POPULATION STRUCTURE

AlleleMatching

Allele Matching

Forensic genetics is concerned with matching of genetic profiles from evidence and from persons of interest. Profile match probabilities rest on the probabilities of matching among the alleles constituting the profiles.

Allele matching can refer to alleles within an individual (inbreeding), between individuals within a population (relatedness) and between populations (population structure). In all these cases there are parameters that describe profile match probabilities, and these parameters can be estimated by comparing the observed matching for a target set of alleles with that between a comparison set.

Allele Matching Within Individuals

The inbreeding coefficient for an individual is the probability it receives two alleles at a locus, one from each parent, that are *identical by descent*.

What can be observed, however, is identity in state. An individual is either homozygous or heterozygous at a locus: the two alleles either match or miss-match at that locus. The proportion of matching alleles at a locus is either zero or one, not a very informative statistic, but the proportion of an individual's loci that are homozygous may be informative for their inbreeding status.

There is still a need for a reference: for a locus such as a SNP with a small number of alleles many loci will be homozygous even for non-inbred individuals. Therefore we compare the proportion of loci with matching alleles for an individual with the matching proportion for pairs of alleles taken one from each of two individuals: is allele matching higher within than between individuals?

AlleleMatching

Inbreeding

If \tilde{M}_j is the observed proportion of loci with matching alleles (i.e. homozygous) for individual j, and if \tilde{M}_S is the observed proportion of matching alleles, one from each of two individuals in the population, then the within-population inbreeding coefficient f_j is estimated as

$$\widehat{f}_j = rac{\widetilde{M}_j - \widetilde{M}_S}{1 - \widetilde{M}_S}$$

Note that this can be negative for individuals with high degrees of heterozygosity.

The average of these estimates over all the individuals in a sample from a population estimates the within-population inbreeding coefficient f:

$$\widehat{f} = \frac{\widetilde{M}_I - \widetilde{M}_S}{1 - \widetilde{M}_S}$$

where $\tilde{M}_I = \sum_{j=1}^n \tilde{M}_j/n$. Hardy-Weinberg equilibrium corresponds to f = 0.

AlleleMatching

Slide 4

SNP-based Inbreeding

From 400,000 SNPs on Chromosome 22 of the 1000 Genomes ACB populations (96 Afro-Caribbeans in Barbados);



Allele Matching Between Individuals

How can we tell if a pair of individuals has a high degree of allele matching? What does "high" mean?

We assess relatedness of individuals within a population by comparing their degree of allele matching with the average degree for all pairs of individuals in that population.

Allele Matching Between Individuals

If $\tilde{M}_{jj'}$ is the observed proportion of loci with matching alleles, one from each of individuals j and j', and if \tilde{M}_S is the average of all the $\tilde{M}_{jj'}$'s, then the within-population kinship coefficient $beta_{jj'}$ is estimated as

$$\widehat{eta}_{jj'} = rac{\widetilde{M}_{jj'} - \widetilde{M}_S}{1 - \widetilde{M}_S}$$

Note that this can be negative for pairs of individuals less related than the average pair-matching in the sample.

The average of these estimates over all pairs of individuals in a sample is zero, but this doesn't allow us to compare populations.

SNP-based Coancestry

From 400,000 SNPs on Chromosome 22 of the 1000 Genomes ACB populations (4560 pairs of Afro-Caribbeans in Barbados);



Allele Matching Between Populations

We calibrated allele matching within individuals by comparison with matching between pairs of individuals.

We calibrate the allele matching between pairs of individuals by comparison with matching between pairs of populations. If $\tilde{M}^{ii'}$ is the observed proportion of loci with matching alleles, one from each of populations *i* and *i'*, and if \tilde{M}_B is the average of all the $\tilde{M}^{ii'}$'s, then the total kinship coefficient $\beta_{jj'}$ is estimated as

$$\hat{\beta}_{jj'} = \frac{\tilde{M}_{jj'} - \tilde{M}_B}{1 - \tilde{M}_B}$$

The average of these estimates over all pairs of individuals in a sample from a population is

$$\hat{\beta} = \frac{\tilde{M}_S - \tilde{M}_B}{1 - \tilde{M}_B}$$

This is the " θ " needed for the "theta correction" discussed below.

AlleleMatching

Slide 9

Within-population Matching

We can get some empirical matching proportions when we have a set of profiles. To simplify this initial discussion, consider the following data for the Y-STR locus DYS390 from the NIST database:

	Population				
Allele	Afr.Am.	Cauc.	Hisp.	Asian	Total
20	4	1	1	0	6
21	176	4	17	1	198
22	43	45	14	17	119
23	36	116	50	17	219
24	56	145	129	21	351
25	23	46	21	36	126
26	3	2	2	4	11
27	0	0	2	0	2
Total	341	359	236	96	1032

Within- and Between-population Matching for DYS390

Within the African-American sample there are $341 \times 340 = 115,940$ pairs of profiles and the number of between individual-pair matches is

 $4 \times 3 + 176 \times 175 + 43 \times 42 + 36 \times 35 + 56 \times 55 + 23 \times 22 + 3 \times 2 = 37,470$

so the within-population matching proportion is 37,470/115,940 = 0.323.

Between the African-American and Caucasian samples, there are $341 \times 359 = 122,419$ pairs of profiles and the number of matches is

 $4 \times 1 + 176 \times 4 + 43 \times 45 + 36 \times 116 + 56 \times 145 + 23 \times 4 + 3 \times 2 = 12,403$ so the between-population matching proportion is 12,403/122,419 = 0.101.

Allele Counts in NIST Data for DYS391

	Population				
Allele	Afr.Am.	Cauc.	Hisp.	Asian	Total
7	0	0	1	0	1
8	0	1	0	1	2
9	2	12	16	3	33
10	238	162	128	79	607
11	93	175	89	13	370
12	7	9	2	0	18
13	1	0	0	0	1
Total	341	359	236	96	1032

The within-population matching proportion for the African-American sample is 65,006/115,940=0.561.

The between-population matching proportion for the African-American and Caucasian samples is 54,918/122,419=0.449.

AlleleMatching

Two-locus counts in NIST African-American Data for DYS390, DYS391

DYS390	DYS391	Count n_g	$n_g(n_g-1)$
22	10	34	1122
22	11	9	72
24	10	15	210
24	11	39	1482
24	12	1	0
24	9	1	0
23	10	19	342
23	11	14	182
23	12	3	6
21	10	157	24492
21	11	15	210
21	12	2	2
21	9	1	0
21	13	1	0
25	10	11	110
25	11	12	132
26	10	1	0
26	11	2	2
20	10	1	0
20	11	2	2
20	12	1	0

Two-locus counts in NIST Caucasian Data for DYS390, DYS391

DYS390	DYS391	Count n_g	$n_g(n_g-1)$
22	10	43	1806
22	11	1	0
22	9	1	0
24	10	48	2256
24	11	88	7656
24	12	4	12
24	9	5	20
23	10	50	2450
23	11	60	3540
23	12	2	2
23	9	3	6
23	8	1	0
21	10	3	6
21	11	1	0
25	10	18	306
25	11	22	462
25	12	3	6
25	9	3	6
26	11	2	2
20	11	1	0

Two-locus Matches

The within-population matching proportion for the African-American sample is 28,366/115,940=0.245.

The within-population matching proportion for the Caucasian sample is 18,536/128,522=0.144.

The between-population matching proportion for the African-American and Caucasian samples is 8,347/122,419=0.068.

There is a clear decrease in matching between populations from within populations.

Will match probabilities keep decreasing?

Panel (number of STR loci)	Unrelated		Parent/child
	$Fst = 0^2$	Fst = 0.01	Fst = 0
New FBI core (24) ³	6.28 × 10 ⁻³⁰	5.12 × 10 ⁻²⁹	3.63 × 10 ⁻¹⁸
New FBI core section A $(20)^3$	9.54 × 10 ⁻²⁵	4.77×10^{-24}	3.83 × 10 ⁻¹⁵
13-loci CODIS core (13)	2.34×10^{-15}	5.83 × 10 ⁻¹⁵	1.74 × 10 ⁻⁹
ldentifiler (15)	5.93 × 10 ⁻¹⁸	1.73 × 10 ⁻¹⁷	5.04×10^{-11}
PowerPlex16 (15)	2.43×10^{-18}	7.48×10^{-18}	3.06×10^{-11}
NGM ⁴ (15)	1.12×10^{-19}	4.15 × 10 ⁻¹⁹	5.68×10^{-12}

Table 2 The expected match probability (EMP) of the kits/panels.¹

¹Caucasian population data were used.

Ge et al. 2012. Investigative Genetics 3:1-14.

Will match probabilities keep decreasing?

How do these match probabilities address the observation of Donnelly:

"after the observation of matches at some loci, it is relatively much more likely that the individuals involved are related (precisely because matches between unrelated individuals are unusual) in which case matches observed at subsequent loci will be less surprising. That is, knowledge of matches at some loci will increase the chances of matches at subsequent loci, in contrast to the independence assumption."

Donnelly P. 1995. Heredity 75:26-64.

Are match probabilities independent over loci?

Is the problem that we keep on multiplying match probabilities over loci under the assumption they are independent? Can we even test that assumption for 10 or more loci?

Or is our standard "random match probability" not the appropriate statistic to be reporting in casework? Is it actually appropriate to report statements such as

The approximate incidence of this profile is 1 in 810 quintillion Caucasians, 1 in 4.9 sextillion African Americans and 1 in 410 quadrillion Hispanics.

Putting "match" back in "match probability"

Let's reserve "match" for a statement we make about two profiles and take "match probability" to mean the probability that *two profiles match*. This requires calculations about *pairs of profiles*.

If the source of an evidence profile is unknown (e.g. is not the person of interest), then the match probability is the probability this unknown person has the profile *already seen in the POI*. No two profiles are truly independent, and their dependence affects match probabilities across loci.

Likelihood ratios use match probabilities

As with many other issues on forensic genetics, the issue of multilocus match probability dependencies is best addressed by comparing the probabilities of the evidence under alternative propositions:

 H_p : the person of interest is the source of the evidence DNA profile.

 H_d : an unknown person is the source of the evidence DNA profile.

Write the profiles of the POI and the source of the evidence as G_s and G_c . The evidence is the pair of profiles G_c , G_c .

Likelihood ratios use match probabilities

The likelihood ratio is

$$-R = \frac{\Pr(E|H_p)}{\Pr(E|H_d)}$$
$$= \frac{\Pr(G_c, G_s|H_p)}{\Pr(G_c, G_s|H_d)}$$
$$= \frac{1}{\Pr(G_c|G_s, H_d)}$$
$$= \frac{1}{\operatorname{Match probability}}$$

providing $G_c = G_s$ under H_p . The match probability is the chance an unknown person has the evidence profile given that the POI has the profile: this is not the profile probability.