



Green, P. J. (2015). Discussion of paper by Cowell, Graversen, Lauritzen and Mortera. *Journal of the Royal Statistical Society: Series C*, 64(1), 41-41.

Peer reviewed version

[Link to publication record in Explore Bristol Research](#)
PDF-document

University of Bristol - Explore Bristol Research

General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available:
<http://www.bristol.ac.uk/pure/about/ebr-terms.html>

**Discussion of ‘Analysis of forensic DNA mixtures with artefacts’
by Cowell, Graversen, Lauritzen and Mortera**

Peter J Green,

**University of Bristol, UK, and University of Technology, Sydney,
Australia; P.J.Green@bristol.ac.uk**

For several reasons, I regret not being able to come to the meeting, including that I understand that there is some connection between what I write here and the discussion by Dr Torben Tvedebrink.

Since this paper (‘CGLM’) was completed, I have with Julia Mortera been exploring the effects of uncertainty in the allele frequencies $q = (q_a)_{a=1}^A$. In earlier work (Green and Mortera, 2009) addressing cases where the DNA traces are of discrete allele-presence indicators rather than continuous peak heights, such questions were explored under an (idealised) Dirichlet model – this leads to a Pólya urn scheme readily implementable in a Bayes net formulation for the inference. More precisely, $q|\rho \sim \text{Dirichlet}((M\rho_a)_{a=1}^A)$, where q are the true, unknown, allele frequencies, $\rho = (\rho_a)_{a=1}^A$ the database frequencies and M the database size; this is typically only a few hundred in practice, so there is considerable uncertainty. We write $\alpha_a = M\rho_a$.

Combining this Dirichlet prior on q with the CGLM set-up, Dirichlet–Multinomial conjugacy then gives the joint distribution for the allele counts n_{ia} , recognising this uncertainty. Recall that n_{ia} is the number of a alleles for the i th individual, $a = 1, 2, \dots, A$, $i = 1, 2, \dots, I$. Conditional on allele frequencies $\{q_a\}$, the vectors $n_{i\cdot} = (n_{ia})_{a=1}^A$ are i.i.d. Multinomial($2, (q_a)_{a=1}^A$). Then

$$n_{1\cdot} \sim \text{DM}(2, (\alpha_a)_{a=1}^A)$$

where DM denotes the Dirichlet–Multinomial distribution: $X \sim \text{DM}(n, (\alpha_a)_{a=1}^A)$ means

$$P(X = x) = \int \frac{n!}{\prod_a x_a!} \prod_a q_a^{x_a} \frac{\Gamma(\sum_a \alpha_a)}{\prod_a \Gamma(\alpha_a)} \prod_a q_a^{\alpha_a - 1} dq = \left\{ \frac{n!}{\prod_a x_a!} \right\} \times \left\{ \prod_a \frac{\Gamma(\alpha_a + x_a)}{\Gamma(\alpha_a)} \right\} \times \frac{\Gamma(\sum_a \alpha_a)}{\Gamma(\sum_a \alpha_a + n)},$$

so long as $\sum_a x_a = n$. Furthermore, again by conjugacy, for $i = 2, 3, \dots, I$,

$$n_{i\cdot} | (n_{j\cdot})_{j=1}^{i-1} \sim \text{DM}(2, (\alpha_a + T_{i-1,a})_{a=1}^A)$$

where $T_{i-1,a} = \sum_{j=1}^{i-1} n_{ja}$.

Factorising these distributions over alleles, we find that individual allele counts have Beta–Binomial conditional distributions:

$$n_{ia} | \{n_{jb}, j < i, \forall b\}, \{n_{ib}, b < a\} \sim \text{BB}(2 - S_{i,a-1}, \alpha_a + T_{i-1,a}, \beta_a + U_{i-1,a}) \quad (1)$$

Here BB is the Beta–Binomial distribution: $\text{BB}(n, \alpha, \beta)$ is the same as $\text{DM}(n, (\alpha, \beta))$, $\beta_a = \sum_{b>a} \alpha_b$, $S_{ia} = \sum_{b=1}^a n_{ib}$ as in CGLM and $U_{i-1,a} = \sum_{b>a} T_{i-1,b}$. Note that $\text{BB}(1, \alpha, \beta)$ is just Bernoulli($\alpha/(\alpha + \beta)$). Equation (1) exhibits association among the n_{ia} that is positive across i and negative across a , as would be expected.

In the large-database limit, $\alpha_a \rightarrow \infty$ but $\alpha_a / \sum_a \alpha_a \rightarrow q_a$, and the Beta–Binomial conditional probabilities (1) become

$$n_{ia} | \{n_{jb}, j < i, \forall b\}, \{n_{ib}, b < a\} \sim \text{Binomial}(2 - S_{i,a-1}, q_a / \sum_{b \geq a} q_b) \quad (2)$$

as in Section 2.4.1 of CGLM.

Graversen's (2013) R package `DNAmixtures` can readily be amended to use (1) instead of (2) in a Bayes net computation to sum the terms in equation (8) of CGLM. The corresponding DAG is now considerably more complex, due to the presence of the additional nodes T_{ia} and U_{ia} , and the computation runs much more slowly. (Therese showed us how to amend our amendment to her code to use a more efficient elimination order, and this improved the times.)

Our limited numerical experiments with casework data using this code reveal a curiously mixed picture: uncertainty in allele frequencies may either increase or decrease the weight of evidence $\log_{10}(\text{LR})$, depending on the example. This is in contrast to all our earlier examples, with either allele-presence indicator traces (in Green and Mortera, 2009) or with the Cowell, Lauritzen and Mortera (2007b) model (unpublished), in which this uncertainty always reduced the weight of evidence. This needs further study, but we surmise that the difference might be attributable to maximising out of parameters, in contrast to a more fully Bayesian approach.

In the literature, other phenomena causing dependence among DNA profiles, such as identity by descent, have been modelled in a way leading to the same probabilistic dependence as in the analysis above.