

## Project 1 Report

With my python script, I take in DNA-seq data from a .bam file and putative SNPs metadata from a .tsv and output two files. The first output file contains the posterior probability of the most probable genotype for each chromosome and is called output.tsv. The second file contains the posterior probability of all possible genotypes for each chromosome in a file called all\_poss\_genotypes\_output.tsv. The first output file is more useful as it contains data about the most probably genotype and less redundant data.

In my script, we first check that the input files are correct. There are some default files for the script to fall back and use if the input files are incorrect. After getting the .bam and .tsv files, we index and parse the data from each file. Then we go through each chromosome in the metadata file and process it. From the processing, we are able to find the posterior probabilities of genotype from the chromosome. We then put all our data in a dataframe that is then output as a .tsv file.

With the processing of the data for each chromosome, first have to obtain the correct sequences to observe from the .bam file. Upon receiving the correct sequences, we need to clean the sequence to only contain the relevant alleles (the major and minor allele which are given to us in the metadata). We then find the posterior probability using math that we derive from our model. We use a biallelic model and have the derivation below.

Below, we have the model used as well as the names for variables and assumptions used.

## Using Bi-allelic Model

### Posterior probability derivations for Bi-allelic Model

#### Variables Used

- A = major allele
- B = minor allele
- $\sigma_A$  = amount of A (proportion of A)
- $\epsilon$  = error
- N = number of observations
- $N_A$  = number of A observations
- maf = minor allele frequency

#### Assumptions made

- $\sigma_A = 1/2$  or 0.5
- prior probabilities
  - 1)  $P(AA) = (1 - \text{maf})^2$
  - 2)  $P(BB) = \text{maf}^2$
  - 3)  $P(AB) = 1 - P(AA) - P(BB)$
- Error  $\rightarrow \epsilon = 0.05$

### Posterior probability of genotype AB

$$P(AB | o_1, \dots, o_N) = \frac{P(o_1, \dots, o_N | AB) P(AB)}{P(o_1, \dots, o_N)}$$

↓

$$\text{Note: } P(o_1, \dots, o_N) = P(o_1, \dots, o_N | AA) P(AA) + P(o_1, \dots, o_N | AB) P(AB) + P(o_1, \dots, o_N | BB) P(BB)$$

↓

$$P(AB | o_1, \dots, o_N) = \frac{P(o_1, \dots, o_N | AB) P(AB)}{P(o_1, \dots, o_N | AA) P(AA) + P(o_1, \dots, o_N | AB) P(AB) + P(o_1, \dots, o_N | BB) P(BB)}$$

Note: Break into smaller sections to keep track of

$$P(o_1, \dots, o_N | AA) P(AA) + P(o_1, \dots, o_N | AB) P(AB) + P(o_1, \dots, o_N | BB) P(BB)$$

$$\begin{aligned} 1) P(o_1, \dots, o_N | AA) P(AA) &= \prod_{i:A} P(o_i = A | AA) \prod_{i:B} P(o_i = B | AA) P(AA) \\ &= \prod_{i:A} P(\epsilon_i = 0) \prod_{i:B} P(\epsilon_i = 1) P(AA) \\ &= (1 - \epsilon)^{N_A} (\epsilon)^{N - N_A} P(AA) \end{aligned}$$

$$\begin{aligned} 2) P(o_1, \dots, o_N | BB) P(BB) &= \prod_{i:A} P(o_i = A | BB) \prod_{i:B} P(o_i = B | BB) P(BB) \\ &= \prod_{i:A} P(\epsilon_i = 1) \prod_{i:B} P(\epsilon_i = 0) P(BB) \\ &= \epsilon^{N_A} (1 - \epsilon)^{N - N_A} P(BB) \end{aligned}$$

$$\begin{aligned} 3) P(o_1, \dots, o_N | AB) P(AB) &= \left( \prod_{i:A} P(\epsilon_i = 0 | \sigma = A) P(\sigma = A) + P(\epsilon_i = 1 | \sigma = B) P(\sigma = B) \right) \\ &\quad \left( \prod_{i:B} P(\epsilon_i = 1 | \sigma = A) P(\sigma = A) + P(\epsilon_i = 0 | \sigma = B) P(\sigma = B) \right) P(AB) \end{aligned}$$

Note:

$$\left( \prod_{i:A} P(\tilde{y}_i = 0 | y = A) P(y = A) + P(\tilde{y}_i = 1 | y = B) P(y = B) \right) \left( \prod_{i:B} P(\tilde{y}_i = 1 | y = A) P(y = A) + P(\tilde{y}_i = 0 | y = B) P(y = B) \right) P(A, B)$$

$\downarrow$                        $\downarrow$                        $\downarrow$                        $\downarrow$                        $\downarrow$                        $\downarrow$                        $\downarrow$                        $\downarrow$   
 $(1-c)$                        $s_A$                        $c$                        $1-s_A$                        $c$                        $s_A$                        $1-c$                        $1-s_A$

raise expression to  $N_A$  power

raise expression to  $N - N_A$  power

$$\left( (1-c) s_A + c (1-s_A) \right)^{N_A} \left( c (s_A) + (1-c) (1-s_A) \right)^{N - N_A} P(A, B)$$

$$P(O_1, \dots, O_n | A, B) P(A, B) = \left( (1-c) s_A + c (1-s_A) \right)^{N_A} \left( c (s_A) + (1-c) (1-s_A) \right)^{N - N_A} P(A, B)$$

$$P(A, B | O_1, \dots, O_n) = \frac{P(O_1, \dots, O_n | A, B) P(A, B)}{P(O_1, \dots, O_n | A, A) P(A, A) + P(O_1, \dots, O_n | A, B) P(A, B) + P(O_1, \dots, O_n | B, B) P(B, B)}$$

$$= \frac{\left( (1-c) s_A + c (1-s_A) \right)^{N_A} \left( c (s_A) + (1-c) (1-s_A) \right)^{N - N_A} P(A, B)}{(1-c)^{N_A} (c)^{N - N_A} P(A, A) + \left( (1-c) s_A + c (1-s_A) \right)^{N_A} \left( c (s_A) + (1-c) (1-s_A) \right)^{N - N_A} P(A, B) + c^{N_A} (1-c)^{N - N_A} P(B, B)}$$

posterior probability for AB  
 ↳ derivations for AA and AB will be similar

posterior probability AB

$$\frac{\left( (1-c) s_A + c (1-s_A) \right)^{N_A} \left( c (s_A) + (1-c) (1-s_A) \right)^{N - N_A} P(A, B)}{(1-c)^{N_A} (c)^{N - N_A} P(A, A) + \left( (1-c) s_A + c (1-s_A) \right)^{N_A} \left( c (s_A) + (1-c) (1-s_A) \right)^{N - N_A} P(A, B) + c^{N_A} (1-c)^{N - N_A} P(B, B)}$$

posterior probability AA

$$\frac{(1-c)^{N_A} (c)^{N - N_A} P(A, A)}{(1-c)^{N_A} (c)^{N - N_A} P(A, A) + \left( (1-c) s_A + c (1-s_A) \right)^{N_A} \left( c (s_A) + (1-c) (1-s_A) \right)^{N - N_A} P(A, B) + c^{N_A} (1-c)^{N - N_A} P(B, B)}$$

posterior probability BB

$$\frac{c^{N_A} (1-c)^{N - N_A} P(B, B)}{(1-c)^{N_A} (c)^{N - N_A} P(A, A) + \left( (1-c) s_A + c (1-s_A) \right)^{N_A} \left( c (s_A) + (1-c) (1-s_A) \right)^{N - N_A} P(A, B) + c^{N_A} (1-c)^{N - N_A} P(B, B)}$$