

# Basic Probability

Probability: Quantifying our belief in the possible outcomes of future event

Event 1: First roll of 6-sided die (DR1)

$$P(DR1 = 2) = \frac{1}{6} \quad P(DR1 = 5) = \frac{1}{6}$$

Independent Events: The outcome of one event doesn't influence outcome of another

Event 2: Second roll of 6-sided die (DR2)

$$P(DR2 = 3 \mid DR1 = 5) = P(DR2 = 3) = \frac{1}{6}$$

"given"

Joint probability for independent events

$$P(DR1 = 5 \text{ AND } DR2 = 2) = P(DR1 = 5) \cdot P(DR2 = 2) = \frac{1}{6} \cdot \frac{1}{6} = \frac{1}{36}$$

Mutually Exclusive Outcomes: when one outcome happens the other can't happen at same time

$$P(DR1 = 5 \text{ OR } DR1 = 4) = P(DR1 = 5) + P(DR1 = 4) = \frac{1}{6} + \frac{1}{6} = \frac{2}{6} = \frac{1}{3}$$

# Maximum Likelihood Parameter Estimation

Parameter: A true (as opposed to estimated) measure of an individual or population

Maximum Likelihood: Method for estimating unknown parameter from data & statistical Model

ML Steps ① Determine mathematical relationship between parameter being estimated and data being observed (build model)

② Calculate Probability of observing data assuming a particular parameter value

③ Repeat step #2 for each possible parameter value

④ Best parameter estimate is value that maximized the probability of data

$$\text{Likelihood (Parameter = } x) = P(\text{Data} \mid \text{Parameter} = x)$$

# Calculating Genotype Likelihoods (1 read)

Parameter: Genotype of individual 01 @  
chr 16 pos 18,716 (C/G SNP)

Data: 1 C read, sequencing error rate = 3%

Possible Parameter Values:  $G=CC$ ,  $G=CG$ ,  $G=GG$

$$\text{Likelihood}(G=CC) = P(\text{Data} | G=CC)$$

$$\text{Likelihood}(G=CG) = P(\text{Data} | G=CG)$$

$$\text{Likelihood}(G=GG) = P(\text{Data} | G=GG)$$

$$P(\text{C read} | G=CG) = 0.49$$

$$\begin{aligned} &\text{sample C allele \& no error} \quad \text{OR} \quad \text{sample G allele \&} \\ &\quad (0.5 \cdot 0.97) \quad + \quad (0.5 \cdot 0.03(\frac{1}{3})) \end{aligned}$$

$$P(\text{C read} | G=CC) = 0.97$$

$$\begin{aligned} &\text{sample C allele \& no error} \\ &\quad (1 \cdot 0.97) \end{aligned}$$

$$\begin{aligned} P(\text{C read} | G=GG) &= 0.01 \quad G \rightarrow C \\ &\text{sample G allele \& } \text{error} \\ &\quad (1 \cdot 0.03(\frac{1}{3})) \end{aligned}$$



## Calculating Genotype Likelihoods (Multiple reads)

Data: 2 c reads, sequencing error rate = 3%

$$P(\text{C, C reads} \mid G = \text{CC}) = 0.97 \cdot 0.97 = 0.94$$

read 1 C & read 2 C

$$P(\text{C read} \mid G = \text{CC}) \cdot P(\text{C read} \mid G = \text{CC})$$

$$P(\text{C, C reads} \mid G = \text{CG}) = 0.49 \cdot 0.49 = 0.24$$

$$P(\text{C read} \mid G = \text{CG}) \cdot P(\text{C read} \mid G = \text{CG})$$

$$P(\text{C, C reads} \mid G = \text{GG}) = P(\text{C read} \mid G = \text{GG})^2 = 0.01^2$$

$$\boxed{\text{Most likely Genotype} = \text{CC}} = 0.0001$$

### More Data

$$P(\text{C, C, G, C reads} \mid G = \text{CC}) =$$

$$P(\text{C read} \mid G = \text{CC}) \cdot P(\text{C read} \mid G = \text{CC}) \cdot P(\text{G read} \mid G = \text{CC}) \cdot$$

$$P(\text{C read} \mid G = \text{CC}) =$$

$$0.97 \cdot 0.97 \cdot 0.01 \cdot 0.97 = 0.0091$$

# Bayesian Genotype Calling

Likelihoods:  
 $P(\text{Data} | \text{Parameter})$

$$P(c, c \text{ reads} | G = cc) = 0.94$$
$$P(c, c \text{ reads} | G = cG) = 0.24$$
$$P(c, c \text{ reads} | G = GG) = 0.0001$$

Don't sum to 1 - not Intuitive

$P(\text{Parameter} | \text{Data})$  - <sup>More</sup> Intuitive

$$P(G = cc | c, c \text{ reads}) = \frac{0.94}{0.94 + 0.24 + 0.0001} = 0.797$$

$$P(G = cG | c, c \text{ reads}) = \frac{0.24}{0.94 + 0.24 + 0.0001} = 0.203$$

$$P(G = GG | c, c \text{ reads}) = \frac{0.0001}{0.94 + 0.24 + 0.0001} = 0$$

Much more Satisfying ~~is~~ Legit?

Individual came from Population X

$$f_x(c) = 0.1 \quad f_x(G) = 0.9$$

$$f_x(cc) = \cancel{0.01} \quad 0.1^2 = 0.01$$
$$f_x(cG) = 2(0.1)(0.9) = 0.18$$
$$f_x(GG) = 0.9^2 = 0.81$$

# Bayes' Theorem

$$P(P_k | D) = \frac{P(D | P_k) P(P_k)}{\sum_{i=1}^n P(D | P_i) P(P_i)}$$

Posterior Probabilities  $\nearrow$   
 Likelihoods  $\swarrow$   
 Priors  $\nwarrow$

allele freq & HWE

Priors (~~allele freq~~)

Likelihoods

$$P(D | G=CC) = 0.94$$

$$P(D | G=CG) = 0.24$$

$$P(D | G=GG) = 0.0001$$

Priors (uniform)

$$P(G=CC) = \frac{1}{3}$$

$$P(G=CG) = \frac{1}{3}$$

$$P(G=GG) = \frac{1}{3}$$

Priors (allele freq)

$$P(G=CC) = 0.01$$

$$P(G=CG) = 0.18$$

$$P(G=GG) = 0.81$$

$$P(G=CC | D) = \frac{0.94 (\frac{1}{3})}{0.94 (\frac{1}{3}) + 0.24 (\frac{1}{3}) + 0.0001 (\frac{1}{3})} = 0.797$$

$$P(G=CG | D) = 0.203$$

$$P(G=GG | D) = 0$$

$$P(G=CC | D) = \frac{0.94 (0.01)}{0.94 (0.01) + 0.24 (0.18) + 0.0001 (0.81)} = 0.178$$

$$P(G=CG | D) = \frac{0.24 (0.18)}{0.94 (0.01) + 0.24 (0.18) + 0.0001 (0.81)} = 0.820$$

$$P(G=GG | D) = \frac{0.0001 (0.81)}{0.94 (0.01) + 0.24 (0.18) + 0.0001 (0.81)} = 0.002$$