

DNA Identification: Quantitative Data Modeling

Mark W Perlin, PhD, MD, PhD
Cybergenetics, Pittsburgh, PA

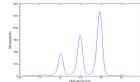
TrueAllele® Lectures
Fall, 2010



Cybergenetics © 2003-2010

Quantitative Data

- PCR is a linear process
- peak heights reflect the underlying DNA quantity
- use quantitative peak heights to explain the observed data



Genotype Model: 1 + 1 = 2

victim genotype 

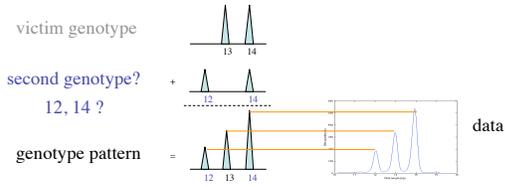
second genotype?  + 

12, 14 ?

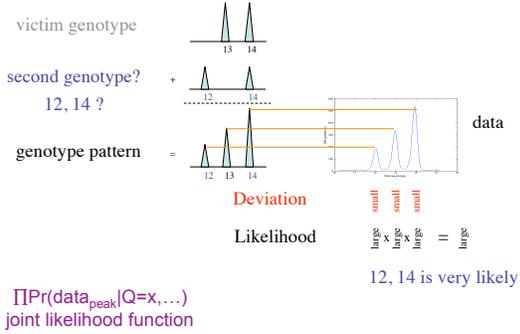
genotype pattern = 

Consider all possible allele pair values by trying out each candidate

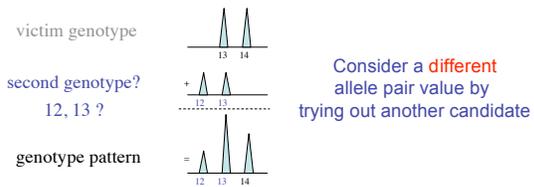
Compare Model to Data



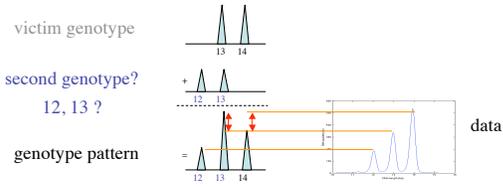
Likelihood Function



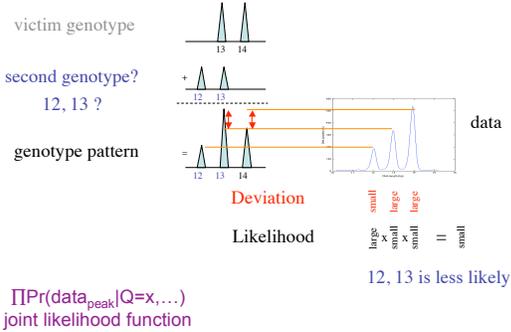
Genotype: Alternative Value



Compare Model to Data



Likelihood Function



All Genotype Possibilities

allele pair	likelihood	probability
12, 12	small	5%
12, 13	small	5%
12, 14	large	90%
13, 13	very small	0
13, 14	very small	0
14, 14	very small	0
12, 15	0	0
13, 16	0	0
sum =		100%

prior \times likelihood \Rightarrow posterior

Genotype inference

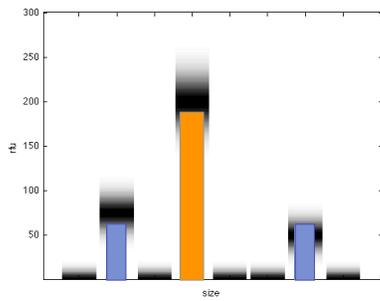
$$\Pr(Q=x|\text{data}, \dots) \propto \prod \Pr(\text{data}|Q=x, \dots) \times \Pr(Q=x)$$

posterior probability joint likelihood function prior probability

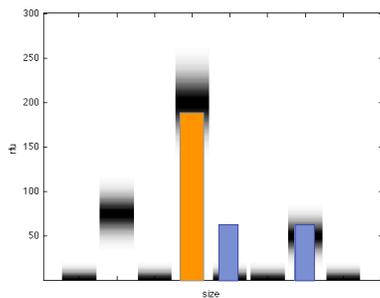
Try out all value possibilities;
better fit's more likely it.

$\prod \Pr(\text{data}_{\text{locus}}|Q=x, \dots)$
joint likelihood function

Genotype probability with data uncertainty



Genotype alternative value



Bayesian probability

- Assess ALL genotype patterns to find the probability of each allele pair.
- Similarly compute the data variance.
- Small data variation is **RESTRICTIVE**: only few genotype values are possible. (more certain)
- Large data variation is **PERMISSIVE**: many genotype values are possible. (less certain)

Likelihood ratio match statistic reflects genotype uncertainty

$$LR = \frac{\Pr(Q=s|data)}{\Pr(Q=s)}$$

Genotype **certainty concentrates** probability on just a few good bets, and focuses LR. (more info)

Genotype **uncertainty diffuses** probability across many candidates, and reduces LR. (less info)

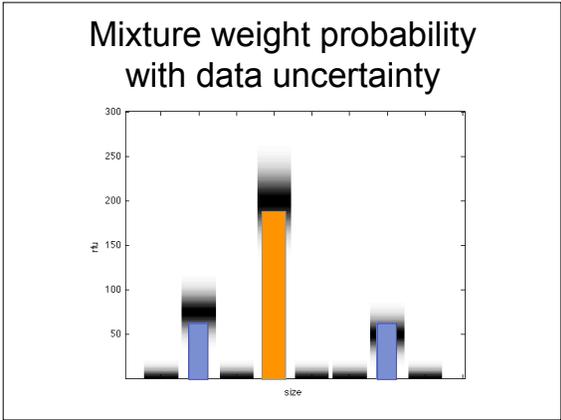
Mixture weight inference

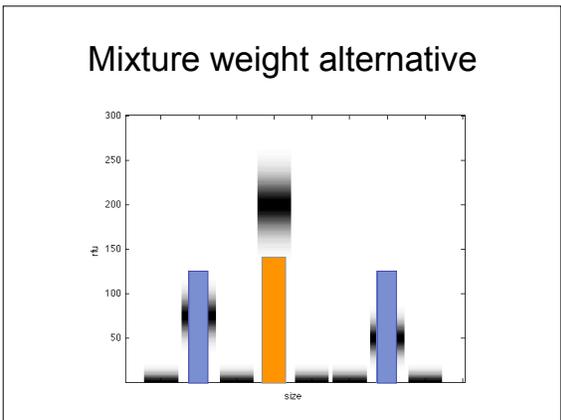
$$\Pr(W=w|data, \dots) \propto \prod \Pr(data|W=w, \dots) \times \Pr(W=w)$$

posterior probability joint likelihood function prior probability

Try out all value possibilities;
better fit's more likely it.

$\prod \Pr(data_{locus} | W=w, \dots)$
joint likelihood function





Data variance inference

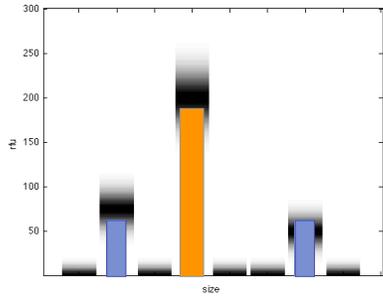
$$\Pr(V=v|\text{data}, \dots) \propto \prod \Pr(\text{data}|V=v, \dots) \times \Pr(V=v)$$

posterior probability joint likelihood function prior probability

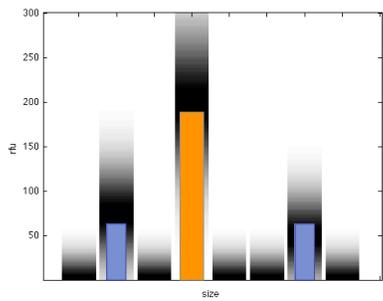
Try out all value possibilities;
better fit's more likely it.

$\prod \Pr(\text{data}_{\text{peak}}|V=v, \dots)$
joint likelihood function

Data variance probability of data peak uncertainty



Data variance alternative



Quantitative data modeling

- genotype is main variable of interest
- genotype gives identification LR
- mixture weight is explanatory variable
- data variance, stochastic effects
- identification information preserved by quantitative modeling
