

Multiallelic calling model in bcftools (-m)

Petr Danecek, Richard Durbin

Version: December 5, 2013

Let f_A be an estimate of the allele frequency of the base A across all samples

$$f_A = \frac{\sum_k q_k \delta_{b_k, A}}{\sum_k q_k},$$

where $b_k \in \{A, C, G, T\}$ is the base at k -th read and q_k the corresponding quality¹ and

$$f_{A|AC} = \frac{f_A}{f_A + f_C}$$

$$f_{A|ACG} = \frac{f_A}{f_A + f_C + f_G}.$$

Calculate likelihoods of all possible combinations of alleles² across all samples i as

$$L_A = \prod_i L_A^i$$

$$L_{AC} = \prod_i L_{AC}^i$$

$$L_{ACG} = \prod_i L_{ACG}^i$$

$$L_{ACGT} = \prod_i L_{ACGT}^i$$

where

$$L_A^i = PL^i(AA) = P(data|AA)$$

$$L_{AC}^i = f_{A|AC}^2 PL^i(AA) + f_{C|AC}^2 PL^i(CC) + 2f_{A|AC} f_{C|AC} PL^i(AC)$$

$$L_{ACG}^i = f_{A|ACG}^2 PL^i(AA) + f_{C|ACG}^2 PL^i(CC) + f_{G|ACG}^2 PL^i(GG)$$

$$+ 2(f_{A|ACG} f_{C|ACG} PL^i(AC) + f_{A|ACG} f_{G|ACG} PL^i(AG) + f_{C|ACG} f_{G|ACG} PL^i(CG))$$

$$L_{a_1, \dots, a_n}^i = \sum_j f_{a_j|a_1, \dots, a_n}^2 PL^i(a_j a_j) + 2 \sum_{j < k} f_{a_j|a_1, \dots, a_n} f_{a_k|a_1, \dots, a_n} PL^i(a_j a_k).$$

¹See `bcf_call_glfgen` in `bam2bcf.c` for calculation of q_k .

²In the current implementation, at most tri-allelic sites are considered.

Select the most likely set of alleles $\{a\}$ with likelihood $L_{\{a\}}$ and the second most likely set of alleles $\{b\}$ with likelihood $L_{\{b\}}$. Accept $\{a\}$ either if the number of alleles in $\{a\}$ is smaller than in $\{b\}$ or if the significance level for a 1 degree of freedom likelihood ratio test

$$\chi^2 = 2 \log \frac{L_{\{b\}}}{L_{\{a\}}}$$

exceeds a given threshold. Assuming HWE, for i -th sample we select the genotype $\{a\} = a_1 a_2$ which maximizes the likelihood

$$P_{\{a\}}^i = (2 - \delta_{a_1, a_2}) f_{a_1|\{a\}} f_{a_2|\{a\}} P L^i(a_1 a_2).$$

The corresponding genotype quality is

$$\text{GQ}_{\{a\}}^i = -10 \log \left[1 - \frac{P_{\{a\}}^i}{\sum_{\{b\}} P_{\{b\}}^i} \right]$$

and the call quality is

$$\text{QUAL} = -10 \log \frac{L_{\{ref\}}}{\sum_{\{b\}} L_{\{b\}}}$$

for variant calls and

$$\text{QUAL} = -10 \log \left[1 - \frac{L_{\{ref\}}}{\sum_{\{b\}} L_{\{b\}}} \right]$$

for non-variant calls, where the sum iterates over all possible genotypes $\{b\}$.