# 1 Genotype Likelihoods from Reads

$$L(g = 1) = p(d \mid g = 1)$$
$$= p(d, g = ra \mid g = 1) + p(d, g = rb \mid g = 1) + p(d, g = rc \mid g = 1)$$
$$= \sum_a p(d \mid g = ra) \times p(g = ra \mid g = 1)$$

Here we do not favour any heterozygous genotype, and all have likelihood 1/3. This may be changed to reflect empirical or dbSNP data

$$L(g = 1) = 1/3 \sum_a p(d \mid g = ra)$$
$$= 1/3 \sum_a p(d \mid g = ra, \text{ado}) \times p(\text{ado}) + p(d \mid g = ra, \text{no ado})(1 - p(\text{ado}))$$
$$p(d \mid g = ra, ado) = p(d \mid g = ra, dropr) * p(dropr) + p(d \mid g = ra, dropa) * p(dropa)$$

Here we assume either allele is equally likely to be dropped in an ado event and $p(\text{drop r}) = p(\text{drop a}) = 0.5$. This is unlikely to change.

$$p(d \mid g = ra, \text{ado}, \text{drop a}) = \prod_i p(d_i \mid g = rr)$$
$$p(d \mid g = ra, \text{ado}, \text{drop r}) = \prod_i p(d_i \mid g = aa)$$

# 2 Cell-locus Posterior probabilities

Using Bayes' rule:

$$p(g = k \mid d) = \frac{p(d \mid g = k) \, p(g = k)}{p(d)}$$
$$= \frac{p(d \mid g = k) \, \mu^k \, (1 - \mu)^{2-k}}{p(d)}$$

where k is the mutation rate, a learnable parameter. Note: this parameter may be overestimated if the algorithm finds more mutations, increases the rate prior, and so finds more mutations. There may be no reason for this to converge.

Since $p(d)$ is the same for all values of $k$ at a cell-locus, we do not need to find it and can simply normalise.

# 3 Probabilistic Hamming distance

The Hamming distance between two sequences $s, p$ both length $n$ is given by

$$\sum_i^n (1 - \delta_{s_i p_i})$$

where $\delta_{ab}$ is the Kronecker delta. Since we have a probabilistic tree, we use a similar metric but weighted by the posterior probabilities of the genotypes at each locus:

$$\sum_i^n \sum_{(a,b) \in \{0,1\} \times \{0,1\}} (1 - \delta_{ab}) p(s_i p_i = ab)$$