

# Vignette to *gammadnamix*: An open source R-package to handle continuous model for evaluating STR DNA mixture evidences with artefacts.

Øyvind Bleka

January 8, 2015

## Abstract

In this vignette we are introducing a forensic software, the package *gammadnamix*, which utilises the peak height information in STR DNA data and takes care of artefacts such as allele drop-out, drop-in and stutters. Our model assumes that the resulting peak heights from the PCR-process are gamma-distributed as presented by Cowell et al. [4]. We will incorporate this model into a Bayesian framework to make it available for users to specify their own prior information about allele drop-in and stutters. *gammadnamix* is an open-source free R-package which does model inference using maximum likelihood estimation or integrated likelihood for any number of known and unknown contributors and for multiple replicates.

## 1 Introduction

Weight-of-evidence evaluation where peak heights or peak areas are incorporated requires a continuous mixture model which contains nuisance parameters. Such parameters in the models are often apriori unknown and are either estimated using maximum likelihood estimator or integrated out by turning into a Bayesian framework [6]. At the moment there exists several commercial softwares that are based on continuous models which take peak heights into account for weighting of evidence [12], [9], [4]. A full Bayesian evaluation of such models are very computationally intensive, even for 2-person mixtures. Taylor et al. [12] and Perlman et al. [9] suggests models that are so complex (high dimensional) that numerical integral methods would be intractable to use. As an approximation, the integral is estimated by Markov Chain Monte Carlo (MCMC) simulations. However, the models considered in [2], [3], [4] and [11] are more simplified and could be calculated with a required accuracy using numerical integration [8].

We consider the model proposed by Cowell et al. [4] which handles different mixture proportion, gamma distributed peak height variability and stutter ratio as nuisance parameters. They evaluated the weight of evidence frequently based on the maximum log-likelihood properties. In our software we do in

addition Bayesian inference by integrating out the nuisance parameters in the model such that the weight of evidence formula becomes similar as the Bayes Factor<sup>1</sup>. Both the frequentistic inference and the Bayesian inference for any given hypotheses are implemented in the open-source R-package *gammadnamix*.

## 2 The continuous model

For a particular locus we observe alleles  $\mathbf{A} = (A_1, \dots, A_J) \subseteq \mathbb{A}$  with corresponding peak heights  $\mathbf{Y} = (Y_1, \dots, Y_J)$  where  $J$  is number of observed alleles. We let the outcome of possible alleles be given as  $\mathbb{A}$ . We introduce an individual  $k \in \{1, \dots, C\}$  to have genotype  $g_k = (a_k/b_k)$ , where alleles  $a_k, b_k \in \mathbb{A}$ . By combining the genotypes for the  $C$  contributors we obtain the genotype vector  $\mathbf{g} = (g_1, \dots, g_C)$ .

### 2.1 The allele peak height distribution

Cowell et al. [4] assumed that the peak height at allele  $A_j$  indexed at  $j \in \{1, \dots, J\}$  for contributor  $k$  is distributed as

$$Y_{jk} \sim_{iid.} \text{gamma}(\rho n_{jk} m_k, \tau)$$

where  $\rho$  is proportional to the total amount of DNA in the mixture before amplification and  $\tau$  is the scale parameter. For a given genotype  $g_k$ ,  $n_{jk} = \mathbb{I}(A_j = a_k) + \mathbb{I}(A_j = b_k)$  gives contribution amount for contributor  $k$  to allele  $A_j$ . Further, *iid.* means that the observed peak heights are drawn independently for each contributor and allele. The mixture proportion parameter  $\mathbf{m} = (m_1, \dots, m_C)$  is defined on *simplex*( $\mathbf{1}_C$ ) (i.e.  $\mathbf{m} \in [0, 1]^C$  with  $m_C = \sum_{k=1}^{C-1} m_k = 1$ ).

By summing the independent peak heights over each contributors the observed peak heights is given to have the following distribution

$$Y_j \sim_{iid.} \text{gamma}(\rho \mathbf{m}^T \mathbf{n}_j, \tau)$$

where  $\mathbf{n}_j = (n_{j1}, \dots, n_{jC})$  is the vector giving the contribution amount for each contributor (found through  $\mathbf{g}$ ).

#### 2.1.1 Model properties

The peak height mean and variance of the model is given as  $\mu_j = E[Y_j] = \rho \tau \mathbf{m}^T \mathbf{n}_j$  and  $Var[Y_j] = \rho \tau^2 \mathbf{m}^T \mathbf{n}_j$ . Hence the coefficient of variation is given as  $CV[Y_j] = \rho^{-\frac{1}{2}}$  which is a strongly decaying function down to 200rfu and hence it will model the low-template DNA effect in that smaller peak heights has greater variability than larger peak heights.

#### 2.1.2 Reparameterization

As in [4] the parameters of the model are reparameterized such that that the parameters are orthogonally spanned and easier to interpret directly. By

---

<sup>1</sup>The Bayes inference analogy to Likelihood Ratio is Bayes Factor (by assuming equal prior for the hypotheses)

requiring  $\mu = \rho\tau \geq 0$  to be the expected mean peak height (for a single heterozygote contributor) and  $\sigma = \rho^{-\frac{1}{2}} \geq 0$  to be the coefficient of variance of the peak heights, we have that  $\rho = \sigma^{-2}$  and  $\tau = \mu\sigma^2$ . However, when presenting the gamma model we will still use  $\rho$  and  $\tau$  for simplicity.

## 2.2 Genotype probabilities

The continuous model requires that we make a suggestion of a combined genotype profiles  $\mathbf{g}$  as a possible contribution to the evidence. Often a hypothesis assumes that one or more profile components in  $\mathbf{g}$  are unknown. When such genotype profiles are unknown, it is reasonable to assign these components to a discrete random variable set, such that a given combined genotype profile has probability  $p(\mathbf{g})$ . In our software we will use the point estimates of the population allele-frequencies to determine  $p(\mathbf{g})$ . However we will extend the model of  $p(\mathbf{g})$  to include the scalar  $f_{st}$  which takes into account that the components of  $\mathbf{g}$  are a sub-population of the population (which the point estimates are based on) [5].

With the Hardy-Weinberg Equilibrium assumption ( $f_{st} = 0$ ), we have that

$$p(g_k) = \begin{cases} 2^{\mathbb{I}(a_k \neq b_k)} P(a_k)P(b_k) & \text{if } k \text{ is unknown contributor} \\ 1 & \text{if } k \text{ is known contributor} \end{cases}$$

The correction formula when  $f_{st} > 0$  is given as:

$$p(g_k) = \begin{cases} 2^{\mathbb{I}(a_k \neq b_k)} \prod_{j \in \{a_k, b_k\}} \left( \frac{u_j f_{st} + (1-f_{st})P(j)}{1+(v-1)f_{st}} \right) & \text{if } k \text{ is unknown contributor} \\ 1 & \text{if } k \text{ is known contributor} \end{cases}$$

where  $P(j)$  is allele frequency of allele  $j \in \{a_k, b_k\}$ ,  $u_j$  denotes previously number of sampled alleles of allele  $j$  and  $v$  denotes previously number of total sampled alleles.

By assuming that each contributors are independent of each others (given the  $f_{st}$  relation) we have that  $p(\mathbf{g}) = \prod_{k=1}^C p(g_k)$ . This is the Hardy-Weinberg Linkage assumption.

## 2.3 Modeling artefacts

In this subsection we follow the models provided by Cowell et al. [4] to handle artefacts as stutter and drop-out. To handle drop-in events, we provide our own method similar to Puch-Solis [10]. We also remind the reader that the expected model peak height of allele  $j$  is denoted as  $\mu_j = E[Y_j]$  as given in earlier sub-section.

### 2.3.1 Incorporate stutter

The stutter ratio parameter  $\xi$  is modelled to be the ratio of non-stuttered contributing peak heights from allele  $a_{j+1}$  to allele  $a_j$ . The non-stuttered and stuttered peak heights are independently modelled as gamma  $((1-\xi)\rho\mathbf{m}^T \mathbf{n}_j, \tau)$

and gamma  $(\xi \rho \mathbf{m}^T \mathbf{n}_{j+1}, \tau)$  respectively such that the summed peak heights on allele  $j$  is independent distributed as

$$Y_j \sim \text{gamma}(\rho \mathbf{m}^T ((1 - \xi \mathbb{I}(a_{j-1} \in \mathbf{A})) \mathbf{n}_j + \xi \mathbb{I}(a_{j+1} \in \mathbf{A}) \mathbf{n}_{j+1}), \tau) = f_j(y_j | \mathbf{g}, \boldsymbol{\theta})$$

where  $f_j$  becomes the density function of the peak height given the set of genotype profiles  $\mathbf{g}$  and the reparameterized model parameters is given as  $\boldsymbol{\theta} = (\mathbf{m}, \mu, \sigma, \xi)$ . Note that the indicator  $\mathbb{I}(a_{j-1} \in \mathbf{A})$  is true only if  $a_{j-1}$  exist in the provided allele outcome. Also note that indicator  $\mathbb{I}(a_{j+1} \in \mathbf{A})$  is true only if allele  $a_{j+1}$  exists in the provided allele-frequencies.

Note that if the peak height of parent allele,  $a_{j-1}$ , is zero, but it still has expected contribution, it is still expected to contribute  $\xi$  amount to  $a_j$ . Also note that if allele  $a_j$  can not be stuttered by any parent allele  $a_{j-1}$  (i.e. there doesn't exist any parent allele),  $\xi = 0$ , and hence the model for the corresponding peak height reduces to

$$Y_j \sim \text{gamma}(\rho \mathbf{m}^T \mathbf{n}_j, \tau)$$

Last, note that from this we have that the conditional distribution for a relative amount of stutter for peak height  $Y_j$  is given as

$$S_j \sim \text{beta}(\xi \rho \mathbf{m}^T \mathbf{n}_j, (1 - \xi) \rho \mathbf{m}^T \mathbf{n}_{j+1})$$

### 2.3.2 Handling allele drop-out

Sometimes, the peak heights of a contributing allele (i.e. the allele of a true donor) is not seen in the evidence because the amplification produces a rf peak height below a certain detection threshold  $T$ . We define this as an allele drop-out. For a given genotype  $\mathbf{g}$  in our model, we have the set of non-observed alleles given in  $\mathbf{g}$  as

$$\Omega = \{A_j \in \mathbf{g} : \mu_j > 0 \cap Y_j < T\}$$

If  $\Omega$  is empty (i.e. we have no allele drop-out),  $p(\Omega | \mathbf{g}, \boldsymbol{\theta}) = 1$ . Otherwise we incorporate the allele drop-out event in the model by including the probability factor

$$p(\Omega | \mathbf{g}, \boldsymbol{\theta}) = \prod_{A_j \in \Omega} Pr(y_j < T | \mathbf{g}, \boldsymbol{\theta}) = \prod_{A_j \in \Omega} \int_0^T f_j(x | \boldsymbol{\theta}, \mathbf{g}) dx$$

### 2.3.3 Handling allele drop-in

The event of an allele drop-in is defined to be a non-explained peak height which is not contributed by any contributors nor as a stuttered peak height. The dropped in allele is assumed to originate from a contaminating cell from the population with probability  $p_C$ . Such contaminations are often identified by running "negative controls" which aims to detect contaminating alleles. By summarising the negative controls for a lab it is possible to both estimate the probability of drop-in,  $p_C$ , and infer a drop-in density function as a function of allele peak height  $Y$  given as  $Pr(Y = y | \text{Drop-in}) = h(y)$ . Puch-Solis [10]

found that  $h$  fitted a gamma distribution well for SGM with threshold given as 25. However we assume that the detection threshold is given above  $T = 50$ . We found that a shifted exponential density function with a rate parameter  $\lambda$  starting from threshold  $T$  was a reasonable choice for  $h$ . We will further assume  $\lambda$  to be prior knowledge to the model such that  $h(y) = h(y|\lambda)$ .

For a given genotype combination  $\mathbf{g}$ , we define the allele drop-in set as

$$\Psi = \{A_j \in \mathbf{A} : (\mu_j = 0 \cap Y_j \geq T)\}$$

which is the set of non-explained alleles (a peak height is observed for an allele which the model assumes no contribution to).

We incorporate the drop-in peak height information to the likelihood with the factor

$$p(\Psi|\mathbf{g}) = \begin{cases} \prod_{A_j \in \Psi} (p_C P(A_j) h(Y_j|\lambda)) & \text{if } \Psi \text{ is non-empty} \\ (1 - p_C) & \text{if } \Psi \text{ is empty} \end{cases}$$

where

$$h(y|\lambda) = \lambda \exp^{-\lambda(y-T)} \mathbb{I}(y \in [T, \infty))$$

is the shifted exponential density function with parameter  $\lambda$ .

#### 2.3.4 Handling non-contributed markers

If the stain is low-template, there may be some markers which have no observed peak heights. The explanation for this will be that at least one peak height is below the detection threshold. This is a part of the observation and should be taken into account to the model to lower the parameter(s) connected to the amount of dna.

### 2.4 Apriori distribution for model parameters

A Bayesian framework requires prior distribution to the unknown parameters in the model to be incorporated. Prior distributions for the parameters of  $\boldsymbol{\theta}$  are given as

$$p(\boldsymbol{\theta}) = p(\mathbf{m})p(\tau)p(\xi)$$

where we specify non-informative prior on each parameters

$$\begin{aligned} p(\mathbf{m}) &= \text{Dirichlet}(\mathbf{1}_C) \\ p(\mu) &= \mathbb{I}(\mu \in [0, \mu_1]) \\ p(\sigma) &= \mathbb{I}(\sigma \in [0, \sigma_1]) \\ p(\xi) &= \mathbb{I}(\xi \in [0, \min(\xi_1, 1)]) \end{aligned}$$

The posterior distribution of  $\mu$  and  $\sigma$  should be investigated closer (using MCMC) in order to find a suitable values of  $\mu_1$  and  $\sigma_1$ . The prior density of the stutter ratio can be choosen flexible in the routine such that it may be inferred by stutter information from the laboratory (i.e. choosing  $p(\xi) \approx p(\xi|\text{stutter-data})$ ).

## 2.5 Example

It is perhaps hard to adapt these formulas and definition without having a concrete example to apply them on. We start with a 2-person mixture (or a possible single-source).

Let  $E = \{\mathbf{A} = (15, 16, 17, 18, 19), \mathbf{Y} = (100, 500, 0, 1000, 0)\}$  be the observed evidence with allele information and peak heights.

## 3 Inference of the model

### 3.1 One replicate

A specified hypothesis  $H$  defines a set  $\mathbb{Q} = \mathbb{Q}(H)$  which consists of all possible combined genotypes  $\mathbf{g}$  which satisfies  $H$ . Let the mixture evidence  $E = (\mathbf{A}, \mathbf{Y})$  be the observed data. When threating  $I$  loci simultaneously we let the vectorized alleles, peak heights and the combined genotype outcome be denoted as  $\mathbf{Y} = (\mathbf{Y}_1, \dots, \mathbf{Y}_I)$ ,  $\mathbf{A} = (\mathbf{A}_1, \dots, \mathbf{A}_I)$  and  $\mathbb{Q} = (\mathbb{Q}_1, \dots, \mathbb{Q}_I)$ . We are marginalizing out all combined genotypes (latent variables) with respect to the possible outcome, such that simultaneous over all loci, the likelihood function of the model parameters  $\boldsymbol{\theta} = (\mathbf{m}, \mu, \sigma, \xi)$  becomes

$$L(\boldsymbol{\theta}|E, H) = \prod_{i=1}^I L_i(\boldsymbol{\theta}|E_i, H) = \prod_{i=1}^I \left( \sum_{\mathbf{g}_i \in \mathbb{Q}_i} p(E_i|\mathbf{g}_i, \boldsymbol{\theta}) p(\mathbf{g}_i|f_{st}, H) \right)$$

For a given locus  $i$  we have that the probability of the observations is given by the expression

$$p(E|\mathbf{g}, \boldsymbol{\theta}) = p(\Psi|\mathbf{g}) p(\Omega|\mathbf{g}, \boldsymbol{\theta}) \prod_{A_j \in \chi} f_j(y_j|\mathbf{g}, \boldsymbol{\theta})$$

where  $\chi = \mathbf{A} \setminus \Psi = \{A_j \in \mathbf{A} : (\mu_j > 0 \cap Y_j \geq T)\}$  is the observed set of alleles where we expect a contributing peak height. The peak height density function  $f_j$ , drop-in density function  $p(\Psi|\mathbf{g})$  and drop-out density function  $p(\Omega|\mathbf{g}, \boldsymbol{\theta})$  were all defined in the previous section.

### 3.2 Multiple replicates

Assume that a sample has been replicated such that each repliated sample is assumed to contain the same contributors and satisfy the same model assumptions. Consider the number of replicates as  $S$ , and let the observed replicated evidences to be given as  $\mathbf{E} = (E^{(1)}, \dots, E^{(S)})$ . Then the formula above is extended to

$$L(\boldsymbol{\theta}|\mathbf{E}, H) = \prod_{i=1}^I \left( \sum_{\mathbf{g}_i \in \mathbb{Q}_i} p(\mathbf{g}_i|f_{st}, H) \prod_{s=1}^S p(E_i^{(s)}|\mathbf{g}_i, \boldsymbol{\theta}) \right)$$

by assuming the conditional replicated observations are independent of each others.

This differs some to how [4] defines the likelihood for multiple replicates in that they assume different parameter sets  $\theta$  for each replicated observation:

$$L(\theta_1, \dots, \theta_S | \mathbf{E}, H) = \prod_{i=1}^I \left( \sum_{\mathbf{g}_i \in \mathbb{Q}_i} p(\mathbf{g}_i | f_{st}, H) \prod_{s=1}^S p(E_i^{(s)} | \mathbf{g}_i, \theta_1, \dots, \theta_S) \right)$$

which makes the model more flexible (but with less power).

### 3.3 Reducing the mixture proportion space

Note that  $\mathbb{Q} = \{\mathbb{Q}_1, \dots, \mathbb{Q}_I\}$  must contain all possible genotype combinations, also symmetrical possibilities for unknown genotypes. The reason for this is that when we consider multiple loci simultaneously, the order of the combined genotype profiles is relevant when scaling with the global mixture proportion parameter. However, we will now show how the symmetry in  $\mathbb{Q}$  may support to reduce the space of the mixture proportions.

Denote the parameter space of the mixture proportions  $\mathbf{m}$  as  $M$  where  $M = \text{simplex}(\mathbf{1}_C)$  (i.e.  $\mathbf{m} \in [0, 1]^C$  with restriction  $m_C = 1 - \sum_{k=1}^{C-1} m_k \geq 0$ ). When there is at least 2 unknown contributors in a hypothesis, it is possible to reduce the space  $M$ . Let  $\mathbf{m} = \{\mathbf{m}_K, \mathbf{m}_U\} \in M$  such that  $\mathbf{m}_K$  is the vector for known contributors, while  $\mathbf{m}_U$  is the vector for unknown contributors. By requiring the elements in  $\mathbf{m}_U$  to be in decreasing order, changing  $M$  to  $\tilde{M}$ , this avoids a multimodal posterior density in the space of  $\mathbf{m}$  when the hypothesis includes at least 2 unknowns.

### 3.4 Estimating mode of the posterior distribution (Frequentistic inference)

Instead of working in the restricted space of  $\theta$  we transform<sup>2</sup> to the unrestricted space of  $\phi$  and make unrestricted optimization of the posterior distribution

$$p(\theta(\phi) | E, H) \propto L(\theta(\phi) | E, H) p(\theta(\phi))$$

to find the restricted maximum argument  $\theta^*$ . The maximum likelihood for the hypothesis is then given as  $L(\theta^* | E, H) p(\theta^*)$ . Because the uncertainty of  $\theta$ , the routine also returns the posterior covariance structure  $\Sigma$  on  $\theta^*$  using the delta method approximation<sup>3</sup>.  $\Sigma$  can also further be used as the covariance of the proposal distribution when running Random Walk Metropolis Hastings. This can be applied to draw samples from  $p(\theta | E, H)$ . It can also be used together with the maximum likelihood value to make an Laplace approximation to the integrated likelihood of the evidence.

<sup>2</sup>See Reparametrization Appendix A

<sup>3</sup>To get the full covariance structure of  $\theta$  we used the derivings in Appendix A

### 3.5 The integrated likelihood of evidence (Bayesian inference)

The maximum likelihood in previous sub-section is depending on the point estimation given as the mode of the posterior distribution, and hence it does not take into account non-gaussian structure in the posterior distribution. We let  $\boldsymbol{\theta} = (\mathbf{m}, \boldsymbol{\theta}_2)$  where  $\boldsymbol{\theta}_2 = (\mu, \sigma, \xi)$  excludes the mixture proportion parameters. Since we are working in a Bayesian framework, the integrated likelihood of the given hypothesis  $H$  can be calculated by integrating out the parameters as

$$L(H|E) = \int_{\boldsymbol{\theta}} L(\boldsymbol{\theta}|E, H) p(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int_{\mathbf{m} \in \mathcal{M}} p(\mathbf{m}|E, H) p(\mathbf{m}) d\mathbf{m} \quad (1)$$

where

$$p(\mathbf{m}|E, H) = \int_{\boldsymbol{\theta}_2} L(\mathbf{m}, \boldsymbol{\theta}_2|E, H) p(\boldsymbol{\theta}_2) d\boldsymbol{\theta}_2$$

Note here that  $\mathcal{M} = \text{simplex}(\mathbf{1}_C)$  (i.e.  $\mathbf{m} \in [0, 1]^C$  with restriction  $m_C = 1 - \sum_{k=1}^{C-1} m_k \geq 0$ ).

#### 3.5.1 Integral over a simplex

To integrate with respect to  $\mathcal{M}$  we need to take the simplex restrictions into account. We have that since  $m_C = 1 - \sum_{k=1}^{C-1} m_k$ , we only need to consider the integral over the variables  $(m_1, \dots, m_{C-1})$ . With the restriction  $\sum_{k=1}^{C-1} m_k \leq 1$  we get the integrated likelihood as

$$L(H|E) = \int_{m_1 \in [0, 1]} \int_{m_2 \in [0, 1-m_1]} \cdots \int_{m_{C-1} \in [0, 1-\sum_{k=1}^{C-2} m_k]} p(\mathbf{m}|E, H) p(\mathbf{m}) dm_{C-1} \cdots dm_1$$

#### 3.5.2 Reduction of mixture proportion space

Recall that we may reduce the space of  $\mathcal{M}$  to  $\tilde{\mathcal{M}}$  if the hypothesis includes at least 2 unknowns. The integrated likelihood can then be written as

$$L(H|E) = U! \int_{\mathbf{m} \in \tilde{\mathcal{M}}} p(\mathbf{m}|E, H) p(\mathbf{m}) d\mathbf{m}$$

where  $U = |\mathbf{m}_U|$ , the number of unknown contributors. Reducing the space from  $\mathcal{M}$  to  $\tilde{\mathcal{M}}$  is efficient since we reduce the space of  $\mathbf{m}_U$  by require the mixture proportions of the unknown contributors to be decreasing ordered.

Estimation of these integrals are done by numerical multidimensional integration.

#### 3.5.3 Laplace approximation of the integrated likelihood

Consider  $\boldsymbol{\theta}^*$  to be the estimated mode and  $\Sigma(\boldsymbol{\theta}^*)$  the inverse negative hessian from previous section. When there is sufficient large enough amount of data, the likelihood will be approximately normal distributed (Central Limit Theorem).



A second order taylor expansion (see Appendix A.4 for technical details) gives a rough approximation to the integrated likelihood as

$$L(H|E) \approx U!(2\pi)^{\frac{p}{2}} |\Sigma(\boldsymbol{\theta}^*)|^{\frac{1}{2}} L(\boldsymbol{\theta}^*|E, H) p(\boldsymbol{\theta}^*)$$

where  $U$  is number of unknowns in hypothesis  $U$  and  $p$  is number of paramaters in the model ( $p = C + 2$ ). If at least 2 unknowns are considered in  $H$ , the posterior of the mixture proportions will be multimodal and hence we must scale with number of possible ordered outcomes for the unknown contributors.

### 3.5.4 Numerical integration methods

The integrated likelihood in equation 1 requires to calculate an integral over the space of the noisance parameters  $\boldsymbol{\theta}$ . Steven G. Johnson implemented the R-package *cubature* to do adaptive multidimensional integration based on algorithms in [7] and [1]. A relative error requirement is given as an input for the accuracy of the integration. The method also estimates the relative error which can again be used to provide an error interval for the integrated likelihood (see Appendix A.3).

## 3.6 Uncertainty of LR over the posterior parameter space

Using a Markov Chain Monte Carlo Metropolis Hastings Random Walk method we are able to provide samples from the posterior distribution of the parameters in the model. This may be done for each of two competing hypothesis  $H_p$  and  $H_d$  such that  $\theta_p^{(m)} \sim p(\theta_p|E)$  and  $\theta_d^{(m)} \sim p(\theta_d|E)$  for a iteration  $m$ . At each iteration we calculate  $LR^{(m)} = \frac{p(\mathbf{E}|\theta_p^{(m)}, H_p)}{p(\mathbf{E}|\theta_d^{(m)}, H_d)}$ . By providing  $M$  iterations we obtain a distribution for  $LR$  where the uncertainty of the parameters under each hypothesis are the distributing values.

## 3.7 Mixture deconvolution

The specified model can be used to propose the most probable genotype combinations for the given evidence(s)  $E$ . Let a suggested combined genotypes over all loci be given as  $G = (\mathbf{g}_1, \dots, \mathbf{g}_I) \in \mathbb{Q}(H)$  where  $H$  is a specified hypothesis.

### 3.7.1 Conditioned on maximum likelihood estimates

Assume that  $\boldsymbol{\theta}^*$  is the maximum likelihood estimates obtained from the optimizing routine. Then by using Bayes theorem we have that

$$p(G|E, \boldsymbol{\theta}^*, H) = \frac{p(E|G, \boldsymbol{\theta}^*) p(G|H)}{p(E|H, \boldsymbol{\theta}^*)}$$

where  $G|H$  is independend of  $\boldsymbol{\theta}^*$  and  $p(E|H, \boldsymbol{\theta}^*)$  is the maximized likelihood value found in the optimizing routine.

We require to return a ranked list of genotypes  $\mathbb{G}$  which ensures that  $\sum_{G \in \mathbb{G}} p(G|E, \boldsymbol{\theta}^*, H) > \alpha$  for a given selection of  $\alpha$ . Because of this, the number of elements in  $\mathbb{G}$  will vary. Details of finding the jointly top ranked genotype combination probability may be found in Appendix.

### 3.7.2 Integrating out parameter uncertainty

We may integrate the uncertainty of  $\theta$  as following:

$$p(G|E, H) = \frac{p(G|H)}{p(E|H)} \int_{\theta} \prod_{i=1}^I p(E_i|g_i, \theta) p(\theta) d\theta$$

where  $p(E|H)$  is the marginalised likelihood value found in the integration routine. However this is not yet implemented...

## 3.8 Weight-of-evidence

Often a specification of two rivaling hypothesis,  $H_p$  which defines the prosecution hypothesis, and  $H_d$  which defines the defence hypothesis, has been proposed. The Likelihood ratio of the likelihood for each of two such hypothesis are commonly used as a metric in order to determine which hypothesis that are most likely true for a given evidence  $E$ . For a given specified hypothesis  $H$  we can construct a likelihood  $L(\theta|E, H)$  which is based on the model described as earlier in the section. The likelihood function describes a model family which is a function of the unknown parameters  $\theta$ . We now introduce two ways to proceed to do weight-of-evidence.

The Likelihood ratio method is often used to determine if a suspect  $S$  contributes to an evidence  $E$  or not (determine an exclusion). The likelihood function requires that we assume the number of individuals  $C$  contributing to the evidence  $E$  in order to calculate the likelihood under each of the hypotheses. In order to do the determination, the two following hypotheses may be specified:

$H_p$  :Suspect  $S$  and  $C - 1$  unknown individuals contributes to evidence  $E$ .

$H_d$  : $C$  unknown individuals contributes to evidence  $E$ .

The logic behind this specification is that we want to determine how much more likely it is that the genotype profile for suspect  $S$  contributes to the evidence  $E$  compared to a random man from the population. Note that the more information we know about the true contributors to the profiles, the better determination can we do. For instance this could be a victim  $V$  which we are absolutely sure contributes to  $E$ . The hypotheses are then modified such that one unknown individual in each of the hypotheses are exchanged with the genotype profile of  $V$ . All references specified under  $H_p$  but not under  $H_d$  is assumed to be a known non-contributor under  $H_d$ . This is necessary information if  $f_{st} > 0$ .

In the following sub section we specify technically how different approaches can be used to get a weight-of-evidence value.

### 3.8.1 Weight-of-evidence based on Maximized Likelihood(frequentistic inference)

$$LR(H_p, H_d|E, \theta_p^*, \theta_d^*) = \frac{L(\theta_p^*|E, H_p)}{L(\theta_d^*|E, H_d)} \quad (2)$$

where  $\theta^*$  is the maximum likelihood estimate given the subscripted hypothesis.

### 3.8.2 Weight-of-evidence based on Integrated Likelihood (bayesian inference)

$$LR(H_p, H_d|E) = \frac{L(H_p|E)}{L(H_d|E)} = \frac{\int_{\theta_p} L(\theta_p|E, H_p) \theta_p}{\int_{\theta_d} L(\theta_d|E, H_d) \theta_d} \quad (3)$$

### 3.8.3 Weight-of-evidence based on Laplace approximation

$$LR(H_p, H_d|E) \approx \frac{L(\theta_p^*|E, H_p) U_p!}{L(\theta_d^*|E, H_d) U_d!} \sqrt{\frac{|H(\theta_d^*|E, H_d)|}{|H(\theta_p^*|E, H_p)|}} \quad (4)$$

where  $\theta^*$  is the maximum likelihood estimate given for the subscripted hypothesis,  $U_p$  and  $U_d$  are number of unknown contributors under  $H_p$  and  $H_d$  respectively and  $|H|$  is the determinant of the hessian matrix of the log-likelihood function under a given hypothesis.

## 3.9 Database search

Instead of considering a suspect  $S$  in  $H_p$  we exchange this to be  $S = j$  where  $j$  is an individual in the database. We name this exchanged hypothesis for  $H_p(j)$ . Under the defence hypothesis  $H_d$ , the individual  $j$  is not explicit considered in the hypothesis. However, with the specification  $f_{st} > 0$ , the information that  $j$  is a known non-contributor requires that the likelihood under the hypothesis  $H_d$  is calculated again for each individual  $j$  in the database. Therefore we define the defence hypothesis as  $H_d(j)$ .

In order to calculate the Likelihood ratio for each individuals in the database we must calculate  $L(\theta_p^*|E, H_p(j))$  and  $L(\theta_d^*|E, H_d(j))$  where  $\theta^*$  is the maximum likelihood estimate given the subscripted hypothesis. From this we have

$$LR(j|E, \theta_p^*, \theta_d^*) = \frac{L(\theta_p^*|E, H_p(j))}{L(\theta_d^*|E, H_d(j))} \quad (5)$$

In a Bayesian framework it is also possible to integrate out the uncertainty of the parameters  $\theta$  in the likelihood formula as described before such that the Likelihood Ratios for each individuals in the hypothesis becomes

$$LR(j|E) = \frac{L(H_p(j)|E)}{L(H_d(j)|E)} = \frac{\int_{\theta_p} L(\theta_p|E, H_p(j)) \theta_p}{\int_{\theta_d} L(\theta_d|E, H_d(j)) \theta_d} \quad (6)$$

such that the Likelihood ratio expression becomes independent of the choice of the parameter  $\theta$ .

Note that if  $f_{st} = 0$  is specified,  $H_d(j) = H_d$  for all individual  $j$  in database such that the denominator of the likelihood ratio expression becomes a constant for all  $j$ . This makes the database search much faster.

### 3.10 Validation of the maximum likelihood based model

For a given locus  $i$  and model parameter  $\theta$  we have that the peak heights  $\mathbf{y}$  has the density function

$$p(\mathbf{y}|\theta, E, H) = \sum_{\mathbf{g} \in \mathbb{Q}} p(\mathbf{y}|\mathbf{g}, \theta) p(\mathbf{g}|f_{st}, H)$$

#### 3.10.1 Distribution of peak heights for each contributor

Given the model parameters  $\theta$ , the distributing (non-stuttering) heterozygote peak height  $y$  for contributor  $k$  is given as  $\text{gamma}(y, \rho * m_k, \text{scale} = \tau)$ . This distributions can be used as an exploratory tool to see what peak height contribution each of the contributors are expected to have when  $\theta$  is the maximum likelihood estimates under a specific hypothesis.

#### 3.10.2 Check that the gamma distribution is a reasonable assumption

If the model fits the data, we would expect that the cumulative probabilities up to the observed peak heights are uniformly distributed. Since we do this on only for the observed peak heights above threshold  $T$ , we need to condition on this. We define the conditional probability density function for allele  $A_j$  given the other alleles  $\mathbf{A}_{-j}$  and that they are observable above threshold  $T$  as

$$p(y_j|\mathbf{y}_{-j}, \theta, E, H, Y_j \geq T) \propto p(y_j, \mathbf{y}_{-j}|\theta, H, Y_j \geq T)$$

Our aim is to calculate

$$p_{Y_j} = \text{Pr}(y_j \leq Y_j | \mathbf{y}_{-j} = \mathbf{Y}_{-j}, \theta, H, Y_j \geq T)$$

We estimate this quantity as

$$\hat{p}_{Y_j} = \frac{\int_T^{Y_j} p(x|\mathbf{y}_{-j} = \mathbf{Y}_{-j}, \theta, H)}{\int_T^U p(x|\mathbf{y}_{-j} = \mathbf{Y}_{-j}, \theta, H)}$$

where  $U$  is a sufficient large number such that  $p(U, \mathbf{y}_{-j}|\theta, H) \approx 0$ .

If the fitted model fits the theoretical assumption,  $p_{Y_j} \sim \text{Unif}(0, 1)$ .

#### 3.10.3 Dropout properties of the fitted model

Consider a single heterozygote allele peak height  $y$ . The dropout distribution for this peak height for a particular contributor  $k$  is given as

$$\text{Pr}(y = 0 | n_j = 1, m_k, \mu, \sigma) = \Gamma(T | \rho * m_k, \tau)$$

Where  $\mu$  and  $\sigma$  are unknown parameters. Using posterior samples of  $\sigma^{(m)} \sim p(\sigma|E)$  we are able to construct dropout density functions as a function of  $\mu$ .

## References

- [1] J. Berntsen, T. O. Espelid, and A. C. Genz. An adaptive algorithm for the approximate calculation of multiple integrals. *ACM Transactions on Mathematical Software*, 17(4):437–451, 1991.
- [2] R. G. Cowell, S. L. Lauritzen, and J. Mortera. A gamma model for dna mixture analysis. *Bayesian Analysis*, 2(2):333–348, 2007.
- [3] R. G. Cowell, S. L. Lauritzen, and J. Mortera. Identification and separation of dna mixtures using peak area information. *Forensic Science International: Genetics*, 166, 2007.
- [4] R. G. Cowell, T. Graversen, S. L. Lauritzen, and J. Mortera. Analysis of forensic dna mixtures with artefacts. *Appl. Statist.*, 64(1):1–32, 2015.
- [5] J. Curran, P. Gill, and R. Bill, M. Interpretation of repeat measurement dna evidence allowing for multiple contributors and population substructure. *Forensic Science International*, 148:47–53, 2005.
- [6] I. W. Evett, P. Gill, and J. A. Lambert. Taking account of peak areas when interpreting mixed DNA profiles. *J Forensic Sci.*, 43:62–69, 1998.
- [7] A. C. Genz and A. A. Malik. An adaptive algorithm for numeric integration over an n-dimensional rectangular region. *Journal of Computational and Applied Mathematics*, 6(4):295–302, 1980.
- [8] T. Hahn. CUBA: A Library for multidimensional numerical integration. *Comput.Phys.Commun.*, 168:78–95, 2005. doi: 10.1016/j.cpc.2005.01.010.
- [9] M. W. Perlin, M. M. Legler, C. E. Spencer, J. L. Smith, W. P. Allan, J. L. Belrose, and B. W. Duceman. Validating trueallele?? dna mixture interpretation. *Journal of Forensic Sciences*, 56:1430–1447, 2011.
- [10] R. Puch-Solis. A dropin peak height model. *Forensic Sci. Int. Genet.*, (11):80–84, 2014.
- [11] R. Puch-Solis, L. Rodgers, A. Mazumder, S. Pope, I. W. Evett, J. Curran, and D. Balding. Evaluating forensic dna profiles using peak heights, allowing for multiple donors, allelic dropout and stutters. *Forensic Sci. Int. Genet.*, 2013.
- [12] D. Taylor, J. A. Bright, and J. Buckleton. The interpretation of single source and mixed dna profiles. *Forensic Science International: Genetics*, 7, 2013.

## A Appendix I: Derivings

### A.1 Expanded covariance structure of restricted mixture proportion

To get the full covariance structure of  $\boldsymbol{\theta} = (\mathbf{m}, \mu, \sigma, \xi)$  we need to calculate  $Cov[m_C, \boldsymbol{\theta}]$ . The reason is that our implemented likelihood function does not

depend on  $m_C$  directly. The formula for the covariance with the other parameters is given as

$$\text{Cov}[m_C, \boldsymbol{\theta}] = \text{Cov}\left[-\sum_{k=1}^{C-1} m_k, \boldsymbol{\theta}\right] = -\sum_{k=1}^{C-1} \text{Cov}[m_k, \boldsymbol{\theta}]$$

with the variance given as

$$\text{Var}[m_C] = \text{Var}\left[\sum_{k=1}^{C-1} m_k\right] = \sum_{k=1}^{C-1} \sum_{l=1}^{C-1} \text{Cov}[m_k, m_l]$$

## A.2 Reparameterization

There are advantages in transforming the restricted space to unrestricted space. For instance it could be a good idea to transform the simplex space  $\mathcal{M} = \text{simplex}(\mathbf{1}_C)$  (i.e.  $\mathbf{m} \in [0, 1]^C$  with restriction  $m_C = 1 - \sum_{k=1}^{C-1} m_k \geq 0$ ) into to domain  $\mathbb{R}^{C-1}$ . Optimization, Laplace approximation and MCMC performance will perhaps be better in an unrestricted space. We now present the transformation from the restricted parameter space of  $\boldsymbol{\theta} = (\mathbf{m}, \mu, \sigma, \xi) \in (\mathcal{M}, \mathbb{R}_+^2, [0, 1])$  to the unrestricted space  $\boldsymbol{\phi} = (\boldsymbol{\nu}, \tilde{\mu}, \tilde{\sigma}, \tilde{\xi}) \in \mathbb{R}^{(C+2)}$ .

For given contributors  $C$ ,  $\mathbf{m} = (m_1, \dots, m_C) \in \mathcal{M}$  are given as the mixture proportion with  $m_C = 1 - (m_1 + \dots + m_{C-1})$ . Transformations for given  $\mathbf{m} \in \mathcal{M}$  (constrained) to unconstrained domain  $\mathbb{R}^{C-1}$  is given as

$$\begin{aligned} \nu_1 &= \log\left(\frac{m_1}{1 - m_1}\right) \\ \nu_2 &= \log\left(\frac{m_2/(1 - m_1)}{1 - m_2/(1 - m_1)}\right) \\ &\vdots \\ \nu_{(C-1)} &= \log\left(\frac{m_{(C-1)}/(1 - \sum_{k=1}^{C-2} m_k)}{1 - m_{(C-1)}/(1 - \sum_{k=1}^{C-2} m_k)}\right) \end{aligned}$$

with  $\boldsymbol{\nu} = (\nu_1, \dots, \nu_{(C-1)}) \in \mathbb{R}^{(C-1)}$ . This transformation also takes care of  $\sum_{k=1}^{C-1} m_k \leq 1$

We also do transformations of  $(\mu, \sigma, \xi)$  such that  $\tilde{\sigma} = \log \sigma$ ,  $\tilde{\mu} = \log \mu$  and  $\tilde{\xi} = \log\left(\frac{\xi}{1-\xi}\right)$

### A.2.1 Inverse transformations

The inverse transformation of the variables  $\boldsymbol{\nu} \in \mathbb{R}^{(C-1)}$  back to  $\mathcal{M}$  is given as

$$\begin{aligned} m_1 &= (1 + e^{-v_1})^{-1} \\ m_2|m_1 &= (1 + e^{-v_2})^{-1}(1 - m_1(v_1)) \\ &\vdots \\ m_{C-1}|(m_1, \dots, m_{(C-2)}) &= (1 + e^{-v_{(C-1)}})^{-1} \left( 1 - \sum_{k=1}^{C-2} m_k(v_1, \dots, v_k) \right) \\ m_C|(m_1, \dots, m_{(C-1)}) &= 1 - \sum_{k=1}^{C-1} m_k(v_1, \dots, v_k) \end{aligned}$$

Also, the transformations back to  $(\mu, \sigma, \xi)$  is given as

$$\begin{aligned} \sigma &= e^{\tilde{\sigma}} \\ \mu &= e^{\tilde{\mu}} \\ \xi &= (1 - e^{-\tilde{\xi}})^{-1} \end{aligned}$$

### A.2.2 Jacobian matrix of the inverse transformation

For a given number of contributors  $C$ , the elements  $(i, j)$  in the  $(C+2) \times (C+2)$ -Jacobian matrix  $J(\boldsymbol{\phi})$  of the inverse transformation of the variables  $\boldsymbol{\phi}$  back to the variables  $\boldsymbol{\theta}$  is given as

$$\begin{aligned} J_{ij} &= \frac{\partial m_i(v_1, \dots, v_i)}{\partial v_j} = -(1 + e^{-v_i})^{-1} \sum_{l=j}^{i-1} \frac{\partial m_l(v_1, \dots, v_l)}{\partial v_j} & 0 < j < i < (C-1) \\ J_{ij} &= \frac{\partial m_i(v_1, \dots, v_i)}{\partial v_i} = e^{-v_i} (1 + e^{-v_i})^{-2} \left( 1 - \sum_{l=1}^{i-1} m_l(v_1, \dots, v_l) \right) & 0 < i = j < C \\ J_{ij} &= \frac{\partial m_i(v_1, \dots, v_i)}{\partial v_j} = 0 & 0 < i < j < C \\ J_{ij} &= \frac{\partial \mu(\tilde{\mu})}{\partial \tilde{\mu}} = e^{\tilde{\mu}} & i = j = C \\ J_{ij} &= \frac{\partial \sigma(\tilde{\sigma})}{\partial \tilde{\sigma}} = e^{\tilde{\sigma}} & i = j = C+1 \\ J_{ij} &= \frac{\partial \xi(\tilde{\xi})}{\partial \tilde{\xi}} = e^{-\tilde{\xi}} (1 + e^{-\tilde{\xi}})^{-2} & i = j = C+2 \\ J_{ij} &= 0 & \text{else} \end{aligned}$$

where else is the set  $\{\{i, j\} : (C-1) < j < i < (C+3) \text{ or } (C-1) < i < j < (C+3)\}$ .

### A.2.3 Delta method

Given the covariance matrix  $\tilde{\Sigma}(\hat{\boldsymbol{\phi}})$  of the maximum likelihood estimates of  $\boldsymbol{\phi}$  (optimized in the space  $\mathbb{R}^{(C+2)}$ ), the covariance matrix of the maximum likelihood estimates of  $\boldsymbol{\theta}$  is approximately  $\Sigma(\boldsymbol{\theta}^*) \approx J(\hat{\boldsymbol{\phi}})^T \tilde{\Sigma}(\hat{\boldsymbol{\phi}}) J(\hat{\boldsymbol{\phi}})$

### A.3 Error interval from relative error

Let  $y$  be a measure for the real  $x$ . With Relative Error given as  $\delta = \frac{|y-x|}{x}$ , the error interval for  $x$  becomes  $[\frac{y}{1+\delta}, \frac{y}{1-\delta}]$ . With the absolute error  $\Delta = |y - x|$  given, the error interval for  $x$  becomes  $[y - \Delta, y + \Delta]$ .

### A.4 Laplace approximation

We now show how to derive the laplace approximation based on a second order taylor expansion and use this in weight-of-evidence. Consider we have a likelihood function  $L(\theta)$  with  $l(\theta) = \log(L(\theta))$  to be the log-likelihood where we consider  $\theta$  to be a  $p$  large parameter vector. The integral we want to calculate can be written as

$$L(E) = \int L(\theta|E) d\theta = \int e^{l(\theta|E)} d\theta = e^{l(\theta_0|E)} \int e^{l(\theta|E) - l(\theta_0|E)} d\theta$$

With a second order taylor approximation around  $\theta_0$  we have that

$$l(\theta|E) \approx l(\theta_0|E) + \Delta l(\theta_0)(\theta - \theta_0) + 0.5(\theta - \theta_0)^T H(\theta_0)(\theta - \theta_0)$$

where  $\Delta l(\theta_0)$  is the partial derivative vector and  $H$  is the hessian matrix for  $l(\theta|E)$ , both evaluated at  $\theta = \theta_0$ .

Note that when  $\theta_0 = \hat{\theta}_{ML}$ ,  $\Delta l(\theta_0|E) = (0, \dots, 0)$  (i.e. a singular point). And hence it follows that

$$L(E) \approx e^{l(\hat{\theta}_{ML}|E)} \int e^{0.5(\theta - \hat{\theta}_{ML})^T H(\hat{\theta}_{ML})(\theta - \hat{\theta}_{ML})} d\theta$$

Further note that if  $H(\hat{\theta}_{ML})$  is positive definite (i.e.  $\hat{\theta}_{ML}$  is a positive maximum),

$$p(\theta) = (2\pi)^{-\frac{p}{2}} |H(\hat{\theta}_{ML})|^{-\frac{1}{2}} e^{-0.5(\theta - \hat{\theta}_{ML})^T (-H(\hat{\theta}_{ML}))(\theta - \hat{\theta}_{ML})}$$

is a multivariate normal density with mean  $\hat{\theta}_{ML}$  and  $-H(\hat{\theta}_{ML})^{-1}$ , and the negative inverted hessian as the covariance matrix. Hence

$$L(E) \approx e^{l(\hat{\theta}_{ML}|E)} (2\pi)^{\frac{p}{2}} |H(\hat{\theta}_{ML})|^{-\frac{1}{2}} = L(\hat{\theta}_{ML}|E) (2\pi)^{\frac{p}{2}} |H(\hat{\theta}_{ML})|^{-\frac{1}{2}}$$

Note that this approximation is only good if the likelihood function is symmetrical and number of samples to the paramters are large (large sample theory convergence).

Further assume that we want to evaluate LR between the hypothesis  $H_p$  and  $H_d$ , with same number of parameters in both models. Then we have that

$$LR(H_p, H_d|E) \approx \frac{L(\hat{\theta}_{ML}|E, H_p)}{L(\hat{\theta}_{ML}|E, H_d)} \sqrt{\frac{|H(\hat{\theta}_{ML}|H_d)|}{|H(\hat{\theta}_{ML}|H_p)|}} \approx \frac{L(\hat{\theta}_{ML}|E, H_p)}{L(\hat{\theta}_{ML}|E, H_d)}$$

where the last identity is true if  $H(\hat{\theta}_{ML}|H_d) \approx H(\hat{\theta}_{ML}|H_p)$  (i.e. the hessian matrix doesn't change very much for the two comparing hypotheses),



Note that if any hypothesis contains at least 2 unknowns, the likelihood function will be multimodal, but symmetrical. Because of this, the number of multimodalities needs to be taken into account in the approximated integral. Let  $U_p$  be number of unknown contributors under  $H_p$  and  $U_d$  be number of unknown contributors under  $H_d$ . Then

$$LR(H_p, H_d|E) \approx \frac{U_p!L(\hat{\theta}_{ML}|E, H_p)}{U_d!L(\hat{\theta}_{ML}|E, H_d)}$$