p} + \frac{(\hat{q} - q)^2}{q} + \frac{(\hat{r} - r)^2}{r}.$$

Find the risk function for the ML estimator. Show that it is actually minimax, under this loss function.

47. A nonparametric minimax estimator for an unknown mean

[xx to be polished. xx] Suppose i.i.d. observations X_1, \dots, X_n are available from an unknown distribution P on the unit interval $[0, 1]$. We only know that $P \in \mathcal{M}$, the set of all distribution functions on $[0, 1]$. We wish to estimate the mean $\theta = \int x dP(x)$, with quadratic loss function $(\hat{\theta} - \theta)^2$. Below I exhibit a minimax estimator (mm) for this problem. Lehmann (1951, Mimeographed Lecture Notes on the Theory of Point Estimation from Berkeley) did this, with similar arguments.

Lehmann also claimed in these Lecture Notes that the estimator given in (mm) is admissible – but his argument was not correct, as it turns out. Nils Lid Hjort's perhaps First Theorem was to prove that the (mm) estimator *is* nonparametrically admissible (in an exam project for Erik N. Torgersen on decision theory, 1975, which consisted in reading, comprehending, and presenting the Ferguson 1973 paper for the exam marker).

- (a) Work out the risk function for the direct sample average \bar{X} :

$$R(\bar{X}, P) = (1/n)\sigma(P)^2, \quad \text{with} \quad \sigma(P)^2 = \int \{x - \theta(P)\}^2 dP(x).$$

- (b) Show that the variance $\sigma(P)^2$ is maximal, over all distributions on $[0, 1]$, when P is concentrated in the end-points 0 and 1, with equal probabilities $\frac{1}{2}, \frac{1}{2}$. Hence

$$\max\{R(\bar{X}, P): P \in \mathcal{M}\} = (1/n)(1/4).$$

- (c) Then consider the cool enough estimator

$$\hat{\theta} = \frac{1}{\sqrt{n}+1} \frac{1}{2} + \frac{\sqrt{n}}{\sqrt{n}+1} \bar{X}. \quad (\text{mm})$$

Show that its risk function can be written

$$\begin{aligned} R(\hat{\theta}, P) &= \left(\frac{\sqrt{n}}{\sqrt{n}+1}\right)^2 \frac{\sigma(P)^2}{n} + \left\{ \frac{\sqrt{n}}{\sqrt{n}+1} \theta(P) + \frac{1}{\sqrt{n}+1} \frac{1}{2} - \theta(P) \right\}^2 \\ &= \frac{1}{(\sqrt{n}+1)^2} [\sigma(P)^2 + \left\{ \frac{1}{2} - \theta(P) \right\}^2]. \end{aligned}$$

- (d) Show that the max risk for the (mm) estimator is

$$\max\{R(\hat{\theta}, P): P \in \mathcal{P}\} = \frac{1/4}{(\sqrt{n}+1)^2}.$$

- (e) Start with the prior $P \sim \text{Dir}(aP_0)$ for P , with P_0 a given distribution on the unit interval, with mean $\theta_0 = \theta(P_0)$ and standard deviation $\sigma_0 = \sigma(P_0)$. Show that the Bayes estimator becomes

$$\hat{\theta} = c_n \theta_0 + (1 - c_n) \bar{x}, \quad \text{with} \quad c_n = a/(a + n).$$

- (f) Show that the risk function for this Bayes estimator becomes

$$R(\hat{\theta}, P) = c_n^2 \{\theta(P) - \theta_0\}^2 + (1 - c_n)^2 \sigma^2(P)/n.$$

- (g) Then work out the corresponding minimum Bayes risk, for the $\text{Dir}(aP_0)$ prior,

$$\text{MBR}(aP_0) = \int R(\hat{\theta}, P) d\mathcal{P}(P) = \sigma_0^2 \frac{a}{(a+1)(a+n)}.$$

Show that this is maximised when $a = \sqrt{n}$ and P_0 has equal mass $\frac{1}{2}$ and $\frac{1}{2}$ at the endpoints 0 and 1.

- (h) Then show that it is minimax (Lehmann 1951, Berkeley Notes, precursor to the Theory of Point Estimation 1983 book). [xx fill in, not too hard. xx]

- (i) Then show that it is actually also admissible; Lehmann made an error here, in these 1951 Berkeley Notes, but Nils 1976 has several proofs. [xx i fill in one of these, perhaps as a separate exercise; this is considerably harder than proving minimaxity. xx]

48. A nonparametric minimax estimator for a distribution function

Let X_1, \dots, X_n be i.i.d. from some unknown distribution function F on the real line. The standard empirical distribution function is $F_n(t) = (1/n) \sum_{i=1}^n I\{X_i \leq t\}$, corresponding to probability weight $1/n$ in each observation, and can be computed using the `ecdf` in R. Here I exhibit a different estimator, namely

$$F^*(t) = \frac{\sqrt{n}}{\sqrt{n}+1} F_n(t) + \frac{1}{\sqrt{n}+1} \frac{1}{2}.$$

We shall see that it is minimax, as proven in Hjort (1976), with the loss function

$$L(F, \hat{F}) = \int (\hat{F} - F)^2 dW = \int \{\hat{F}(t) - F(t)\}^2 dW(t),$$

with W some weight measure with finite mass.

- (a) Let F have a $\text{Dir}(aF_0)$ prior. Show that the Bayes estimator becomes

$$\hat{F}(t) = \frac{a}{a+n} F_0(t) + \frac{n}{a+n} F_n(t) = c_n F_0(t) + (1 - c_n) F_n(t).$$

- (b) Find its risk function and its max-risk. Find also the max-risk for F_n .
 (c) Work out that the minimum Bayes risk for the Dirichlet process prior can be expressed as

$$\text{MBR}(\text{Dir}(aF_0)) = \int F_0(1 - F_0) dW \frac{a}{(a+1)(a+n)}.$$

When is this as its biggest, i.e. Nature creating the worst possible situation for the statistician?

- (d) Show that F^* indeed is minimax.
 (e) I recall attempting to prove, back in 1976, that F^* is also admissible, without fully succeeding – so that particular problem is up for grabs, I think.

49. Survival models via Gamma process boundary hitting

Part of the point of this exercise is to show that the tools and machines of Bayesian nonparametrics may be utilised for different purposes, including model building. Let Z be a Gamma process with $Z(t) \sim \text{Gamma}(aM(t), 1)$, and suppose an event takes place as soon as Z crosses the threshold c . The distribution of

$$T = \min\{t: Z(t) \geq c\}$$

can then be worked with, for many purposes.

- (a) Show that its survival distribution becomes

$$S(t) = \Pr\{Z(t) < c\} = G(c, aM(t), 1) = \frac{1}{\Gamma(aM(t))} \int_0^c x^{aM(t)-1} \exp(-x) dx.$$

- (b) Program such curves $S(t)$ for a few choices of a and c , with $M(t) = t$, with $M(t) = \log(1+t)$, and with $M(t) = t^{1/2}$. For each case, compute and display also the hazard rate $h(t) = f(t)/S(t)$, where the density f can be computed numerically via

$$f(t) = -\{G(c, aM(t+\varepsilon), 1) - G(c, aM(t-\varepsilon), 1)\}/(2\varepsilon)$$

for a small ε .

- (c) For the Egyptian life times data set, fit the model which takes $M(t) = \exp(\kappa t) - 1$, making it into a three-parametric model, with parameters (a, κ, c) . Do this via the log-likelihood function

$$\ell_n(a, \kappa, c) = \sum_{i=1}^n \log f(t_i, a, \kappa, c).$$

For numerical simplicity, divide lifelengths with 100, so that these are recorded on a scale from 0 to 1 (actually from 0.5/100 to 96/100).

- (d) Then distinguish the men and the women in that dataset, and fit the model which takes the same a and the same $M(t) = \exp(\kappa t) - 1$, but two different thresholds c_w and c_m . Compute the $\ell_{n,\max}$ value and the AIC, and compare with the best models for these data, used in Claeskens and Hjort (2008, Ch. 2).
- (e) Then try the Gamma process boundary hitting time model which takes $aM(t)$ for the men, with $M(t) = \exp(\kappa t) - 1$, and $aM_w(t)$ for the women, with

$$M_w(t) = M(t) + dW(t),$$

with $W(t)$ the cdf for the uniform distribution of $[15/100, 40/100]$, i.e. for the age window $[15, 40]$. Estimate a, κ, d, c , and assess the uncertainty of the the estimates. Compute ℓ_{\max} and show that the AIC value $2\ell_{\max} - 2 \dim$ is very high, with \dim the length of the parameter vector. Give an interpretation of the fitted model and spend a minute speculating about lives lived two thousand years ago. [xx nils put in the figure here. xx]

50. Survival models via biggest Gamma process jumps

In Exercise XX we build survival models from threshold hitting times for Gamma processes. Another version is via the sizes of the biggest jumps. Let $Z(t) \sim \text{Gamma}(aM(t), 1)$, and suppose an event takes place when a jump exceeds a threshold d .

- (a) Show that the survival time in question can be characterised as $T = \min\{t: J(t) > d\}$, where $J(t)$ is the biggest jump experienced over the time window $[0, t]$. Hence show that the survival curve becomes

$$S(t) = \Pr\{J(t) < d\} = \exp\{-aM(t) E_1(v)\} \quad \text{for } t > 0,$$

using results of Exercise 24.

- (b) In a Cox type regression situation, with survival data (t_i, δ_i, x_i) for individuals $i = 1, \dots, n$, suppose individual i has a distribution corresponding to the biggest jump size crossing level v_i . Show that this means hazard rates

$$h_i(s) = am(s)E_1(v_i) \quad \text{for } i = 1, \dots, n.$$

With model parametrisation $E_1(v_i) = \exp(x_i^\dagger \beta)$, we have reinvented the Cox regression model, complete with proportional hazards.

- (c) [xx a bit more. frailty. fixed frailty and frailty along the way. competing risks. additive risks and the Aalen model. xx]

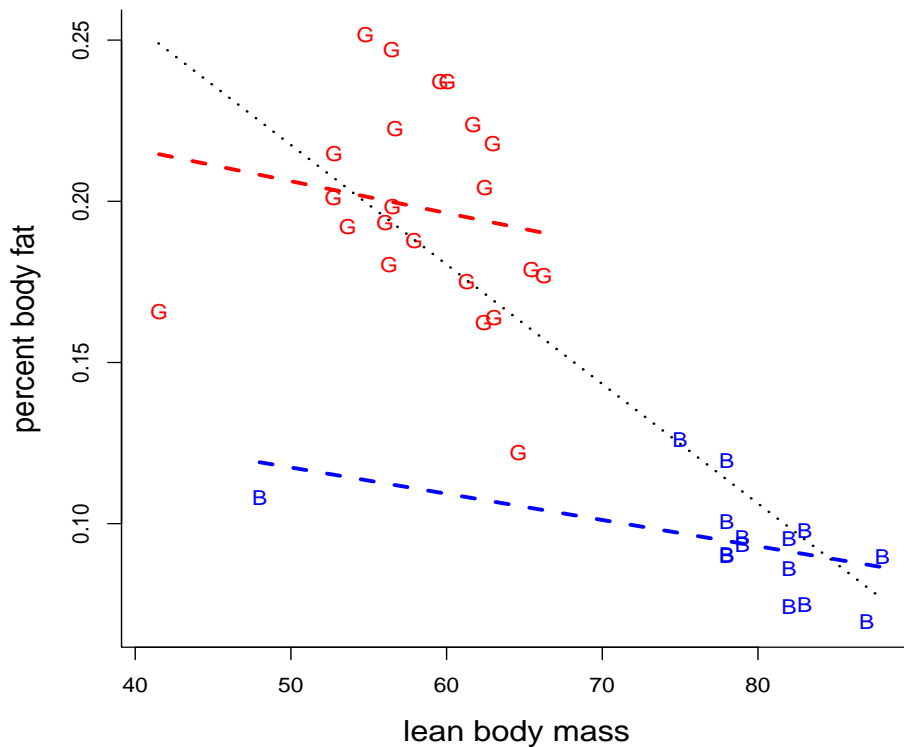


Figure 0.10: For 37 Australian rowers (22 girls and 15 boys), the figure displays x , lean body mass, and y , percent body fat (converted to the unit scale). The three simple and direct regression lines are for boys and for girls separately, and for the overall data ignoring gender information.

51. Density estimation via log-linear expansions

For 37 Australian rowers (the first 22 girls, the next 15 boys), these are the lean body mass x (in kg) and the percent body fat y (converted to unit interval scale), respectively:

66.24 57.92 56.52 54.78 56.31 62.96 56.68 62.39 63.05 56.05 53.65 65.45
 64.62 60.05 56.48 41.54 52.78 52.72 61.29 59.59 61.70 62.46 78.00 75.00
 78.00 87.00 78.00 79.00 79.00 48.00 82.00 82.00 82.00 83.00 88.00 83.00 78.00

and

0.177 0.188 0.198 0.252 0.180 0.218 0.222 0.162 0.164 0.194 0.192 0.179
 0.122 0.237 0.247 0.166 0.215 0.201 0.175 0.237 0.224 0.204 0.090 0.126
 0.090 0.070 0.100 0.096 0.094 0.108 0.086 0.095 0.074 0.098 0.090 0.075 0.120

The data, taken from Australian Institute of Sport, via the wondrous dataset collection `ozdas1`, are plotted in Figure 0.10. In this exercise we shall use the data to estimate the density of x , so far ignoring the gender information, partly since the bimodality offers a little challenge; also, the two very slim rowers, perhaps cox-swains, contribute to a mild statistical confusion. Given the density estimation method one may of course also estimate the x density for boys and for girls separately, and similarly for y .

For convenience we divide x by 100, with the resulting 37 data points x_1, \dots, x_n to be seen as i.i.d. from a density $f(x)$ on the unit interval. To estimate the density we shall use the log-linear expansion

$$f(x, \theta) = \exp\left\{\sum_{j=1}^m \theta_j T_j(x) - c(\theta_1, \dots, \theta_m)\right\},$$

with basis functions

$$T_j(x) = \sqrt{2} \cos(j\pi x) \quad \text{for } x \in [0, 1],$$

and with

$$c(\theta_1, \dots, \theta_m) = c(\theta) = \log\left(\int_0^1 \exp\left\{\sum_{j=1}^m \theta_j T_j(x)\right\} dx\right).$$

- (a) Show that $f(x, \theta)$ indeed defines a density function on the unit interval. Show also that the basis functions are orthonormal, in the sense that $\int_0^1 T_j^2 dx = 1$ and $\int_0^1 T_j T_k dx = 0$ for $j \neq k$. This helps both numerical accuracy and interpretation, though the model above, and the methods to follow, make sense and can be made to work also for general sequences of basic functions, like the polynomial $T_j(x) = (x - \frac{1}{2})^j$.
- (b) For θ small, where the model can be seen as leading to small perturbations of the uniform density, show that

$$c(\theta_1, \dots, \theta_m) = \frac{1}{2} \sum_{j=1}^m \theta_j^2 + O\left(\sum_{j=1}^m |\theta_j|^3\right).$$

- (c) Show that the log-likelihood function for the data becomes

$$\ell_n(\theta) = n \sum_{j=1}^m \theta_j \bar{T}_j - nc(\theta_1, \dots, \theta_m),$$

with $\bar{T}_j = (1/n) \sum_{i=1}^n T_j(x_i)$. Find the ML estimator $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_m)$, and plot the associated ML estimated density, with length-of-expansion m set equal to $m = 2, 3, 4, 5, 6$. You may also compute AIC values $2\ell_{n, \max} - 2m$ to check which model order it prefers. You need to compute the $c(\theta)$ function via numerical integration, e.g. using `integrate` in R.

- (d) Then to a Bayesian handling of such a model. Place independent normal priors $\theta_j \sim N(0, \tau^2/j^2)$ on the parameters. Here τ is a fine-tuning parameter; in particular, if τ is small, the random density is close to a uniform with high probability. Show that

$$Z(x) = \sum_{j=1}^{\infty} \theta_j T_j(x)$$

is a well-defined normal process, even with infinitely many basis functions. Find its mean, variance, and covariance function.

- (e) Take initially $m = 6$ (a rather modest approximation to infinity). Set up an MCMC to sample from the posterior distribution of $(\theta_1, \dots, \theta_m)$. Use this to compute and display the Bayes estimator

$$\hat{f}(x) = E\{f(x, \theta) \mid \text{data}\},$$

along with a pointwise 90% credibility interval. Repeat this exercise with a higher number for m .

- (f) Now that you have a Bayesian nonparametric density estimation scheme, apply it to the boys and the girls separately, and comment. Also use it for estimating the density of y , the percent body fat.

52. More on density estimation

The method of Exercise 50 starts uses a perturbation model around the uniform density, and can sometimes perform well, depending on both the underlying correct density and the sample size. It often pays off to ‘give it a good start’, however, for example as follows. With $f_0(x)$ such a start estimate, perhaps the prior mean, work with the parametric class

$$f(x, \theta) = \frac{f_0(x) \exp\{\sum_{j=1}^m \theta_j T_j(x)\}}{\int_0^1 f_0(x') \exp\{\sum_{j=1}^m \theta_j T_j(x')\} dx'} = f_0(x) \exp\left\{\sum_{j=1}^m \theta_j T_j(x) - c_m(\theta)\right\},$$

where

$$c_m(\theta) = \log\left(\int_0^1 f_0(x) \exp\left\{\sum_{j=1}^m \theta_j T_j(x)\right\} dx\right).$$

Simulate a dataset of size n from a Beta distribution with parameters $(4, 1)$, say, and select a $f_0(x)$ which is not so far away from this, say the Beta with parameters $(5, 2)$. Use independent priors $\theta_j \sim N(0, \tau^2/j^2)$ and m equal to perhaps 10 in this construction. Simulate 10^4 realisations from the posterior distribution of $(\theta_1, \dots, \theta_m)$, via MCMC, and plot the 0.05, 0.50, 0.95 pointwise quantiles of $f(x, \theta)$ given data.

53. Regression via linear expansions

[xx to be written and polished. xx] Consider the pairs (x_i, y_i) of lean body mass (divided by 100, ranging from 0.415 to 0.880) and percent body fat for the $n = 37$ Australian rowers above. We are to carry out nonparametric regression for $m(x) = E(Y | x)$, so far ignoring genders. Consider the model

$$y_i = m(x_i) + \varepsilon_i \quad \text{for } i = 1, \dots, n,$$

with the ε_i i.i.d. with mean zero and standard deviation σ , and where

$$m(x) = \beta_0 + \sum_{j=1}^m \beta_j (x - \frac{1}{2})^j.$$

For simplicity in this exercise take $\sigma = 0.040$ known, and use independent priors $\beta_j \sim N(0, \tau^2/j^2)$ for $j = 1, 2, 3, \dots$ but $\beta_0 \sim N(0.50, 1)$. The τ works as a tuning parameter, and you may start by setting it equal to 1.3579.

Now work out the posterior distribution for $\beta = (\beta_0, \beta_1, \dots, \beta_m)$. This can be done directly, via Gaussian calculations, or via MCMC. With $m = 10$, simulate 10^4 realisations for β given data. Compute and display the Bayes estimated curve

$$\widehat{m}(x) = E\left\{\beta_0 + \sum_{j=1}^m \beta_j (x - \frac{1}{2})^j \mid \text{data}\right\} = \widehat{\beta}_0 + \sum_{j=1}^m \widehat{\beta}_j (x - \frac{1}{2})^j.$$

Compare with least squares analysis and AIC.

XX. The Beta- and Gamma-process Police Department Tweetery

A Nils-Emil story will be told here, once we've understood things well enough. Data are non-homogeneous Poisson counts $Z_i(t)$ for weekends $i = 1, \dots, n$, with time t running through $[0, 60]$ hours. There's a covariate vector x_i available for weekend i , possibly also time-dependent, with $x_i(s)$ previsible at time s . One of the models we're aiming for takes

- (i) β from $\pi(\beta)$;
- (ii) G is an extended Gamma $(dA_0(s), b(s))$, where $dA_0(s) = a_0(s) ds$, i.e. $A_0(t)$ is the integral $\int_0^t a_0(s) ds$;
- (iii) event processes Z_1, \dots, Z_n are independent and Poisson, with

$$dZ_i(s) \sim \text{Pois}(\exp(x_i^t \beta) dG(s)) \quad \text{for } i = 1, \dots, n.$$

We take $a_0(s)$ and $b(s)$ known, and wish to carry Bayesian inference for (β, G) . The illustration story will be the Beta- and Gamma-process Police Department's tweets over a sequence of Oslo weekends, 2018.

References

- Billingsley, P. (1968). *Convergence of Probability Measures*. Wiley, New York.
- Clauset, A. (2017). The enduring threat of a large interstate war. Technical Report, One Earth Foundation.
- Clauset, A. (2018). Trends and fluctuations in the severity of interstate wars. *Science Advances* **4**, xx-xx.
- Cunen, C., Hjort, N.L., and Nygård, H. (2019). Statistical sightings of better angels. *Journal of Peace Research* [to appear].
- Ferguson, T.S. (1973). A Bayesian analysis of some nonparametric problems. *Annals of Statistics* **1**, 209–230.
- Ferguson, T.S. (1974). Prior distributions on spaces of probability measures. *Annals of Statistics* **2**, 615–629.
- Ferguson, T.S. and Klass, M.J. (1972). A representation of independent increment processes without Gaussian components. *Annals of Mathematical Statistics* **43**, 1634–1643.
- Ghosal, S. and van der Vaart, A. (2017). *Fundamentals of Nonparametric Bayesian Inference*. Cambridge University Press, Cambridge.
- Ghosh, A. (2017). Robust inference under the Beta regression model with application to health care studies. [arXiv](#)
- Heger, A. (2007). Jeg og jordkloden. Dagsavisen.
- Hjort, N.L. (1976). *The Dirichlet Process Applied to Some Nonparametric Problems*. Cand. real. thesis [in Norwegian], Department of Mathematics, Nordlysobservatoriet, University of Tromsø.
- Hjort, N.L. (1985). Discussion contribution to P.K. Andersen and Ø. Borgan's 'Counting process models for life history data: A review'. *Scandinavian Journal of Statistics* **12**, xx–xx.
- Hjort, N.L. (1985). An informative Bayesian bootstrap. Technical Report, Department of Statistics, Stanford University.
- Hjort, N.L. (1986). Discussion contribution to P. Diaconis and D. Freedman's paper 'On the consistency of Bayes estimators'. *Annals of Statistics* **14**, 49–55.

- Hjort, N.L. (1990). Nonparametric Bayes estimators based on Beta processes in models for life history data. *Annals of Statistics* **18**, 1259–1294.
- Hjort, N.L. (1991). Bayesian and empirical Bayesian bootstrapping. Statistical Research Report, Department of Mathematics, University of Oslo.
- Hjort, N.L. (2003). Topics in nonparametric Bayesian statistics [with discussion]. In *Highly Structured Stochastic Systems* (eds. P.J. Green, N.L. Hjort, S. Richardson). Oxford University Press, Oxford.
- Hjort, N.L. (2018). Towards a More Peaceful World [Insert ‘!’ or ‘?’ Here]. FocuStat Blog Post.
- Hjort, N.L. (2010). An invitation to Bayesian nonparametrics. In *Bayesian Nonparametrics* (by Hjort, N.L., Holmes, C.C., Müller, P., and Walker, S.G.), 1–21.
- Hjort, N.L., Holmes, C.C., Müller, P., and Walker, S.G. (2010). *Bayesian Nonparametrics*. Cambridge University Press, Cambridge.
- Hjort, N.L. and Kim, Y. (2013). Beta processes and their applications and extensions. Statistical Research Report, Department of Mathematics, University of Oslo.
- Hjort, N.L. and Ongaro, A. (2005). Exact inference for random Dirichlet means. *Statistical Inference for Stochastic Processes* **8**, 227–254.
- Hjort, N.L. and Ongaro, A. (2006). On the distribution of random Dirichlet jumps. *Metron* **LXIV**, 61–92.
- Hjort, N.L. and Petrone, S. Nonparametric quantile inference using Dirichlet processes. In *Festschrift for Kjell Doksum* (ed. V. Nair).
- Hjort, N.L. and Walker, S.G. (2009). Quantile pyramids for Bayesian nonparametrics. *Annals of Statistics* **37**, 105–131.
- Lehmann, E.L. (1951). Notes on the Theory of Point Estimation. (Mimeographed by C. Blyth.) Department of Statistics, University of Berkeley, California.
- Müller, O., Quintana, F.A., Jara, A., and Hanson, T. (2015). *Bayesian Nonparametric Data Analysis*. Springer-Verlag, Berlin.
- Rubin, D. (1981). The Bayesian bootstrap. *Annals of Statistics* **9**, 130–134.
- Sethuraman, J. (1994). A constructive definition of Dirichlet priors. *Statistica Sinica* **4**, 639–650.
- Sethuraman, J. and Tiwari, R. (1982). Convergence of Dirichlet measures and the interpretation of their parameter. In: *Proceedings of the Third Purdue Symposium on Statistical Decision Theory and Related Topics* (eds. S.S. Gupta and J. Berger), 305–315. Academic Press, New York.
- Stoltenberg, E.Aa. and Hjort, N.L. (2019a). Simultaneous estimation of Poisson parameters. *Journal of Multivariate Analysis*, in its way.
- Stoltenberg, E.Aa. and Hjort, N.L. (2019b). Modelling and analysing the Beta- and Gamma Police Tweetery data. [Manuscript, in progress.]
- Wolpert, R.L. and Ickstadt, K. (1998). Poisson/gamma random field models for spatial statistics. *Biometrika* **85**, 251–267.