Name: Chan Kei Lam SID: 1155221419

Data & Data types (20 pts)- done

1a. Text

1b. convert to data matrix

Student	Gender	Height (cm)	Weight (kg)
1	Male	180	75
2	Female	165	50
3	Female	164	49

2.

Sequential data (associated with regression) – prediction of mutations in infectious microbiome

Networks- to show the correlation between gene sequence and survival rate under different types of antibiotics

Basic concept of Python programming (20 pts)- done

- 1a. import numpy which is an additional plug-in for Python
- 1b. Calculate the mean of the list [1,2,3] using numpy and assign the result to variable x
- 1c. Print the string 'x is ' followed by the value of x
- 2a. import numpy
- 2b. z = numpy.array([1, 2, 3])

2c.

mean	mean_z = numpy.mean(z)
max	max_z = numpy.max(z)

min	min_z = numpy.min(z)
SD	std_z = numpy.std(z)

Sequence alignment & Dynamic programming (20 pts)-done

1a. (A-C), (C-C), (C-T), (G-_)
$$\rightarrow$$
 (-7) + 2 + (-5) + (-10) = -20

1b. (A-_), (C-C), (C-C), (G-T)
$$\rightarrow$$
 (-10) + 2 + 2 + (-7) = -13

1c. Since -13 > -20... Thus, sequence in (a.) would not be the best alignment we are looking for.

2. The best alignment score (similarity) between the two sequences is -10.

Best alignment: F(ACCG, CCT): -13 = (-10) + 2+2+(-7)

	_	A	C	C	G
_	0 →	-10	-20	-30	-40
	v				
С	-10	-7	-8	-18	-28
С	-20	-17	-5	-6	-16
T	-30	-27	-15	-10	-13

Sequence assembly & Sequence mapping (20 pts) -un

1. CGA one repeated, cannot be read: GGA, ATCG, ATC

Т	Α	Α	Т	G	С	G	Α	Т	G	G	С	Т	G	G	G	T	T	Α	Α	T
Т	Α	Α																		
				G	С	G														
					С	G	Α													
					С	G	Α													
						G	Α	Т												
					С	G	Α	Т												
													G	G	G					
														G	G	T				
													G	G	G	T				
															G	T	T			

2. gene expression matrix

Sample	count number		
Red	4 (or 5 with one repeated)		
Blue	4		

Data cleaning & Data visualization (20 pts) done

1. Sample-wise, to remove difference due to biological and technological variations, e.g 10000000

merger

Sample name	Gene- 1	Gene-2	Gene-3
Sample 1	679	-	260
Sample 2	448	515	211
Sample 3	873	621	-
Sample 4	408	365	164
Sample 5	401	499	202
Sample 6	800	605	699

.

2. sample-wise, as we want to see the difference between gene expressions under similar conditions more than the difference induced by samples, e.g. biological and technological variations

For the missing data:

Mean of Gene-2 (except Sample-1):
$$(515+621+365+499+605)/5 = 521$$

Mean of Gene-3: $(260+211+164+202+699)/5 = 307.2$

١

Sample name	Gene- 1	Gene-2	Gene-3
Sample 1	679	- → 521	260
Sample 2	448	515	211

Sample 3	873	621	- → 307.2
Sample 4	408	365	164
Sample 5	401	499	202
Sample 6	800	605	699

3. Use Min-max normalization: v-Vmin / vmax-vmin

Sample name	Calculation
Sample 1	Gene-2: (521 - 260) / (679 – 260) ≈ 0.623
Sample 2	Gene-1: (448 - 211) / (515 - 211) ≈ 0.780
Sample 3	Gene-2: (621 - 307.2) / (873 - 307.2) ≈ 0.555
Sample 4	Gene-2: (365 - 164) / (408 - 164) ≈ 0.824
Sample 5	Gene-1: (401 - 202) / (499 - 202) ≈ 0.670
Sample 6	Gene-2: (699-605) / (800 - 605) ≈ 0.482

Last:

Sample name	Gene- 1	Gene-2	Gene-3
Sample 1	1	0.623	0
Sample 2	0.780	1	0
Sample 3	1	0.555	0
Sample 4	1	0.824	0
Sample 5	0.670	1	0
Sample 6	1	0	0. 482