

AALBORG UNIVERSITY

DENMARK

On the exact distribution of the number of alleles in DNA mixtures

Torben Tvedebrink - tvede@math.aau.dk

Department of Mathematical Sciences, Aalborg University, Denmark

Section of Forensic Genetics, Department of Forensic Medicine, Faculty of Health and Medical Sciences, University of Copenhagen, Denmark



COPENHAGEN

Introduction

When more than one individual contributes biological material to a forensic stain, the resulting DNA type is termed a DNA mixture. DNA mixtures occur frequently in forensic genetic casework, and in recent years much research has been devoted to this subject.

This poster presents a derivation of the exact distribution of the number of alleles for any number of profiles and investigated loci. The

Computing $P(N_{I}(m) = n)$ based on α_{m}

In order to evaluate $P(N_l(m)=n)$, we must consider the set \mathcal{A}_m^n of all α_m -vectors of dimension n. Each of these vectors, $\alpha_m = (\alpha_1, \dots, \alpha_n)$, gives the relevant powers of allele frequencies, which represent the count of the specific allele present in the DNA mixture. For example, for $\alpha_2=(2,2)$ both alleles are represented twice, e.g. both mixtures $(A_i A_i, A_i A_i)$ and $(A_i A_i, A_i A_i)$ are represented by (2, 2).

The fact that more mixtures are represented by the same α_m -vector motivates the need for a weightfunction, $c(\alpha_m)$. The $c(\alpha_m)$ -function counts the number of times α_m is formed by the recursions (see the *Recursion*-box), e.g. $\alpha_2 = (2, 1, 1)$ is formed twice in Fig. 1. Furthermore, if the added profile is heterozygous of type h_1 or h_2 (sharing one or both alleles), the formed α_m -vector is counted twice (due to unordered observations). Hence, to obtain $P(N_l(m) = n)$, we evaluate

per locus number of observed alleles is of interest as it indicates the plausible range of the number of contributors. Furthermore, the total number of alleles across all loci are used by some forensic geneticists to estimate the probability that an allele has not been detected.

$P(N_{l}(m) = n) = \sum c(\alpha_{m}) \sum^{T} p_{i_{1}}^{\alpha_{1}} p_{i_{2}}^{\alpha_{2}} \dots p_{i_{n}}^{\alpha_{n}},$

where i_1, \ldots, i_n are *n* different indices, and p_{i_1}, \ldots, p_{i_n} are the allele frequencies of locus *I*.

Recursion – how to obtain α_m

The total number of alleles observed for *m* contributors, N(m), depends on the number of loci, L, through the locus specific allele counts, N_l , by $N(m) = \sum_{l=1}^{L} N_l(m)$. Hence, we focus on computing the distribution of $N_{l}(m)$ below.

The expression for $P(N_l(m) = n)$ gets complicated for increased *m*. Hence, let α_m denote a vector of the numbers of unique alleles that the *m* DNA profiles carry at a given locus. In the simplest case m = 1, the DNA profile may either be heterozygous, $\alpha_1 = (1, 1)$, or homozygous, $\alpha_1 = (2)$. In order to compute P(N(m) = n), the only information necessary is α_m .

Plots of some results

Fig. 2 shows the distributions of P(N(m) = n) for $m = 1, \dots, 8$ of the ten SGM Plus loci, L = 10.



Drop-out

Allelic drop-out will cause fewer alleles to be detected in a sample. This led Gill et al. [1] and more recently Haned et al. [2, 3] to use only the number of observed alleles to estimate the drop-out probability.

However, a relatively low number of observed alleles may also be caused by shared ancestry and subpopulation stratification (often modelled by the θ -correction). In Fig. 4, the estimated drop-out probabilities for a two-person DNA mixture based on Monte Carlo simulations [1] are plotted against θ .

Extending a set of profiles can happen in a limited number of ways. The added profile may be

- lacktriangleright heterozygous, sharing none (h_0) , one (h_1) or both (h_2) alleles with the previous profiles;
- ▶ homozygous, sharing none (H_0) or one (H_1) allele with the previous profiles.

Mathematically, this can be formulated in terms of updating α_m to obtain α_{m+1} , where I is the dimension of α_m :

 $\alpha_m + \boldsymbol{e}_i + \boldsymbol{e}_j, \quad 1 \leq i < j \leq l$ (h_2) $\alpha_m + \boldsymbol{e}_i + \boldsymbol{e}_j, \quad 1 \leq i \leq l, \ j = l+1$ (h_1) $\alpha_m + e_i + e_j, \quad i = l+1, \ j = l+2$ (h_0) $\alpha_{m+1} = \langle$ $\alpha_m + 2\boldsymbol{e}_i, \quad 1 \leq i \leq I$ (H_1) $\alpha_m + 2\boldsymbol{e}_i, \qquad i = l+1$ (H_0) Obtaining α_2 from α_1 is shown in Fig. 1. For example, in the topmost path (green), the first profile is homozygous, $\alpha_1 = (2)$, and adding a heterozygous profile that shares no allele, h_0 , we obtain $\alpha_2 = (2, 1, 1)$.



Fig. 4: Monte Carlo based estimate of P(D) plotted against θ for various numbers of observed alleles, *n*.

Fig. 4 shows that the estimated drop-out probability, $\hat{P}(D)$, is correlated with θ . The declining trend shows that the Monte Carlo approach will overestimate P(D) when $\theta > 0$. For e.g. n = 25, the estimates are $\hat{P}(D; \theta = 0.05) = 0.19$ and $\hat{P}(D; \theta = 0) = 0.25$, respectively.



Fig. 1: Recursion showing how to obtain α_2 from α_1 .

5 6 7 8 9 10 11 12 Number of observed alleles, n

Fig. 3: Distribution of the number of contributors, *m*, given that *n* alleles are observed, P(m|n), for vWA.

Conclusion

The exact distribution of the number of alleles was derived using recursion relations.

Based on $P(N_l(m) = n)$ and a prior distribution, P(m), the posterior of the number of contributors, P(m|n), can be computed.

It was demonstrated that the drop-out probability based on the number of alleles, $\hat{P}(D)$, was correlated with θ .

References

[1] Gill, P., A. Kirkham, and J. Curran (2007). LoComatioN: A software tool for the analysis of low copy number DNA profiles. Forensic Sci Int 166(2-3): 128-138.

Haned, H. (2011). Forensim: An open-[2] source initiative for the evaluation of statistical methods in forensic genetics. Forensic Sci Int Genet 5(4): 265-268.

[3] Haned, H., K. Slooten, and P. Gill (2012). Exploratory data analysis for the interpretation of low template DNA mixtures. Forensic Sci Int Genet 6(6): 762-774.