# Multiallelic calling model in bcftools (-m)

Petr Danecek, Stephan Schiffels, Richard Durbin
Version: June 23, 2014

Let $x$ and $y$ denote alleles. For simplicity of notation we work with SNPs, $x, y \in \{A, C, G, T\}$, but the method is identical for indels. Let's denote the number of samples $N$. In the pileup we observe the set of bases $S \subseteq \{A, C, G, T\}$, each base $x$ is observed $D_x$ times with the qualities $Q_1^x, \ldots, Q_{D_x}^x$. As a simple estimate of allele frequencies we take

$$f_x = \frac{\sum_k Q_k^x}{\sum_{k,y} Q_k^y}. \tag{1}$$

When calling jointly on multiple samples with varying coverage, lower-coverage samples would contribute less to the estimate. Therefore we calculate frequencies $f_x^i$ for each sample $i$ as above and then calculate the site frequency as

$$f_x = \frac{\sum_i f_x^i}{N}. \tag{2}$$

Now, given a particular allele set $S$, we introduce the relative frequencies

$$f_{x|S} = \frac{f_x}{\sum_{y \in S} f_y}. \tag{3}$$

We calculate the likelihood of observing the set of alleles $S$ for each sample

$$L_S^i = \sum_{x,y \in S} f_{x|S} \, f_{y|S} \, G_i(xy), \tag{4}$$

where $G_i(xy)$ are the genotype likelihoods PL of $i$-th sample calculated by mpileup[1]. Given the prior probability $\theta$, the number of non-reference alleles $r$ observed across all samples and using the Watterson factor $W_N$

$$W_N = \sum_{k=1}^{2N-1} \frac{1}{k}, \tag{5}$$

we calculate the overall likelihood for all samples given the allele set $S$ as

$$L_S = (W_n \theta)^r \prod_i L_S^i. \tag{6}$$

Finally we select the most likely set of alleles $X \subseteq S$ so that

$$X = \arg\max_S L_S. \tag{7}$$

---

[1]PL $= -10 * log_{10} P(\text{data}|\text{genotype})$

The site quality of variant sites is given by

$$\text{QUAL} = \frac{L_{\{ref\}}}{\sum_S K_S} \tag{8}$$

and the quality of non-variant sites

$$\text{QUAL} = 1 - \frac{L_{\{ref\}}}{\sum_S K_S}. \tag{9}$$

Assuming HWE, the most likely genotype $(xy)_i$ of $i$-th sample is

$$(xy)_i = \underset{a,b \in X}{\arg\max} \ L_X^i \tag{10}$$

and the corresponding genotype quality is

$$\text{GQ} = \frac{L_X^i}{\sum_Y L_Y^i}. \tag{11}$$