

Deconvolution of MPS STR DNA Mixtures

Søren B. Vilsen

2017-11-03

Contents

1 Introduction

1

1 Introduction

The central equation to statisticians in forensic genetics is the posterior odds of two competing hypotheses:

$$\frac{\mathbb{P}(\mathcal{H}_d|\mathcal{E})}{\mathbb{P}(\mathcal{H}_p|\mathcal{E})} = \frac{\mathbb{P}(\mathcal{E}|\mathcal{H}_d) \mathbb{P}(\mathcal{H}_d)}{\mathbb{P}(\mathcal{E}|\mathcal{H}_p) \mathbb{P}(\mathcal{H}_p)} \quad (1)$$

where \mathcal{E} is the evidence (quantitative and qualitative), \mathcal{H}_d , and \mathcal{H}_p are the hypotheses of the defence and prosecution, respectively, $\mathbb{P}(\mathcal{E}|\mathcal{H}_d)/\mathbb{P}(\mathcal{E}|\mathcal{H}_p)$ is the likelihood ratio of the evidence under the two hypotheses, and $\mathbb{P}(\mathcal{H}_d)/\mathbb{P}(\mathcal{H}_p)$ are the prior odds of the hypotheses. We cannot and should not infer the prior odds, they should be determined by the court, which leaves the likelihood ratio:

$$LR(\mathcal{H}_d, \mathcal{H}_p) = \frac{\mathbb{P}(\mathcal{E}|\mathcal{H}_d)}{\mathbb{P}(\mathcal{E}|\mathcal{H}_p)}. \quad (2)$$

Assuming that the evidence \mathcal{E} consists of quantitative information \mathcal{Q} (the coverage in MPS or peak height/area in CE) about genetic information g , then we can factorise $\mathbb{P}(\mathcal{E}|\mathcal{H}_i)$ as:

$$\mathbb{P}(\mathcal{E}|\mathcal{H}_i) = \mathbb{P}(\mathcal{Q}, g|\mathcal{H}_i) = \mathbb{P}(\mathcal{Q}|g) \mathbb{P}(g|\mathcal{H}_i). \quad (3)$$

Thus, the probability of the evidence, given a hypothesis, as two parts:

- (1) **The probability of the quantitative information given allelic information**, determined by the likelihood under a model.
- (2) **The probability of the allelic information under the hypothesis in question**, assuming either HWE or θ correction [?].

Let \mathcal{U}_p and \mathcal{U}_d be the set of unknown profiles under the prosecutors and defence hypotheses, respectively, then we can generally write the LR as:

$$LR(\mathcal{H}_d, \mathcal{H}_p) = \frac{\sum_{\mathbf{g}_{U_d} \in \mathcal{U}_d} L(\mathcal{Q}|\mathbf{g}_{K_d}, \mathbf{g}_{U_d}) \mathbb{P}(\mathbf{g}_{U_d}|\mathbf{g}_K)}{\sum_{\mathbf{g}_{U_p} \in \mathcal{U}_p} L(\mathcal{Q}|\mathbf{g}_{K_p}, \mathbf{g}_{U_p}) \mathbb{P}(\mathbf{g}_{U_p}|\mathbf{g}_K)}, \quad (4)$$

where $\mathbf{g}_K = \{\mathbf{g}_{K_d}, \mathbf{g}_{K_p}\}$.

The sums in the numerator and denominator of the general LR , seen in Eq. (4), are intractable. Even the sum of the numerator in Eq. (??) can be computationally intensive to evaluate. In order to get around this problem we will approximate the LR , using both simple step approximations and more complicated sampling methods, see Section ?? . We will sample from a distribution centred at the most probable combination of profiles. Therefore, we will find the most probable combination of profiles under the model introduced in Section ?? , given the number of known and unknown contributors. [?] finds this combination of profiles using a greedy algorithm on a database of possible profile combinations. We will find the optimal combination of profiles using an Evolutionary Algorithm described in Section ?? .

In the remainder of these notes, we will assume the following (unless otherwise stated):

- (1) We have sequenced a sample using an \mathcal{M} marker multiplex panel and denote an arbitrary locus by m .
- (2) We observe \mathcal{A}_m alleles at locus m and denote an allele (using MPS these are represented as sequenced DNA strings) by $a \in \{a_{m1}, \dots, a_{m\mathcal{A}_m}\}$. For ease of notation we will sometimes denote this as n_m .
- (3) The DNA sample is a mixture of \mathcal{C} DNA profiles, i.e. that the sample contains exactly \mathcal{C} contributors.
- (4) The data is correct, i.e. contains no drop-outs/ins or other artefacts.