

6

CM124 - Spring 2018 - Midterm

April 25, 2018

Name: _____
Student ID: _____

- Exam questions are multiple choice, only one answer is correct.
- Any questions without an answer selection, or multiple answer selections, will be marked as incorrect.
- Each question must have a written justification in the box below the answer selection. This justification must convince the reader that you had a valid reason for choosing the answer you chose, and that your choice was not random. If you do not justify your answer, or if your justification is not convincing, the problem will be marked as incorrect.
- If you answer all questions correctly your grade will be 100.
- If you answer correctly 5 out of 6 questions your grade will be 95.
- If you answer correctly 4 out of 6 questions your grade will be 90.
- If you answer correctly 3 out of 6 questions your grade will be 85.
- If you answer correctly 2 out of 6 questions your grade will be 75.
- If you answer correctly 1 out of 6 questions your grade will be 60.
- If you answer correctly 0 out of 6 questions your grade will be 0.
- The exam is open notes, but electronics of any kind are not allowed.

1. **MLE of Allele Frequency.** Suppose you have the following alleles and their respective counts:
0 - 600
1 - 200

What is the maximum likelihood estimate of the frequency of the 1 allele?

Select from the following options and justify briefly in the box why you choose your answer.

- (a) 0.1
- (b) 0.75
- (c) 0.9
- (d) 0.25

The MLE in class shows that it's $\frac{m}{n}$ where $m = \# 1s$, $n = \text{total } \# \text{ alleles}$.
So: $\frac{m}{n} = \frac{200}{600+200} = \frac{200}{800} = \frac{2}{8} = \frac{1}{4} = 0.25$ (D)

2. MLE question of Haplotype Frequency.

Suppose you have the following haplotypes and their respective counts:

1001 - 300
1101 - 1200
0011 - 400
0101 - 100

What is the maximum likelihood estimate of the haplotype 1101?

Select from the following options and justify briefly in the box why you choose your answer.

- (a) 0.1
- (b) 0.25
- (c) 0.6
- (d) 0.2

similar to previous problem, the MLE is $\frac{m}{N}$, as derived in class, where m is count of 1101 haplotype & $N =$ total counts of all haplotypes. so:

$$\frac{m}{n} = \frac{1200}{\underbrace{1200+300}_{1500} + \underbrace{400+100}_{500}} = \frac{1200}{2000} = \frac{12}{20} = \frac{6}{10} = 0.6$$

C

3. **Trio Phasing.** Suppose you have the following trio with the genotypes as below

Parent1 = 122102

Parent2 = 011120

Child = 112011

Assume there is no recombination. What is the phase of the child in this trio?

Select from the following options and justify briefly in the box why you choose your answer.

- (a) 011001, 101010
- (b) 111001, 001010
- (c) 011001, 111010
- (d) 110010, 001010

Parent1: 122102 < $\begin{matrix} ?11?01 \\ ?11?01 \end{matrix}$] To Child

Parent2: 011120 < $\begin{matrix} 0???10 \\ 0???10 \end{matrix}$] to child

child: 112011 < $\begin{matrix} ??10??? \\ ??10??? \end{matrix}$ From Parent 1

From Parent 2

~~We discern~~

We discern:

parent1: $\begin{matrix} \underline{1} \underline{1} \underline{1} \underline{0} \underline{0} \underline{1} \\ \underline{0} \underline{1} \underline{1} \underline{1} \underline{0} \underline{1} \end{matrix}$ *

parent2: $\begin{matrix} \underline{0} \underline{0} \underline{1} \underline{0} \underline{1} \underline{0} \\ \underline{0} \underline{1} \underline{0} \underline{1} \underline{1} \underline{0} \end{matrix}$ *

child: $\begin{matrix} \underline{1} \underline{1} \underline{1} \underline{1} \underline{0} \underline{0} \underline{1} & P1 \\ \underline{0} \underline{0} \underline{1} \underline{0} \underline{1} \underline{0} & P2 \end{matrix}$

} B

4. Clark's algorithm.

Suppose you have the following genotypes:

AGHAA ^{AGAAA}
 GAGHA
 HHGHA
 AHAAA
 HHHHA

Perform Clark's algorithm, starting from the first genotype, what is the final haplotype(s) that were added to your known set.

Select from the following options and justify briefly in the box why you choose your answer.

- (a) AAAAA
- (b) GGGGA
- (c) AGAAA
- (d) AGGAA

A is 0 (Homozygous)	knowns!
G is 2	AGGAA
H is 1 (Heterozygous)	AGAAA
A < AGGAA	GAGAA
A < AGAAA	GAGGA
	AAAAA

GAGHA < GAGAA
 GAGHA < GAGGA

HHGHA < ??G?A } can be solved with: $\frac{AGGAA}{GAGGA}$
 HHGHA < ??G?A } HHGHA

AHAAA < AAAAA
 AHAAA < AGAAA * already added

HHHHA < ????A } can be solved with: $\frac{AGAAA}{GAGGA}$
 HHHHA < ????A } HHHHA

Since we only use knowns to find the unknowns, the last one we added to the known bank was AAAAA A

$$1 \cdot 4 + \frac{1}{2} \cdot 4 = 4 + 2 = 6$$

5. Haplotype phasing with the EM algorithm.

Suppose you have genotypes $G = \{01210, 10222, 00110\}$. Perform one round of the EM algorithm for haplotype phasing (assume that the haplotypes are equally probable $p_1 = p_2 \dots = p_n = 1/n$). After one round of EM, what haplotype(s) have the highest estimated probability?

Select from the following options and justify briefly in the box why you choose your answer.

- (a) 00110, 10111
- (b) ~~00100, 00111, 00010, 00000~~
- (c) 00111, 00111, 00110, 00100
- (d) 10111, 00111, ~~00100~~

<table style="width: 100%; border-collapse: collapse;"> <tr> <td style="width: 30%;">01210</td> <td style="width: 30%; border-left: 1px solid black;"> <table style="width: 100%; border-collapse: collapse;"> <tr><td>*00100</td><td rowspan="4" style="font-size: 2em; vertical-align: middle;">}</td><td rowspan="4">0.5</td></tr> <tr><td>*01110</td></tr> <tr><td>*01100</td></tr> <tr><td>*00110</td></tr> </table> </td> <td style="width: 30%; border-left: 1px solid black;"> <table style="width: 100%; border-collapse: collapse;"> <tr><td>00100</td><td>$\frac{1}{2} + \frac{1}{2} = 1$</td><td>$\rightarrow \frac{1}{6} \checkmark$</td></tr> <tr><td>01110</td><td>$\frac{1}{2}$</td><td>$\rightarrow \frac{1}{12}$</td></tr> <tr><td>01100</td><td>$\frac{1}{2}$</td><td>$\rightarrow \frac{1}{12}$</td></tr> <tr><td>00110</td><td>$\frac{1}{2} + \frac{1}{2} = 1$</td><td>$\rightarrow \frac{1}{6} \checkmark$</td></tr> <tr><td>10111</td><td>1</td><td>$\rightarrow \frac{1}{6} \checkmark$</td></tr> <tr><td>00111</td><td>1</td><td>$\rightarrow \frac{1}{6} \checkmark$</td></tr> </table> </td> </tr> <tr> <td>10222</td> <td style="border-left: 1px solid black;"> <table style="width: 100%; border-collapse: collapse;"> <tr><td>*10111</td><td rowspan="2" style="font-size: 2em; vertical-align: middle;">}</td><td rowspan="2">1</td></tr> <tr><td>*00111</td></tr> </table> </td> <td style="border-left: 1px solid black;"> <table style="width: 100%; border-collapse: collapse;"> <tr><td>00000</td><td>$\frac{1}{2} = \frac{1}{2}$</td><td>$\rightarrow \frac{1}{12}$</td></tr> <tr><td>00010</td><td>$\frac{1}{2} = \frac{1}{2}$</td><td>$\rightarrow \frac{1}{12}$</td></tr> </table> </td> </tr> <tr> <td>00110</td> <td style="border-left: 1px solid black;"> <table style="width: 100%; border-collapse: collapse;"> <tr><td>*00110</td><td rowspan="4" style="font-size: 2em; vertical-align: middle;">}</td><td rowspan="4">0.5</td></tr> <tr><td>*00000</td></tr> <tr><td>*00010</td></tr> <tr><td>*00100</td></tr> </table> </td> <td style="border-left: 1px solid black;"> <table style="width: 100%; border-collapse: collapse;"> <tr><td colspan="3" style="text-align: center;">+</td></tr> <tr><td colspan="3" style="text-align: center; font-size: 2em;">6</td></tr> </table> </td> </tr> </table>	01210	<table style="width: 100%; border-collapse: collapse;"> <tr><td>*00100</td><td rowspan="4" style="font-size: 2em; vertical-align: middle;">}</td><td rowspan="4">0.5</td></tr> <tr><td>*01110</td></tr> <tr><td>*01100</td></tr> <tr><td>*00110</td></tr> </table>	*00100	}	0.5	*01110	*01100	*00110	<table style="width: 100%; border-collapse: collapse;"> <tr><td>00100</td><td>$\frac{1}{2} + \frac{1}{2} = 1$</td><td>$\rightarrow \frac{1}{6} \checkmark$</td></tr> <tr><td>01110</td><td>$\frac{1}{2}$</td><td>$\rightarrow \frac{1}{12}$</td></tr> <tr><td>01100</td><td>$\frac{1}{2}$</td><td>$\rightarrow \frac{1}{12}$</td></tr> <tr><td>00110</td><td>$\frac{1}{2} + \frac{1}{2} = 1$</td><td>$\rightarrow \frac{1}{6} \checkmark$</td></tr> <tr><td>10111</td><td>1</td><td>$\rightarrow \frac{1}{6} \checkmark$</td></tr> <tr><td>00111</td><td>1</td><td>$\rightarrow \frac{1}{6} \checkmark$</td></tr> </table>	00100	$\frac{1}{2} + \frac{1}{2} = 1$	$\rightarrow \frac{1}{6} \checkmark$	01110	$\frac{1}{2}$	$\rightarrow \frac{1}{12}$	01100	$\frac{1}{2}$	$\rightarrow \frac{1}{12}$	00110	$\frac{1}{2} + \frac{1}{2} = 1$	$\rightarrow \frac{1}{6} \checkmark$	10111	1	$\rightarrow \frac{1}{6} \checkmark$	00111	1	$\rightarrow \frac{1}{6} \checkmark$	10222	<table style="width: 100%; border-collapse: collapse;"> <tr><td>*10111</td><td rowspan="2" style="font-size: 2em; vertical-align: middle;">}</td><td rowspan="2">1</td></tr> <tr><td>*00111</td></tr> </table>	*10111	}	1	*00111	<table style="width: 100%; border-collapse: collapse;"> <tr><td>00000</td><td>$\frac{1}{2} = \frac{1}{2}$</td><td>$\rightarrow \frac{1}{12}$</td></tr> <tr><td>00010</td><td>$\frac{1}{2} = \frac{1}{2}$</td><td>$\rightarrow \frac{1}{12}$</td></tr> </table>	00000	$\frac{1}{2} = \frac{1}{2}$	$\rightarrow \frac{1}{12}$	00010	$\frac{1}{2} = \frac{1}{2}$	$\rightarrow \frac{1}{12}$	00110	<table style="width: 100%; border-collapse: collapse;"> <tr><td>*00110</td><td rowspan="4" style="font-size: 2em; vertical-align: middle;">}</td><td rowspan="4">0.5</td></tr> <tr><td>*00000</td></tr> <tr><td>*00010</td></tr> <tr><td>*00100</td></tr> </table>	*00110	}	0.5	*00000	*00010	*00100	<table style="width: 100%; border-collapse: collapse;"> <tr><td colspan="3" style="text-align: center;">+</td></tr> <tr><td colspan="3" style="text-align: center; font-size: 2em;">6</td></tr> </table>	+			6		
01210	<table style="width: 100%; border-collapse: collapse;"> <tr><td>*00100</td><td rowspan="4" style="font-size: 2em; vertical-align: middle;">}</td><td rowspan="4">0.5</td></tr> <tr><td>*01110</td></tr> <tr><td>*01100</td></tr> <tr><td>*00110</td></tr> </table>	*00100	}			0.5	*01110	*01100	*00110	<table style="width: 100%; border-collapse: collapse;"> <tr><td>00100</td><td>$\frac{1}{2} + \frac{1}{2} = 1$</td><td>$\rightarrow \frac{1}{6} \checkmark$</td></tr> <tr><td>01110</td><td>$\frac{1}{2}$</td><td>$\rightarrow \frac{1}{12}$</td></tr> <tr><td>01100</td><td>$\frac{1}{2}$</td><td>$\rightarrow \frac{1}{12}$</td></tr> <tr><td>00110</td><td>$\frac{1}{2} + \frac{1}{2} = 1$</td><td>$\rightarrow \frac{1}{6} \checkmark$</td></tr> <tr><td>10111</td><td>1</td><td>$\rightarrow \frac{1}{6} \checkmark$</td></tr> <tr><td>00111</td><td>1</td><td>$\rightarrow \frac{1}{6} \checkmark$</td></tr> </table>	00100	$\frac{1}{2} + \frac{1}{2} = 1$	$\rightarrow \frac{1}{6} \checkmark$	01110	$\frac{1}{2}$	$\rightarrow \frac{1}{12}$	01100	$\frac{1}{2}$	$\rightarrow \frac{1}{12}$	00110	$\frac{1}{2} + \frac{1}{2} = 1$	$\rightarrow \frac{1}{6} \checkmark$	10111	1	$\rightarrow \frac{1}{6} \checkmark$	00111	1	$\rightarrow \frac{1}{6} \checkmark$																											
*00100	}	0.5																																																					
*01110																																																							
*01100																																																							
*00110																																																							
00100	$\frac{1}{2} + \frac{1}{2} = 1$	$\rightarrow \frac{1}{6} \checkmark$																																																					
01110	$\frac{1}{2}$	$\rightarrow \frac{1}{12}$																																																					
01100	$\frac{1}{2}$	$\rightarrow \frac{1}{12}$																																																					
00110	$\frac{1}{2} + \frac{1}{2} = 1$	$\rightarrow \frac{1}{6} \checkmark$																																																					
10111	1	$\rightarrow \frac{1}{6} \checkmark$																																																					
00111	1	$\rightarrow \frac{1}{6} \checkmark$																																																					
10222	<table style="width: 100%; border-collapse: collapse;"> <tr><td>*10111</td><td rowspan="2" style="font-size: 2em; vertical-align: middle;">}</td><td rowspan="2">1</td></tr> <tr><td>*00111</td></tr> </table>	*10111	}	1	*00111	<table style="width: 100%; border-collapse: collapse;"> <tr><td>00000</td><td>$\frac{1}{2} = \frac{1}{2}$</td><td>$\rightarrow \frac{1}{12}$</td></tr> <tr><td>00010</td><td>$\frac{1}{2} = \frac{1}{2}$</td><td>$\rightarrow \frac{1}{12}$</td></tr> </table>	00000	$\frac{1}{2} = \frac{1}{2}$	$\rightarrow \frac{1}{12}$	00010	$\frac{1}{2} = \frac{1}{2}$	$\rightarrow \frac{1}{12}$																																											
*10111	}	1																																																					
*00111																																																							
00000	$\frac{1}{2} = \frac{1}{2}$	$\rightarrow \frac{1}{12}$																																																					
00010	$\frac{1}{2} = \frac{1}{2}$	$\rightarrow \frac{1}{12}$																																																					
00110	<table style="width: 100%; border-collapse: collapse;"> <tr><td>*00110</td><td rowspan="4" style="font-size: 2em; vertical-align: middle;">}</td><td rowspan="4">0.5</td></tr> <tr><td>*00000</td></tr> <tr><td>*00010</td></tr> <tr><td>*00100</td></tr> </table>	*00110	}	0.5	*00000	*00010	*00100	<table style="width: 100%; border-collapse: collapse;"> <tr><td colspan="3" style="text-align: center;">+</td></tr> <tr><td colspan="3" style="text-align: center; font-size: 2em;">6</td></tr> </table>	+			6																																											
*00110	}	0.5																																																					
*00000																																																							
*00010																																																							
*00100																																																							
+																																																							
6																																																							

Answer: (C)

Actual answer

$$p' = p + \epsilon u \quad u \approx \nabla f = \left[\frac{df}{dp_1} - \frac{\sum_{i=1}^n \frac{df}{dp_i}}{n} \right]$$

$$\log_x L(p_1, \dots, p_n) = \sum_{i=1}^n C_i \log(p_i)$$

$$p_1 = p_2 = \dots = \frac{1}{n}$$

$$\frac{dL}{dp_1} = \frac{C_1}{p_1} = \frac{C_1}{1/n} = C_1 \cdot n$$

$$\sum_{i=1}^n \frac{C_i n}{n} = \sum_{i=1}^n C_i = C$$

$$p_i^{(1)} - p_j^{(1)} = \frac{1}{n} - \frac{1}{n} + \epsilon(C_i n - C)$$

6. Likelihood optimization

Consider the likelihood function of the haplotype frequency estimation problem with no missing data. In this setting, the parameters are the haplotype allele frequencies p_1, \dots, p_n , satisfying $\sum_{i=1}^n p_i = 1$, and $p_i \geq 0$, and the data is the haplotype counts c_1, \dots, c_n , and let $c = c_1 + \dots + c_n$. We showed in class that the likelihood function $L(p_1, \dots, p_n)$ satisfies

$$\log L(p_1, \dots, p_n) = \sum_{i=1}^n c_i \log(p_i) = \epsilon(C_i n - C) = \epsilon(C_i - C_j)$$

Assume we are running a gradient ascent with projections algorithm, and we start from the guess $p_1^{(0)} = p_2^{(0)} = \dots = p_n^{(0)} = \frac{1}{n}$. Let $p_1^{(1)}, \dots, p_n^{(1)}$ be the next point that the algorithm reaches. Which of the following is true:

- (a) There is $\epsilon > 0$ such that for every i, j we have $p_i^{(1)} - p_j^{(1)} = (c_i - c_j)\epsilon$.
- (b) There is $\epsilon > 0$ such that for every i, j we have $p_i^{(1)} - p_j^{(1)} = (c_j - c_i)\epsilon$.
- (c) There is $\epsilon > 0$ such that for every i we have $p_i^{(1)} = (c_i - \frac{c}{n})\epsilon$.
- (d) There is $\epsilon > 0$ such that for every i, j we have $p_i^{(1)} + p_j^{(1)} = \frac{2}{n} + (c_i + c_j)\epsilon$.

Let current point = Cp

Select from the above options and justify briefly in the box why you choose your answer.

Let $g(p) = \sum_{i=1}^n c_i \log(p_i)$, s.t. $\sum p_i = 1, p_i \geq 0$
 $CP: p = (p_1, \dots, p_n)$
 $u = \left(\frac{dg}{dp_1} - \lambda, \dots, \frac{dg}{dp_n} - \lambda \right) \quad \frac{dg}{dp_i} = \sum \frac{c_i}{p_i}$

For small $\epsilon > 0$:

$$p' = p + \epsilon u = (p_1, \dots, p_n) + \epsilon \left(\frac{c_1}{p_1} - \lambda, \frac{c_2}{p_2} - \lambda, \dots, \frac{c_n}{p_n} - \lambda \right)$$

where $\lambda = \frac{\sum_{i=1}^n \frac{dg}{dp_i}}{n}$

$$p' = \left(p_1 + \epsilon \left(\frac{c_1}{p_1} - \frac{\sum_{i=1}^n \frac{dg}{dp_i}}{n} \right), \dots, p_n + \epsilon \left(\frac{c_n}{p_n} - \frac{\sum_{i=1}^n \frac{dg}{dp_i}}{n} \right) \right) = (p_1^{(1)}, \dots, p_n^{(1)})$$

For each i, j : $p_i^{(1)} - p_j^{(1)} = \frac{p_i^{(0)}}{1/n} + \epsilon \left(\frac{c_i}{p_i^{(0)}} - \frac{\sum_{i=1}^n \frac{dg}{dp_i}}{n} \right) - \left[\frac{p_j^{(0)}}{1/n} + \epsilon \left(\frac{c_j}{p_j^{(0)}} - \frac{\sum_{i=1}^n \frac{dg}{dp_i}}{n} \right) \right]$

$$\left[\frac{p_i^{(0)}}{1/n} + \epsilon \left(\frac{c_i}{p_i^{(0)}} - \frac{\sum_{i=1}^n \frac{dg}{dp_i}}{n} \right) \right] - \left[\frac{p_j^{(0)}}{1/n} + \epsilon \left(\frac{c_j}{p_j^{(0)}} - \frac{\sum_{i=1}^n \frac{dg}{dp_i}}{n} \right) \right]$$

$$\approx \epsilon(C_i - C_j) = \epsilon(C_i - C_j)$$