

On finding the best matching pair of profiles in two-person DNA mixtures

Torben Tvedebrink - tvede@math.aau.dk

Ass	sumptions	Model	Algorithm	Implementation	Example	Likelihood ratio
С	outline					
	Assump	otions				
	Model					
	Algorith	nm				
	Implem	entation				
	Exampl	e				
	Likeliho	ood ratio				
						2/32
						2/ 32

On finding the best matching pair of profiles in two-person DNA mixtures

Torben Tvedebrink - tvede@math.aau.dk



Controlled experiments with pair-wise mixtures of four different profiles in varying ratios conducted at the Section of Forensic Genetics, Department of Forensic Medicine, Faculty of Health Sciences, University of Copenhagen.





On finding the best matching pair of profiles in two-person DNA mixtures

Torben Tvedebrink - tvede@math.aau.dk





On finding the best matching pair of profiles in two-person DNA mixtures

Torben Tvedebrink - tvede@math.aau.dk



From the data exploration we observe the following proportionalities:

- peak height and peak area.
- peak area and amount of DNA.
- mean and variance of peak areas.

We model the peak areas conditioned on the loci peak area sums, $\mathbf{A}_{+} = (A_{s,+})_{s} \in S$, where S is the set of all loci used for discrimination (included in the typing kit).

We assume conditional independence of the peak areas on different loci conditioned on the loci sums,

$$\mathbf{A}_{s} \perp \mathbf{A}_{t} | (A_{s,+}, A_{t+}).$$

We note that \mathbf{A}_+ is an ancillary statistic as it is fixed for all pairs of profiles and thus contains no information on the profiles.

Assumptions Model Algorithm Implementation Example Likelihood ratio

For each locus $s \in S$ we assume the peak areas follow a conditional normal distribution with the following mean and variance

$$\begin{split} \mathbb{E}(\mathbf{A}_{s}|A_{s,+}) &= \frac{A_{s,+}}{2} \left[\alpha \mathbf{P}_{s,1} + (1-\alpha) \mathbf{P}_{s,2} \right] \\ \mathbb{C}\mathsf{ov}(\mathbf{A}_{s}|A_{s,+}) &= \tau^{2} C_{s} \mathsf{diag}(\mathbf{h}_{s}) C_{s}^{\top} = \tau^{2} W_{s}, \end{split}$$

where $\mathbf{P}_{s,k}$ is an indicator vector with elements 0,1 or 2 indicating whether person k has none, one or two copies of the particular allele.

The α -parameter models the share of person 1 in the mixture and common for all loci.

Furthermore, C_s is given as $I_{n_s} - n_s^{-1} \mathbf{1}_{n_s} \mathbf{1}_{n_s}^{\top}$ where diag(\mathbf{h}_s) in the covariance ensures proportionality of mean and variance.



The estimators of α and τ^2 are found as

$$\hat{\alpha} = \frac{\sum_{s \in \mathcal{S}} \mathbf{x}_0^{s \top} W_s^{-} (\mathbf{A}_s - \mathbf{x}_1^s)}{\sum_{s \in \mathcal{S}} \mathbf{x}_0^{s \top} W_s^{-} \mathbf{x}_0^s}$$

$$\hat{\tau}^2 = N^{-1} \sum_{s \in \mathcal{S}} (\mathbf{A}_s - \hat{\alpha} \mathbf{x}_0^s - \mathbf{x}_1^s)^{\top} W_s^{-} (\mathbf{A}_s - \hat{\alpha} \mathbf{x}_0^s - \mathbf{x}_1^s),$$

where
$$\mathbf{x}_{0}^{s} = \frac{A_{s,+}}{2} (\mathbf{P}_{s,1} - \mathbf{P}_{s,2})$$
 and $\mathbf{x}_{1}^{s} = \frac{A_{s,+}}{2} \mathbf{P}_{s,2}$ and $N = \sum_{s \in S} (n_{s} - 1) - 1 = n_{+} - S - 2$.

The matrix W_s^- is the generalised inverse since W_s is not of full rank due to conditioning on the sum $A_{s,+}$.

Assumptions Model Algorithm Implementation Example Likelihood ratio

We assume that $DNA_1 < DNA_2$ such that $\alpha < 0.5$. Furthermore, we sort the observed peak areas $A_{s,(1)} < \cdots < A_{s,(n_s)}$.

This reduces the number of possible combinations for the best matching pair of profiles.

\mathcal{J}_1 :	$\boldsymbol{P}_1\boldsymbol{P}_2$		\mathcal{J}_2 :	$\boldsymbol{P}_1\boldsymbol{P}_2$	$\bm{P}_1\bm{P}_2$	$\mathbf{P}_1\mathbf{P}_2$	$\mathbf{P}_1\mathbf{P}_2$
$A_{s,(1)}$	2 2		$A_{s,(1)}$	1 1	2 0	1 0	0 1
			$A_{s,(2)}$	1 1	0 2	1 2	2 1
σ	D D	D D		D D		a	
\mathcal{J}_3 :	$\mathbf{P}_1\mathbf{P}_2$	$\mathbf{P}_1\mathbf{P}_2$	$\mathbf{P}_1\mathbf{P}_2$	$\mathbf{P}_1\mathbf{P}_2$		\mathcal{J}_4 :	$\mathbf{P}_1\mathbf{P}_2$
$A_{s,(1)}$	2 0	1 0	1 0	0 1		$A_{s,(1)}$	1 0
$A_{s,(2)}$	0 1	1 0	0 1	0 1		$A_{s,(2)}$	1 0
$A_{s,(3)}$	0 1	0 2	1 1	2 0		$A_{s,(3)}$	0 1
						$A_{s,(4)}$	0 1

Assume that for locus s we observe the following peak areas $\mathbf{A}_s = (100, 200, 300)$. Then there exists 12 possible combinations:

$A_{s,(i)}$	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	${\bm P}_2$	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	\mathbf{P}_2
100	1	0	1	1	1	0	1	1	1	0	1	0
200	1	1	1	0	1	0	0	1	0	1	0	2
300	0	1	0	1	0	2	1	0	1	1	1	0
	-	_	I —	_	I —	_	I —	_	I —			_ 1
$A_{s,(i)}$	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	P ₂	\mathbf{P}_1	P ₂
A _{s,(i)} 100	P ₁ 2	P ₂	P ₁	P ₂	P ₁ 0	P ₂ 1						
<i>A_{s,(i)}</i> 100 200	P ₁ 2 0	P ₂ 0 1	P ₁ 0 2	P ₂ 1 0	P ₁ 0 0	P ₂ 1 1	P ₁ 0 1	P ₂ 2 0	P ₁ 0 1	P ₂ 1 1	P ₁ 0 1	P ₂ 1 0

Assume that for locus s we observe the following peak areas $\mathbf{A}_s = (100, 200, 300)$. Then there exists 12 possible combinations:

$A_{s,(i)}$	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	${\bm P}_2$	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	${\bm P}_2$	\mathbf{P}_1	\mathbf{P}_2
100	1	0	1	1	1	0	1	1	1	0	1	0
200	1	1	1	0	1	0	0	1	0	1	0	2
300	0	1	0	1	0	2	1	0	1	1	1	0
	-											
$A_{s,(i)}$	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	\mathbf{P}_2	\mathbf{P}_1	P ₂	\mathbf{P}_1	P ₂
$A_{s,(i)}$ 100	P ₁ 2	P ₂	P ₁	P ₂	P ₁ 0	P ₂ 1						
A _{s,(i)} 100 200	P ₁ 2 0	P ₂ 0 1	P ₁ 0 2	P ₂ 1 0	P ₁ 0 0	P ₂ 1 1	P ₁ 0 1	P ₂ 2 0	P ₁ 0 1	P ₂ 1 1	P ₁ 0 1	P ₂ 1 0



For $s \in S_i = \{s : s \in S \text{ and } n_s = i\}$

Choose combination $j \in \mathcal{J}_i$ minimising $\hat{\tau}^2$

Set $\mathcal{T} = \{\mathcal{T} \setminus (s, \cdot)\} \cup (s, j)$ and compute $\hat{\alpha}$

Return $\hat{\alpha}$, $\hat{\tau}$ and \mathcal{T} .

		Algorithm		Likelihood ratio
	C . I		5.1	
Propertie	s of the	greedy al	Jorithm	

The success of the greedy algorithm depends on the presence of loci with four alleles.

The fact that the algorithm is greedy only ensures that it finds the local maximum in each iteration. However, for practical purposes it will in most cases also find the global maximum (this depends on the presence of loci with full information).



The greedy algorithm has been implemented as an online tool available at www.math.aau.dk/ \sim tvede/dna.

Features:

- Upload of csv-files
- Graphical evaluation of the fit
- \blacksquare Estimates of α and τ^2 together with trace
- Fixing a suspect profile
- Testing goodness of fit for each locus
- Generation of alternatives *close* to the best matching profiles

Assumptions		Example	Likelihood ratio
Data			

Blu	Blue flourescent dye band				Green flourescent dye band				Yellow flourescent dye band					
Locus	Alle	le	Height	Area	Locus	Alle	le	Height	Area	Locus	Alle	le 🧹	Height	Area
D3	15	•	1802	15410	D8	8	•	1284	10782	D19	13	•	1332	10534
D3	16	• •	1939	16282	D8	12	•	1232	10359	D19	14	0	416	3478
	14	~	710	6100	D8	13	0	903	7891	D19	15	0	504	3968
VVVA	14	0	712	6620	D8	16	0	638	5291	TUO	c		000	6720
VVVA	15	•	125	0020						THU	0	0 •	820	0739
vWA	16	0	626	5637	D21	29	•	1073	9454	TH0	8	•	668	5573
vWA	17	•	830	7362	D21	30	0	1469	12828	TH0	9	0	486	4004
D16	10	0	824	7910	D21	31	•	798	6992	FGA	19	0	490	4415
D16	11	•	1772	17231	D18	13	0	1247	12302	FGA	23	•	865	7968
D16	12	0	586	6101	D18	15	•	899	9104	FGA	24	•	527	5036
D2	17	0	434	4558	D18	17	•	726	7549					
 D2	19		612	6563										
D2	25	•	843	9257										
02	23	••	045	5251										

Data from a controlled two-person mixture experiment. The true profiles are marked by \circ (profile 1) and \bullet (profile 2).

Assumptions		Example	Likelihood ratio
Data			



17/32

On finding the best matching pair of profiles in two-person DNA mixtures

Torben Tvedebrink - tvede@math.aau.dk

Assumptions		Example	Likelihood ratio
Data			

Analysing the mixture using the online DNA mixture separator tool: http://www.math.aau.dk/~tvede/dna

On finding the best matching pair of profiles in two-person DNA mixtures

18/32

When estimating the weight of evidence in a forensic genetic setting it is common to use a likelihood ratio of two contradictory hypotheses, H_p and H_d .

Data from STR DNA analyses comprises both qualitative (genetic stain, \mathcal{G}) and quantitative (peak heights and areas, \mathcal{Q}) evidence. Hence, the evidence $\mathcal{E} = (\mathcal{G}, \mathcal{Q})$. Using the definition of conditional probability we have,

$$P(\mathcal{E}|H) = P(\mathcal{G}, \mathcal{Q}|H) = P(\mathcal{Q}|\mathcal{G}, H)P(\mathcal{G}|H).$$

Let $(G', G'') \equiv \mathcal{G}$ be the set of all pairs of profiles G' and G''consistent with \mathcal{G} . Let further $G' : (G', G'') \equiv \mathcal{G}$ be the set of all G' that together with G'' is consistent with \mathcal{G} .



The likelihood ratio can then be written as,

$$LR = \frac{P(\mathcal{E}|H_p)}{P(\mathcal{E}|H_d)}$$





The likelihood ratio can then be written as,

L

$$R = \frac{P(\mathcal{E}|H_p)}{P(\mathcal{E}|H_d)}$$
$$= \frac{\sum_{G_U:(G_U,G_S)\equiv\mathcal{G}} P(\mathcal{Q}|G_U,G_S)P(G_U)}{\sum_{(G_{U_1},G_{U_2})\equiv\mathcal{G}} P(\mathcal{Q}|G_{U_1},G_{U_2})P(G_{U_1},G_{U_2})}$$

On finding the best matching pair of profiles in two-person DNA mixtures

21/32



The likelihood ratio can then be written as,

$$LR = \frac{P(\mathcal{E}|H_{p})}{P(\mathcal{E}|H_{d})}$$

=
$$\frac{\sum_{G_{U}:(G_{U},G_{S})\equiv\mathcal{G}}P(\mathcal{Q}|G_{U},G_{S})P(G_{U})}{\sum_{(G_{U_{1}},G_{U_{2}})\equiv\mathcal{G}}P(\mathcal{Q}|G_{U_{1}},G_{U_{2}})P(G_{U_{1}},G_{U_{2}})}$$

=
$$\frac{\sum_{G_{U}:(G_{U},G_{S})\equiv\mathcal{G}}L(\mathbf{A}|G_{U},G_{S})P(G_{U})}{\sum_{(G_{U_{1}},G_{U_{2}})\equiv\mathcal{G}}L(\mathbf{A}|G_{U_{1}},G_{U_{2}})P(G_{U_{1}})P(G_{U_{2}})}$$

Likelihood ratio

Estimation of *LR* using simulation

- The number of combinations consistent with \mathcal{G} in the denominator of LR equals $1^{S_1}7^{S_2}12^{S_3}6^{S_4}$, where S_i is the number of loci with i observed alleles. For the example this yields 9,029,615,616 terms in the sum.
- A possible methodology for estimating the numerator and denominator of $LR = P(\mathcal{E}|H_p)/P(\mathcal{E}|H_d)$ is importance sampling.

Importance sampling is a method for estimating (often) a mean that is impossible or difficult to compute analytically. Importance sampling increases the probability of sampling the more important points of the sample space for a decrease in the estimate variance.

The denominator of LR may be interpreted as a mean,

$$P(\mathcal{E}|H_d) = \mathbb{E}(f(\mathcal{E}); P) = \sum_{G \equiv \mathcal{G}} L(\mathbf{A}|G)P(G),$$

where $G \equiv G$ is the set of all combinations consistent with G. However, it is not appropriate to sample G_i , i = 1, ..., M using the the allele probabilities, $P(G_i)$, as these does not take the quantitative evidence into account. Defining a proposal function, q, we may improve the sampling using importance sampling. Let $q(G) = \prod_{s \in S} q_s(G_s)$, where

$$q_s(G_s) = \frac{L(\mathbf{A}|G_s, \hat{G}_{-s})P(G_s)}{\sum_{i=1}^{N_s} L(\mathbf{A}|G_{si}, \hat{G}_{-s})P(G_{si})},$$

where (G_s, \hat{G}_{-s}) is the particular combination on locus s and fixing all loci $t \in S \setminus s$ on best match level, and N_s are the number of combinations for locus s.

I.e. $L(\mathbf{A}|G_s, \hat{G}_{-s})$ are 'marginal' likelihood values.

Assumptions Model Algorithm Implementation Example Likelihood ratio

Using q we may rewrite the expression $P(\mathcal{E}|H_d)$ as,

$$P(\mathcal{E}|H_d) = \sum_{G \equiv \mathcal{G}} \frac{L(\mathbf{A}|G)P(G)}{q(G)}q(G) = \sum_{G \equiv \mathcal{G}} L(\mathbf{A}|G)W(G)q(G),$$

which is just another mean, $\mathbb{E}(f(\mathcal{E})h(\mathcal{E});q)$. Note that the importance weights W(G) simplifies

$$W(G) = \frac{P(G)}{\prod_{s \in S} \frac{L(\mathbf{A}|G_s, \hat{G}_{-s})P(G_s)}{\sum_{i=1}^{N_s} L(\mathbf{A}|G_{si}, \hat{G}_{-s})P(G_{si})}}$$
$$= \frac{\prod_{s \in S} \sum_{i=1}^{N_s} L(\mathbf{A}|G_{si}, \hat{G}_{-s})P(G_{si})}{\prod_{s \in S} L(\mathbf{A}|G_s, \hat{G}_{-s})} = \frac{B}{\prod_{s \in S} L(\mathbf{A}|G_s, \hat{G}_{-s})}$$

Assumptions Model Algorithm Implementation Example Likelihood ratio

This ensures that

$$L(\mathbf{A}|G)W(G) = \frac{L(\mathbf{A}|G)B}{\prod_{s\in\mathcal{S}}L(\mathbf{A}|G_s, \hat{G}_{-s})}$$

is rather constant which induces lower variance of the estimator.

Hence, in order to estimate $P(\mathcal{E}|H_d)$ we need to draw combinations $G_i \sim q(\cdot), i = 1, ..., M$, and compute the estimate

$$\hat{P}(\mathcal{E}|H_d) = \frac{1}{M} \sum_{i=1}^{M} L(\mathbf{A}|G_i) W(G_i).$$

Similar arguments apply for $P(\mathcal{E}|H_p)$ where we use $\hat{G}^{(S)}$ in place of \hat{G} for the restricted set $\{G : G \equiv \mathcal{G}, G_S \subseteq G\}$.

Assumptions Model Algorithm Implementation Example Likelihood ratio Example - Important sampling

In order to demonstrate the applicability of importance sampling we computed the exact value of $P(\mathcal{E}|H_d)$ for the blue loci (D3, vWA, D16, D2). This reduced the exhaustive list of combinations to $7^112^26^1 = 6048$ and yielded $P(\mathcal{E}|H_d) = 4.81335 \times 10^{-11}$ when we assume an uniform distribution of the allele probabilities.



On finding the best matching pair of profiles in two-person DNA mixtures

Assumptions Model Algorithm Implementation Example Likelihood ratio

We computed 1000 estimates based on 10,000 samples each. The estimates and true value are plotted in the histogram.



On finding the best matching pair of profiles in two-person DNA mixtures

Assumptions Model Algorithm Implementation Example Likelihood ratio Work in progress - Forward sampling

An alternative sampling strategy is based on forward sampling.

1. Sample a list of the loci in random order, $\mathcal{S}^* \sim \pi_\mathcal{S}(\cdot)$

2. Sample
$$s^*_{(1)} \sim q\left(\hat{G}_{-s_{(1)}}\right)$$

• Sample
$$s^*_{(2)}|s^*_{(1)} \sim q\left(\hat{G}_{-\{s_{(1)},s_{(2)}\}},s^*_{(1)}\right)$$

• : • Sample $s_{(S)}^* | s_{(1)}^*, \dots, s_{(S-1)}^* \sim q\left(s_{(1)}^*, \dots, s_{(S-1)}^*\right)$

This method allows for combinations further away from the best matching combination to be sampled compared to the importance sampling scheme.

Assumptions			Likelihood ratio
Conclusior			

- The method is based on a statistical model for the peak areas.
- The algorithmic implementation automates the work of the forensic scientist.
- Parameter estimates indicates the goodness-of-fit.
- The free online implementation finds the best matching pair of profiles in a few seconds.
- Fast graphical assessment of the validity of H_p by fixing a suspect profile.
- The model and greedy algorithm are easily extended for more than two contributors.
- The model features the possibility of using sampling techniques for assessing the likelihood ratio.



This is work in collaboration with

 Poul Svante Eriksen Associate Professor, MSc Department of Mathematical Sciences, Aalborg University

Helle Smidt Mogensen
 Forensic Geneticist, MSc PhD
 Section of Forensic Genetics, Department of Forensic Medicine,
 Faculty of Health Sciences, University of Copenhagen

 Niels Morling Professor, MD DMSC Section of Forensic Genetics, Department of Forensic Medicine, Faculty of Health Sciences, University of Copenhagen