

# Homework 3

Michael Li

1. (a) code submitted online

List of high G+C regions: (66, 415), (527, 720), (951, 1000)

(b) code submitted online

screenshot of graph on next page

$\mu$	$\ln P(x \mu)$	$\mu$	$\ln P(x \mu)$
0.0001	-1309.121	0.05	-1308.874
0.001	-1298.423	0.1	-1325.005
0.006	-1292.987	0.2	-1344.422
0.007	-1292.887	0.5	-1369.246
0.008	-1292.896	0.75	-1381.872
0.01	-1293.142	0.99	-1409.911

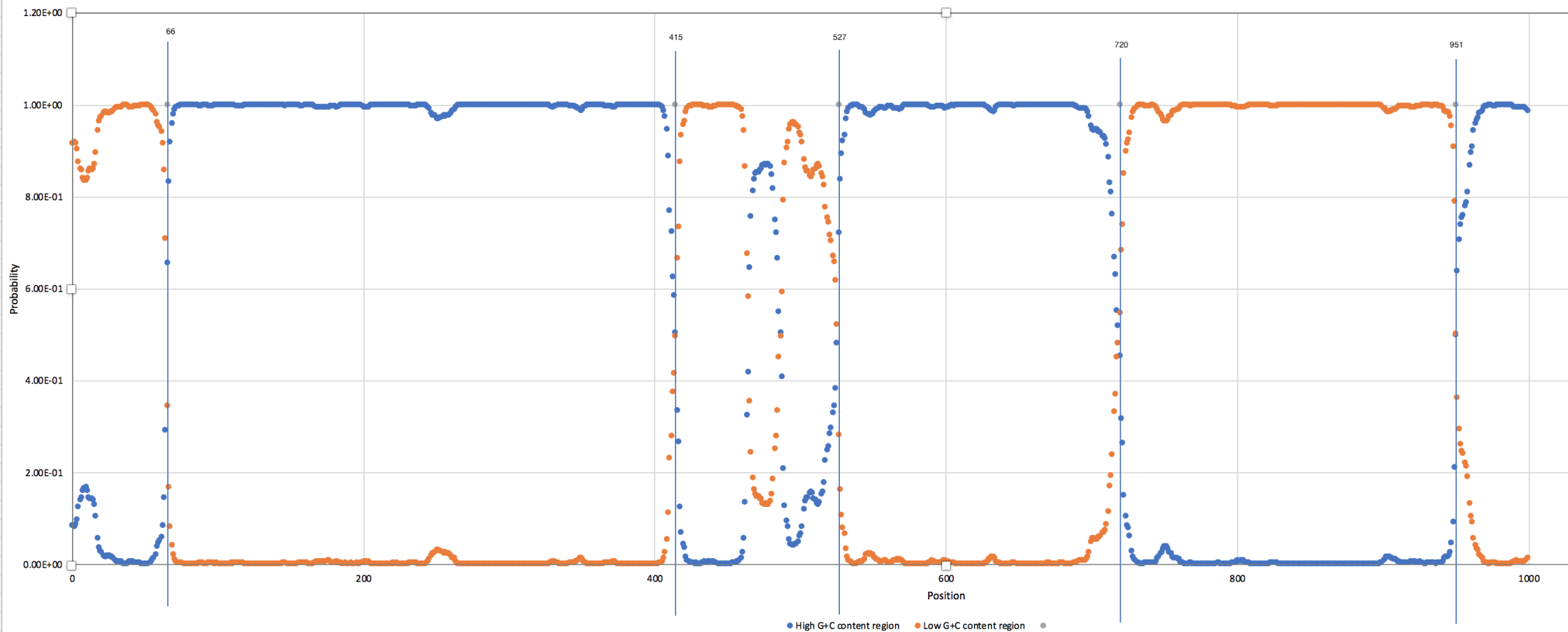
plot attached as screenshot

It looks like the MLE of  $\mu$  is around 0.007, to 0.0075.

The graph has a pretty sharp peak around that area before gently tapering off at high values.

If I used smaller values for  $\mu$  the viterbi path should have fewer switches between the two states, while for a larger  $\mu$  it should have more switches.

Probabilities vs Sequence Position Graph



$\mu$	$\ln(P)$
0.0001	-1309.121
0.001	-1298.423
0.006	-1292.987
0.007	-1292.887
0.008	-1292.896
0.01	-1293.142
0.05	-1308.874
0.1	-1325.005
0.2	-1344.422
0.5	-1369.246
0.75	-1381.872
0.99	-1409.911



$$2. P(x_i | z_i = h) = \begin{cases} \frac{\theta_h}{2} & x_i = G \text{ or } C \\ \frac{1-2\theta_h}{2} & x_i = A \text{ or } T \end{cases} \quad P(x_i | z_i = l) = \begin{cases} \frac{\theta_l}{2} & x_i = G \text{ or } C \\ \frac{1-2\theta_l}{2} & x_i = A \text{ or } T \end{cases}$$

$$(a) \log P(X, Z | \mu, \theta_h, \theta_l) = ?$$

$$P(X, Z) = \underset{\substack{\uparrow \\ \text{initial} \\ \text{prob}}}{a_{0,z_1}} \cdot \underset{\substack{\uparrow \\ \text{initial} \\ \text{emission} \\ \text{prob}}}{e_{z_1}(x_1)} \cdot \prod_{i=2}^n \underset{\substack{\uparrow \\ \text{emission} \\ \text{at } i}}{e_{z_i}(x_i)} \cdot \underset{\substack{\uparrow \\ \text{transition} \\ \text{prob}}}{a_{z_{i-1}, z_i}}$$

$$P(X, Z | \mu, \theta_h, \theta_l) = 0.5 \cdot (1 - \frac{\theta_l}{2}) \cdot \prod_{i=2}^{1000} P(x_i | z_i) \cdot P(z_i | z_{i-1}) \leftarrow \begin{cases} \mu & z_i = z_{i-1} \\ 1-\mu & z_i \neq z_{i-1} \end{cases}$$

$$= 0.5 \cdot e_l(x_1) \cdot \prod_{i=2}^{65} e_l(x_i) \cdot (1-\mu) \cdot e_h(x_{66}) \cdot \mu \cdot \prod_{i=67}^{417} e_h(x_i) \cdot (1-\mu) \cdot e_l(x_{418}) \cdot \mu \cdot \prod_{i=419}^{467} e_l(x_i) \cdot (1-\mu) \cdot e_h(x_{468}) \cdot \mu$$

$$\cdot \prod_{i=469}^{1000} e_h(x_i) \cdot (1-\mu)$$

across all intervals

Since all products, we can move variables around

$$= 0.5 \cdot \mu^{C_b} \cdot (1-\mu)^{C_s} \cdot \prod_{i=1}^{65} e_l(x_i) \cdot \prod_{i=66}^{417} e_h(x_i) \cdot \prod_{i=418}^{467} e_l(x_i) \cdot \prod_{i=468}^{1000} e_h(x_i)$$

$$= 0.5 \cdot \mu^{C_b} \cdot (1-\mu)^{C_s} \cdot \prod_{i \text{ in state } l} e_l(x_i) \cdot \prod_{i \text{ in state } h} e_h(x_i)$$

$$= 0.5 \cdot \mu^{C_b} \cdot (1-\mu)^{C_s} \cdot \left(\frac{\theta_l}{2}\right)^{d_{eL}} \cdot \left(\frac{1-2\theta_l}{2}\right)^{d_{eA}} \cdot \left(\frac{\theta_h}{2}\right)^{d_{eG}} \cdot \left(\frac{1-2\theta_h}{2}\right)^{d_{eA}}$$

$$\log P(X, Z | \mu, \theta_h, \theta_l) = \log($$

$$= \log(0.5) + C_b \cdot \log(\mu) + C_s \cdot \log(1-\mu) + d_{eL} \cdot \log\left(\frac{\theta_l}{2}\right)$$

$$+ d_{eA} \cdot \log\left(\frac{1-2\theta_l}{2}\right) + d_{eG} \cdot \log\left(\frac{\theta_h}{2}\right) + d_{eA} \cdot \log\left(\frac{1-2\theta_h}{2}\right)$$



(b) To maximize  $\log P(X, Z | \mu, \theta_h, \theta_e)$  take the derivative  
 Since we have 3 variables, take the partial derivative of each

$$\log P(X, Z | \mu, \theta_h, \theta_e) = \log(0.5) + C_b \cdot \log(\mu) + C_s \cdot \log(1-\mu) + d_{LB} \cdot \log\left(\frac{\theta_e}{2}\right) + d_{LA} \cdot \log\left(\frac{1}{2} - \theta_e\right) + d_{hB} \cdot \log\left(\frac{\theta_h}{2}\right) + d_{hA} \cdot \log\left(\frac{1}{2} - \theta_h\right)$$

$$\frac{d}{d\mu} \log P(X, Z | \mu, \theta_h, \theta_e) = \frac{C_b}{\mu} - \frac{C_s}{1-\mu} = 0 \Rightarrow C_b/\mu = C_s/(1-\mu)$$

$$C_b(1-\mu) = C_s \cdot \mu$$

$$C_b - C_b \cdot \mu = C_s \cdot \mu$$

$$C_b = C_s \cdot \mu + C_b \cdot \mu$$

$$\mu = \frac{C_b}{C_s + C_b}$$

$$\mu = \frac{7}{993+7} = \frac{7}{1000} = 0.007$$

$$\frac{d}{d\theta_e} \log P(X, Z | \mu, \theta_h, \theta_e) = \frac{d_{LB}}{\frac{\theta_e}{2}} - \frac{d_{LA}}{\frac{1}{2} - \theta_e} = 0 \Rightarrow \frac{d_{LB} \cdot 2}{\theta_e} = \frac{d_{LA}}{\frac{1}{2} - \theta_e}$$

$$d_{LB}(1-2\theta_e) = d_{LA} \cdot \theta_e$$

$$d_{LB} = d_{LA} \theta_e + 2 d_{LB} \theta_e$$

$$\theta_e = \frac{d_{LB}}{d_{LA} + 2 d_{LB}}$$

$$\theta_e = \frac{115}{234+230} = \frac{115}{464} = 0.2478$$

$$\frac{d}{d\theta_h} \log P(X, Z | \mu, \theta_h, \theta_e) = \frac{d_{hB}}{\frac{\theta_h}{2}} - \frac{d_{hA}}{\frac{1}{2} - \theta_h} = 0 \Rightarrow d_{hB}(1-2\theta_h) = d_{hA} \cdot \theta_h$$

$$d_{hB} = d_{hA} \theta_h + 2 d_{hB} \theta_h$$

$$\theta_h = \frac{d_{hB}}{d_{hA} + 2 d_{hB}}$$

$$\theta_h = \frac{495}{156+990} = \frac{495}{1146} = 0.4319$$

We decomposed the parameters by taking partial derivatives which led to 3 equations containing one parameter variable each, allowing us to maximize the parameters.

$$\hat{\mu} = 0.007, \hat{\theta}_e = 0.2478, \hat{\theta}_h = 0.4319$$



### 3(a) Code submitted online

Sequence outputs also submitted separately

The log likelihoods generated are pretty variable, with the range being about 150. This may not be too much when considering the mean in the -1500s, but since a difference of 23 corresponds to a 10 fold difference in likelihood, there is a big difference.

The transition mean and variance were 9.86 and 10.0004. I ran the code a few more times and the mean + variance was consistently around those numbers

This is very close to the mean and variance from the binomial model for 999 possible positions

$$\mu = n \cdot p = 999 \cdot 0.01 = 9.99$$

close to 9.86

$$\sigma^2 = n \cdot p(1-p) = 999 \cdot 0.01 \cdot 0.99 = 9.8901$$

10.0004

Sometimes the generated mean and variance were off by much more but this is reasonable from a small sample size of 50.



more but this is reasonable from a small sample size of 50.

(b) Code Submitted online

sequence output also submitted separately

The log likelihoods for the 50 sampled state paths are all close, but not quite as large as the likelihood given in 1a by viterbi. This makes sense since viterbi maximizes the likelihood. The locations of the high and low GTC regions are all very similar, with some variation around the transition areas.

Plot on next page

Note: the plot was generated by a different runthrough of the code than the rest of the outputs were, so values are slightly different.

The plot of #states at each position is very similar to the posterior state probabilities plot from 1(b)

Number of observations for each state at each position

