

Krish Patel

CS C121

Assignment 2

Q1)

- There are 9 different possible genotypes that can occur. They are: GG, GA, GC AA, AG, AC, CC, CG, CA
- By increasing the number of possible nucleotides in our "sample" we are going to introduce more observation possibilities, and thus the probability for observing the genomes would decrease due to these considerations
- Probabilities for homozygous genotypes:

c) $P(O_i | G = xx) \rightarrow$ where $x = A, G, \text{ or } C$.

$P(O_i = A | G = AA)$

Summing to $(1-\epsilon)$

$$\begin{aligned} &P(O_i = A | E_i = 0, C_i = A, G = AA) \cdot P(E_i = 0) \cdot P(C_i = A | G = AA) \\ &= \underbrace{1}_{\downarrow} \cdot \underbrace{\frac{1}{1-\epsilon}}_{\downarrow} \cdot \underbrace{1}_{\downarrow} \\ &P(O_i = A | E_i = 0, C_i = C, G = AA) \cdot P(E_i = 0) \cdot P(C_i = C | G = AA) \\ &= \underbrace{0}_{\downarrow} \cdot \underbrace{\frac{1}{1-\epsilon}}_{\downarrow} \cdot \underbrace{1}_{\downarrow} \\ &P(O_i = A | E_i = 0, C_i = G, G = AA) \cdot P(E_i = 0) \cdot P(C_i = G | G = AA) \\ &= \underbrace{0}_{\downarrow} \cdot \underbrace{\frac{1}{1-\epsilon}}_{\downarrow} \cdot \underbrace{1}_{\downarrow} \end{aligned}$$

Summing to 0

$$\begin{aligned} &P(O_i = A | E_i = 1, C_i = A, G = AA) \cdot P(E_i = 1) \cdot P(C_i = A | G = AA) \\ &= \underbrace{0}_{\downarrow} \cdot \underbrace{\frac{1}{1-\epsilon}}_{\downarrow} \cdot \underbrace{1}_{\downarrow} \\ &P(O_i = A | E_i = 1, C_i = C, G = AA) \cdot P(E_i = 1) \cdot P(C_i = C | G = AA) \\ &= \underbrace{1}_{\downarrow} \cdot \underbrace{\frac{1}{1-\epsilon}}_{\downarrow} \cdot \underbrace{0}_{\downarrow} \\ &P(O_i = A | E_i = 1, C_i = G, G = AA) \cdot P(E_i = 1) \cdot P(C_i = G | G = AA) \\ &= \underbrace{1}_{\downarrow} \cdot \underbrace{\frac{1}{1-\epsilon}}_{\downarrow} \cdot \underbrace{0}_{\downarrow} \end{aligned}$$

$$P(C|g=AA)$$

$$P(G|g=AA)$$

$$P(C|g=AA)$$

$$=$$

$$E_i=0$$

$$P(O_i=C|E_i=0, C_i=A, G=AA) = \underset{0}{\downarrow} P(E_i=0) = \underset{1-\epsilon}{\downarrow} P(C_i=A|G=AA) = \underset{1}{\downarrow}$$

$$P(O_i=C|E_i=0, C_i=C, G=AA) = \underset{0}{\downarrow} P(E_i=0) = \underset{1-\epsilon}{\downarrow} P(C_i=C|G=AA) = \underset{0}{\downarrow}$$

$$P(O_i=C|E_i=0, C_i=G, G=AA) = \underset{0}{\downarrow} P(E_i=0) = \underset{1-\epsilon}{\downarrow} P(C_i=G|G=AA) = \underset{0}{\downarrow}$$

$$= 0$$

$$E_i=1$$

$$P(O_i=C|E_i=1, C_i=A, G=AA) = \underset{0.5}{\downarrow} P(E_i=1) = \underset{\epsilon}{\downarrow} P(C_i=A|G=AA) = \underset{1}{\downarrow}$$

$$P(O_i=C|E_i=1, C_i=C, G=AA) = \underset{0}{\downarrow} P(E_i=1) = \underset{\epsilon}{\downarrow} P(C_i=C|G=AA) = \underset{0}{\downarrow}$$

$$P(O_i=C|E_i=1, C_i=G, G=AA) = \underset{0.5}{\downarrow} P(E_i=1) = \underset{\epsilon}{\downarrow} P(C_i=G|G=AA) = \underset{0}{\downarrow}$$

$$= 0.5\epsilon$$

SAME FOR G!!!

Thus, given above samples.

$$P(O=x|G=xx) = 1-\epsilon$$

$$P(O=y|G=xx) = \frac{\epsilon}{2} \text{ (for } x \neq y)$$

d)

$$\textcircled{d} \quad \epsilon=0$$

$$P(O_i=A|G=GA) = P(O_i=A|E_i=0, C_i=A, G=GA) = \underset{1}{\downarrow} P(E_i=0) = \underset{1-\epsilon}{\downarrow} P(C_i=A|G=GA) = \underset{0.5}{\downarrow}$$

$$P(O_i=A|G=GA) = P(O_i=A|E_i=0, C_i=C, G=GA) = \underset{0}{\downarrow} P(E_i=0) = \underset{1-\epsilon}{\downarrow} P(C_i=C|G=GA) = \underset{0}{\downarrow}$$

$$P(O_i=A|G=GA) = P(O_i=A|E_i=0, C_i=G, G=GA) = \underset{0}{\downarrow} P(E_i=0) = \underset{1-\epsilon}{\downarrow} P(C_i=G|G=GA) = \underset{0.5}{\downarrow}$$

$$\epsilon=1 \hookrightarrow$$

$$P(O_i=A|G=GA) = P(O_i=A|E_i=1, C_i=A, G=GA) = \underset{0}{\downarrow} P(E_i=1) = \underset{\epsilon}{\downarrow} P(C_i=A|G=GA) = \underset{0.5}{\downarrow}$$

$$P(O_i=A|G=GA) = P(O_i=A|E_i=1, C_i=C, G=GA) = \underset{0.5}{\downarrow} P(E_i=1) = \underset{\epsilon}{\downarrow} P(C_i=C|G=GA) = \underset{0}{\downarrow}$$

$$P(O_i=A|G=GA) = P(O_i=A|E_i=1, C_i=G, G=GA) = \underset{0.5}{\downarrow} P(E_i=1) = \underset{\epsilon}{\downarrow} P(C_i=G|G=GA) = \underset{0.5}{\downarrow}$$

$\epsilon = 0$ Sums to 0

$$\begin{aligned}
 P(O_i = C | G = GA) &= P(O_i = C | E_i = 0, c_i = A, G = GA) \cdot P(E_i = 0) \cdot P(c_i = A | G = GA) \\
 &\quad \downarrow \quad \quad \quad \downarrow \quad \quad \quad \downarrow \\
 &\quad 0 \quad \quad \quad 1 - \epsilon \quad \quad \quad 0 \\
 P(O_i = C | G = GA) &= P(O_i = C | E_i = 0, c_i = C, G = GA) \cdot P(E_i = 0) \cdot P(c_i = C | G = GA) \\
 &\quad \downarrow \quad \quad \quad \downarrow \quad \quad \quad \downarrow \\
 &\quad 0 \quad \quad \quad 1 - \epsilon \quad \quad \quad 0 \\
 P(O_i = C | G = GA) &= P(O_i = C | E_i = 0, c_i = G, G = GA) \cdot P(E_i = 0) \cdot P(c_i = G | G = GA) \\
 &\quad \downarrow \quad \quad \quad \downarrow \quad \quad \quad \downarrow \\
 &\quad 0 \quad \quad \quad 1 - \epsilon \quad \quad \quad 1
 \end{aligned}$$

$\epsilon = 1$ Sums to $\frac{1}{2}\epsilon$

$$\begin{aligned}
 P(O_i = C | G = GA) &= P(O_i = C | E_i = 1, c_i = A, G = GA) \cdot P(E_i = 1) \cdot P(c_i = A | G = GA) \\
 &\quad \downarrow \quad \quad \quad \downarrow \quad \quad \quad \downarrow \\
 &\quad 0.5 \quad \quad \quad \epsilon \quad \quad \quad 0.5 \\
 P(O_i = C | G = GA) &= P(O_i = C | E_i = 1, c_i = C, G = GA) \cdot P(E_i = 1) \cdot P(c_i = C | G = GA) \\
 &\quad \downarrow \quad \quad \quad \downarrow \quad \quad \quad \downarrow \\
 &\quad 0 \quad \quad \quad \epsilon \quad \quad \quad 0 \\
 P(O_i = C | G = GA) &= P(O_i = C | E_i = 1, c_i = G, G = GA) \cdot P(E_i = 1) \cdot P(c_i = G | G = GA) \\
 &\quad \downarrow \quad \quad \quad \downarrow \quad \quad \quad \downarrow \\
 &\quad 0.5 \quad \quad \quad \epsilon \quad \quad \quad 0
 \end{aligned}$$

Thus, this is equal to

$$\frac{1}{2}(1 - \epsilon) + \frac{1}{4}\epsilon = \frac{1}{2} - \frac{1}{4}\epsilon$$

same for $P(O_i = G | G = GA) = \frac{1}{2} - \frac{1}{4}\epsilon$

e) Consider the following case

$$P(O_i = X | G = XY) \text{ and } P(O_i = X | G = YX)$$

This should have the same value for the error rates, due to the symmetries of the genome. It is equally likely to observe a X from a XY and YX due to the assumptions we have made.

Thus

$$P(O_i = A | G = AG) = P(O_i = A | G = GA)$$

$$P(O_i = A | G = AG) = P(O_i = A | G = GA) \text{ etc for all genotypes of the form XY}$$

Now consider the following

$$P(O_i = X | G = XY) \text{ and } P(O_i = Z | G = ZY)$$

These should also have the same probabilities considering that both of them are calculated the same way, and the only difference is the nucleotide which doesn't really differ (the error rates are calculated the same way) Thus we can calculate the probabilities for $P(O_i = A | G = GA)$ and $P(O_i = C | G = GC)$, etc and thus we could find all the probabilities associated with all genotype

Q2)

a)

The likelihood probability of O given $P(O|g=AB)$ reduces to $\frac{1}{2}$ for both $A=O$ and $B=O$ regardless of what Epsilon is. Thus this simplifies our calculations.

b)

The code for estimating the posterior probability is:

```
def calculate_posterior(sample):
    length = len(sample)
    AA_errors = []
    AT_errors = []
    TT_errors = []

    for i in range(length):
        if sample[i][0] == 'T':
            AA_errors.append(float(sample[i][1]))
            AT_errors.append(0.5)
            TT_errors.append(1- float(sample[i][1]))
        else:
            AA_errors.append(1- float(sample[i][1]))
            AT_errors.append(0.5)
            TT_errors.append(float(sample[i][1]))
    return AA_errors, AT_errors, TT_errors

def calculate_posterior_probabilities(AA_errors, AT_errors, TT_errors, n=5):
    multAA = 1
    multAT = 1
    multTT = 1
    for i in range(n):
        multAA = multAA * AA_errors[i]
        multAT = multAT * AT_errors[i]
        multTT = multTT * TT_errors[i]
    return multAA, multAT, multTT

#testing function

AA, AT, TT = calculate_posterior(random_reads)
AAprobs_sample, ATprobs_sample, TTprobs_sample = calculate_posterior_probabilities(AA, AT, TT)

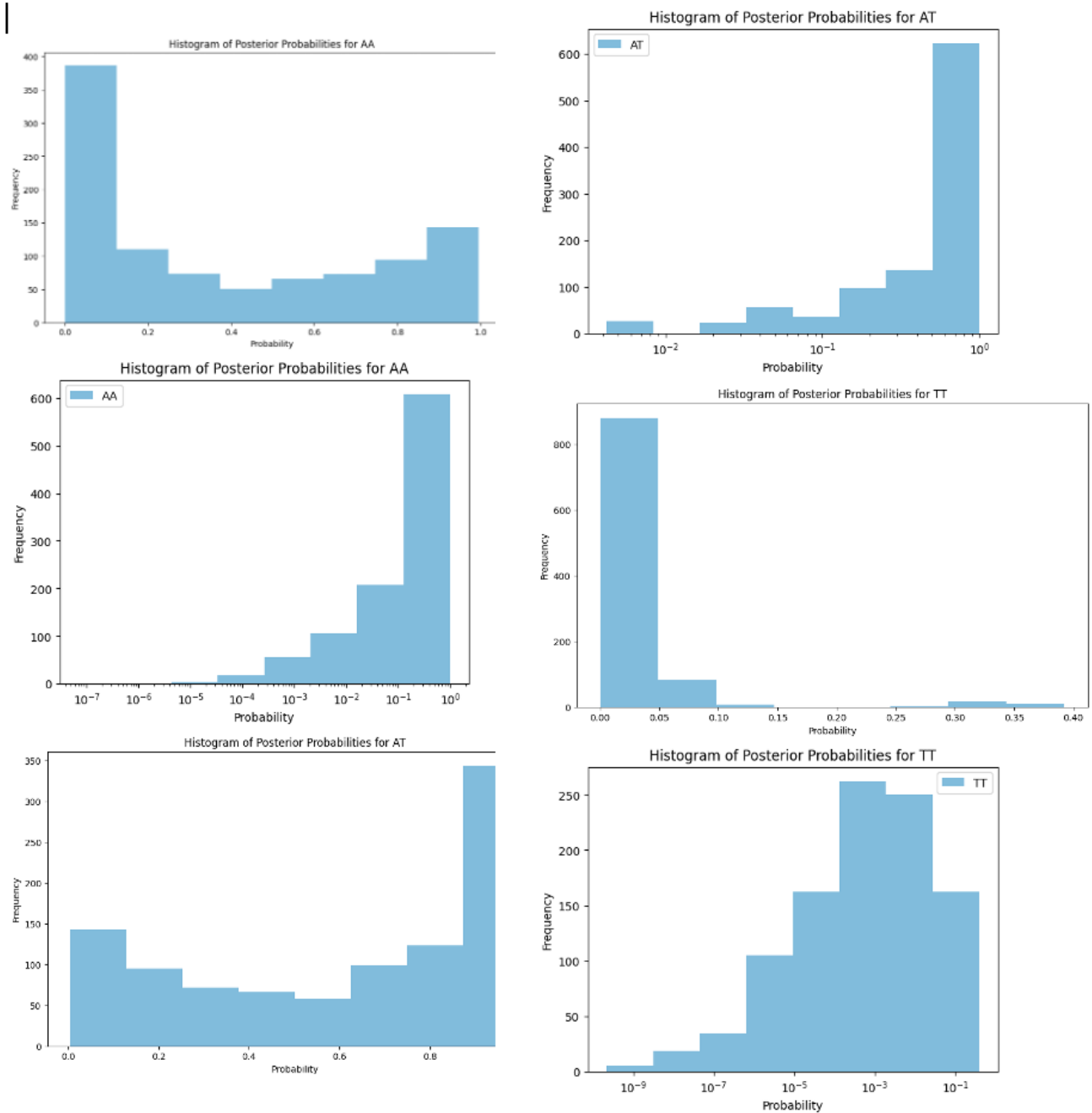
8
9 AA, AT, TT = calculate_posterior(random_reads)
10 AAprobs_sample, ATprobs_sample, TTprobs_sample = calculate_posterior_probabilities(AA, AT, TT)
11
12 def normalized_probs(AAprobs_sample, ATprobs_sample, TTprobs_sample):
13     AAprs = 0.95**2
14     TTprobs= 0.05**2
15     ATprobs = 1- AAprs - TTprobs
16     AAprs_sample = AAprs_sample * AAprs
17     ATprobs_sample = ATprobs_sample * ATprobs
18     TTprobs_sample = TTprobs_sample * TTprobs
19     total = AAprs_sample + ATprobs_sample + TTprobs_sample
20     return AAprs_sample/total, ATprobs_sample/total, TTprobs_sample/total
21
22 print(normalized_probs(AAprs_sample, ATprobs_sample, TTprobs_sample))
23
```

c) The values for the probabilities when using the seed of 1 given that the sample size is 5 is (0.19616364841610875, 0.8036409729284506, 0.00019537865544056428)

Thus, according to this sample it is 80% likely that the genotype is AT

The results are in the following order (probs AA, probs AT, probs TT)

d) For the histograms, I've used two different scales, a linear scale and a logarithmic scale for visualization and thus get a total of 6 histograms
Histogram for 5:



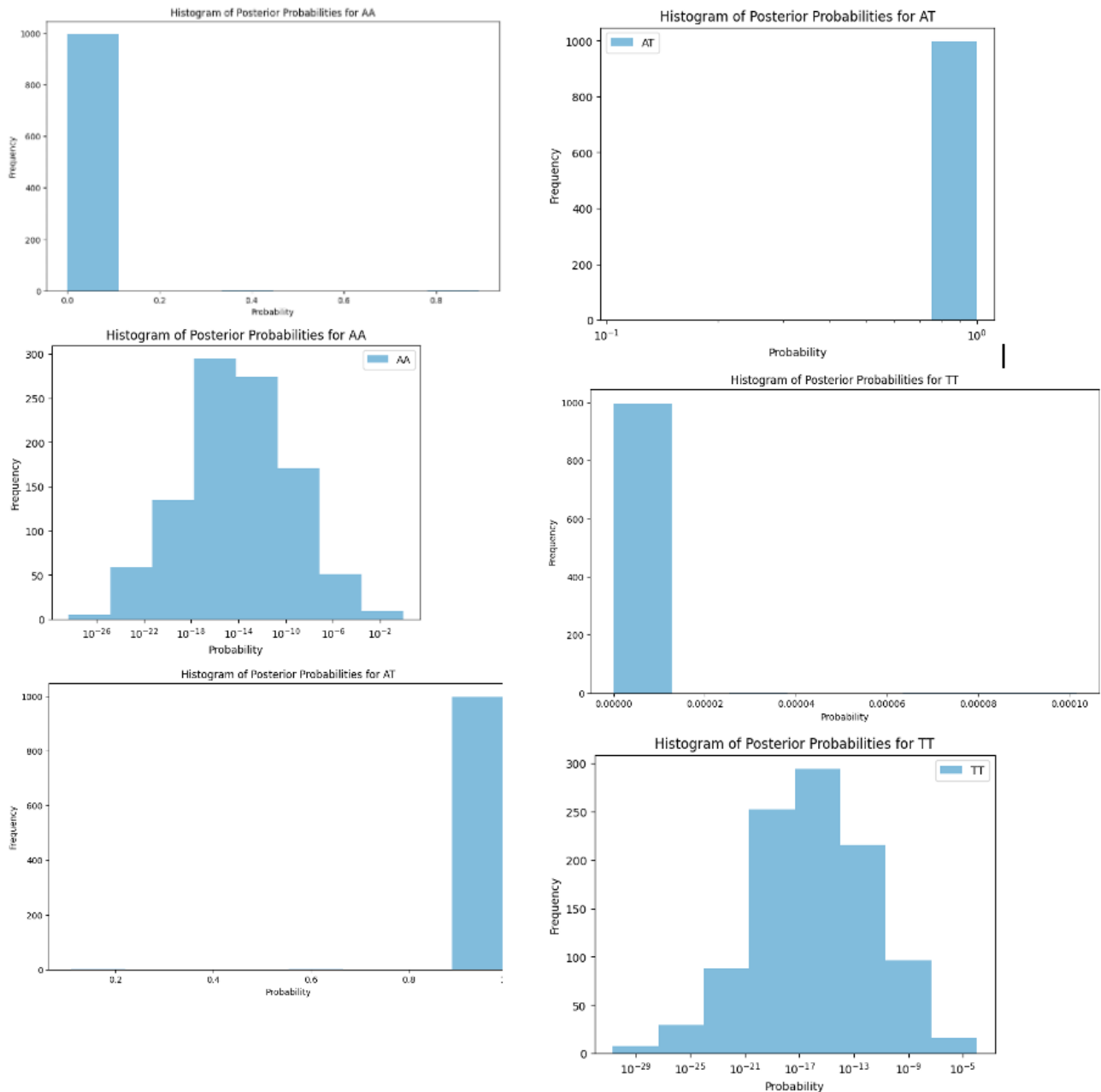
With 5 observations per sample:

AA mean: 0.3770430849784912 AA std: 0.3542507442428913

AT mean: 0.6028430218143397 AT std: 0.338547742866706

TT mean: 0.020113893207169248 TT std: 0.05967340498218455

e) Histogram for 50



By increasing the number of observations in the sample, the histogram becomes much more neater as the noise is diminished (outlier probabilities diminish). We get more variation in the sample. This can be attributed to the fact that with 5 samples, you could never get equal number of As and Ts, however with 50 you usually get a ratio much closer to 50%.

With 50 observations per sample:

AA mean: 0.0014676411022330423 AA std: 0.03136012709503697

AT mean: 0.998532056032999 AT std: 0.03136011327689436

TT mean: 3.0286476786699527e-07 TT std: 4.724017025574451e-06

- f) From the above statistics, it is clear that with an increased number of samples, the standard deviation reduces by a significant margin, and the mean probabilities becomes much more clear and definitive. For instance, mean probability for AT is 60% with 5 reads, but with 50 it becomes ~99.9% thus giving much more accurate results.
- g) The assumptions being made here are that it is equally likely to observe both genotypes, which may not always be true. In real-world scenarios, the likelihood of observing different genotypes may vary based on various factors such as genetic mutations, environmental influences, or selective pressures. Therefore, while we may assume symmetry in the probabilities of observing certain genotypes, it's important to acknowledge that this assumption might not hold universally.

Additionally, another assumption that we make here is that the error rates and the probabilities are true; however, this is also not necessarily the case. In genetic sequencing or any experimental context, there's always a degree of uncertainty associated with measurements and observations. Error rates in sequencing technologies, variation in experimental conditions, or biases in data collection processes can all introduce inaccuracies or biases into the observed probabilities.