

Integrating Biology and Statistics at the freshman level-SYMBIOSIS I Exploring DNA sequences using probability

ETSU



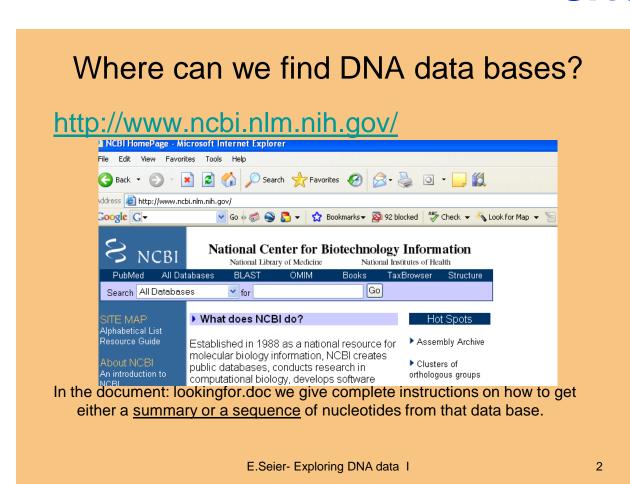
HHMI Grant # 52005872

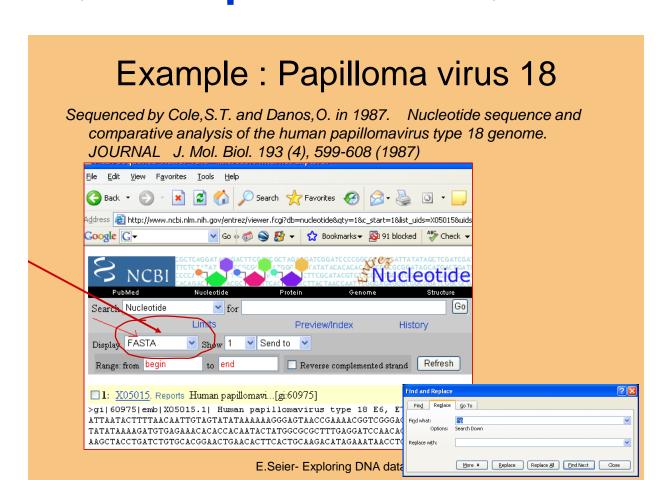
Edith Seier & Karl Joplin seier@etsu.edu; joplin@etsu.edu East Tennessee State University

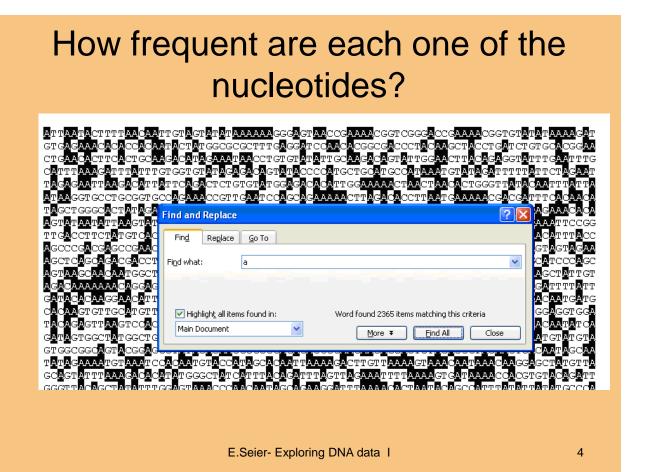
Symbiosis I is the first of a sequence of three Biology, Mathematics and Statistics integrative courses. Students who pass Symbiosis I get credit for Biology I and Introductory Statistics. Module 5 of Symbiosis I covers the basics of DNA genetics from the biological point of view. The module takes DNA as a sequence of the 4 nucleotides C,A,G, T; questions about sequences are posed and the probability tools to answer them are developed or reviewed. This module prepares the students, at an elementary level, for the future study of bioinformatics. At the same time students put into practice their knowledge of Probability in a relevant context.

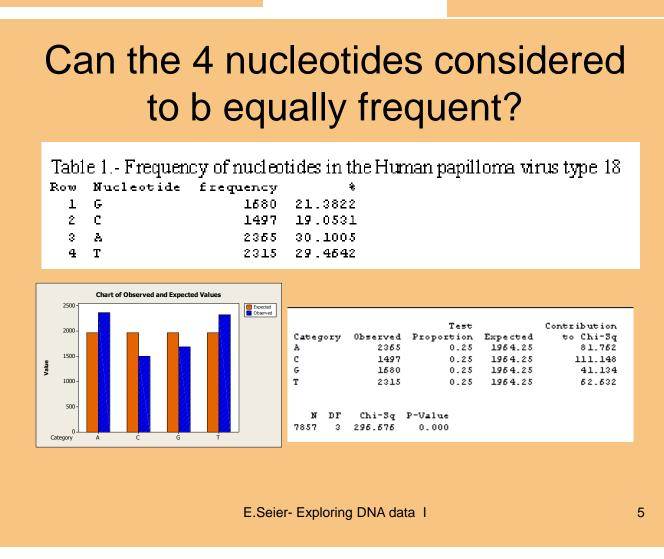
From the mathematical point of view the only pre-requisites are the basic understanding of probability, independence and conditional probability, which are covered in Module 4. The topics covered are: DNA as nucleotide sequences, nucleotide frequency, GC content. Independence and conditional probability in the DNA environment. Transition matrix, graph to represent transition matrices. Probability of a given sequence of nucleotides, repeats of a single nucleotide, length of the repeat, geometric distribution. Palindromes, probability of any palindrome and of specific palindromes, space in between palindromes. Comparing two sequences of nucleotides. Similarities that happen just by chance. Random walks (and their use in testing for similarities). Part of the teaching material created is displayed below. Microsoft word and free software available in the internet are used in the exploration of sequences.

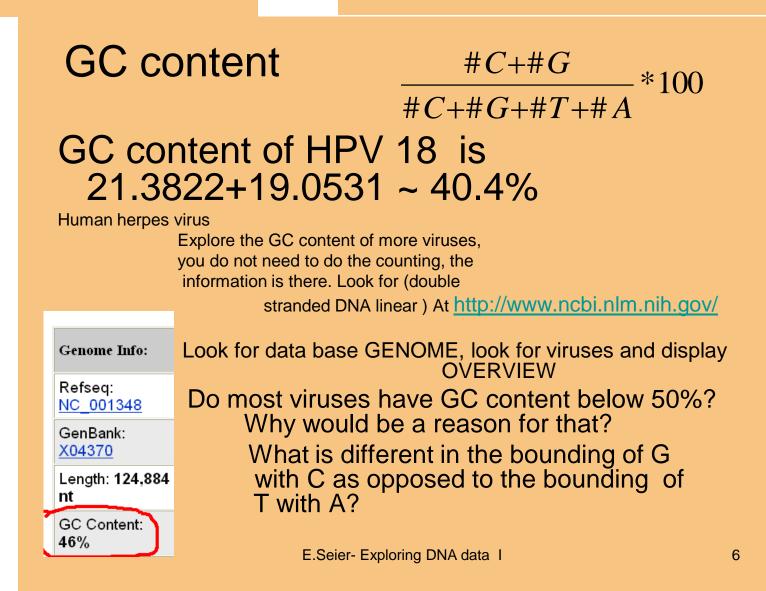
Data bases, frequencies, GC content











Dinucleotides frequency • These are the first 100 nucleotides of the BRCA1 gene

CTTTGTGCAACAGTACTTTCCCAGGATCCACA GGAAATACTCAGAGTCCACCTGGACATTTTA CTTATATTCAGTTTCCAAGTGTCAGAGGGGT **TCAGGA** How many times you find the pair CG? How many times you find the pair CA? How many times you find the pair CT? How many times

you find the pair CC?

E.Seier- Exploring DNA data I

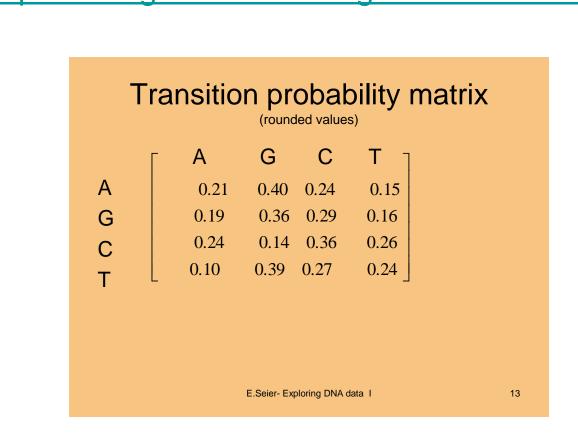
Conditional probabilities P(G|C)=Probability of a G given Overall the that the previous letter was a C proportion of G is 31%, so if we pick one location at random the probability that it is a G is 0.31 but if we know that the previous one is a C, the $P(C \text{ first and G next}) = \overline{100718}$ probability is only

E.Seier- Exploring DNA data I

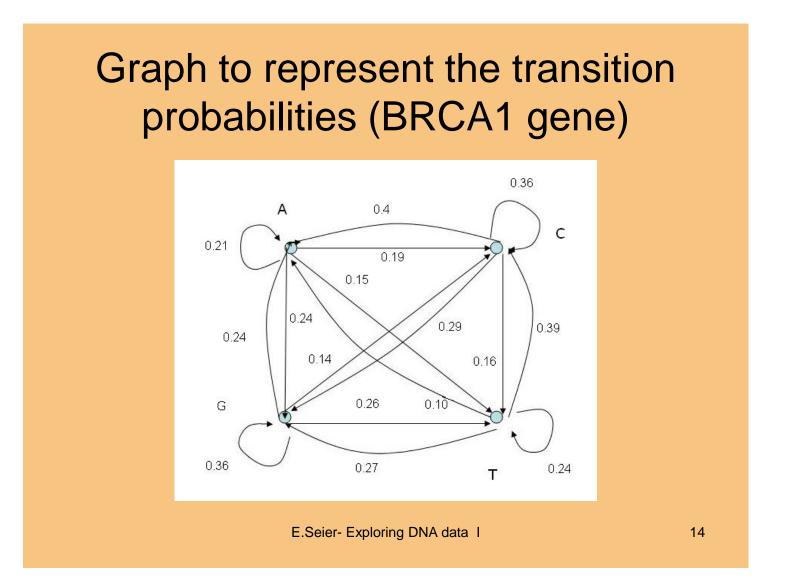
Dinucleotides

Dinucleotide frequency for the whole BRCA1 gene				
F	Row DN	frequency P(first) P(second) Expected ratio	· · · · · · · · · · · · · · · · · · ·
	1 AG	7706 0.19		P(G)=0.31 P(C)=0.295 P(A)=0.19 P(T)=0.205
	2 CT	7868 0.295		There are 100719 nucleotides in
	3 CA	7098 0.29		BRCA1 gene sequence. Thus there
	4 TG	7997 0.20		are 100718 pairs of consecutive
	5 CC	10615 0.29		letters.
	6 GG	11388 0.31		The frequencies of each 'dinucleotide' are displayed in the
	7 Π	4828 0.205		table.
	8 AA	4048 0.19		Why the frequencies differ from what
	9 GA	5901 0.31		we would expect assuming
	10 GC	8935 0.31		independence?
	11 TC	5627 0.20		Portions of the DNA (coding region) gets 'transcripted' into mRNA (T is
1	12 AC	4557 0.19		replaced by U) and mRNA is
	13 GT	5020 0.31		replaced by U) and mRNA is 'translated' into aminoacids (which in
	14 AT	2856 0.19		turn go into proteins)
	15 TA	2120 0.20		Thus the 'writing' of aminoacids would put conditions into the pairs of
	16 CG	4154 0.29		consecutive letters or 'dinucleotides'
	TOTAL	100718		as opposed as if DNA was just a
E.Seier- Exploring DNA data randomly formed sequence of 4 ₁₀				

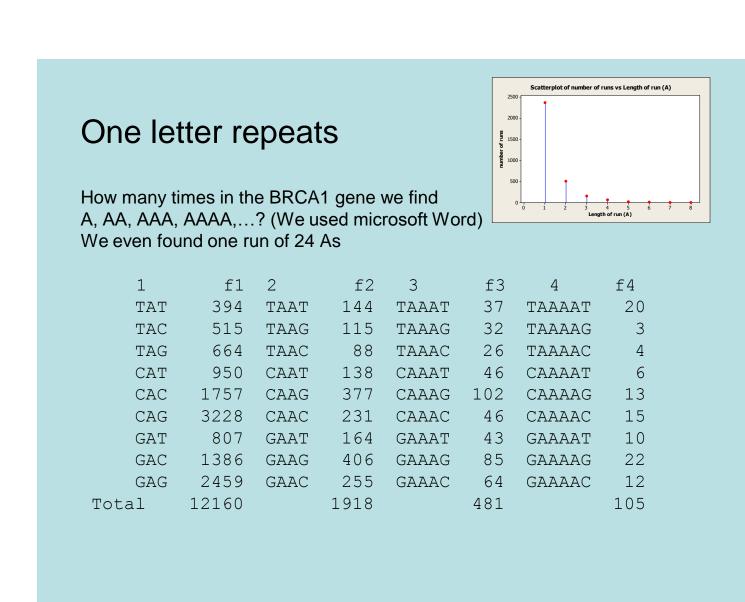
Software useful to count frequencies of nucleotides, dinucleotides and trinucleotides http://www.genomatix.de/cgi-bin/tools/tools.pl

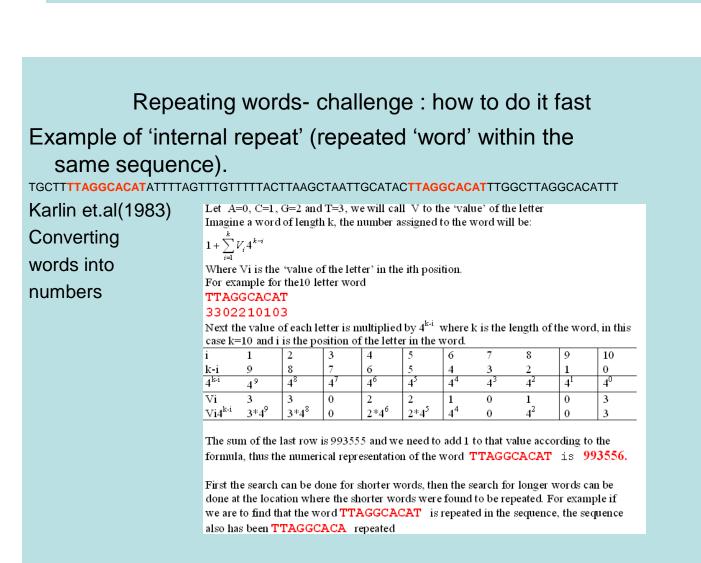


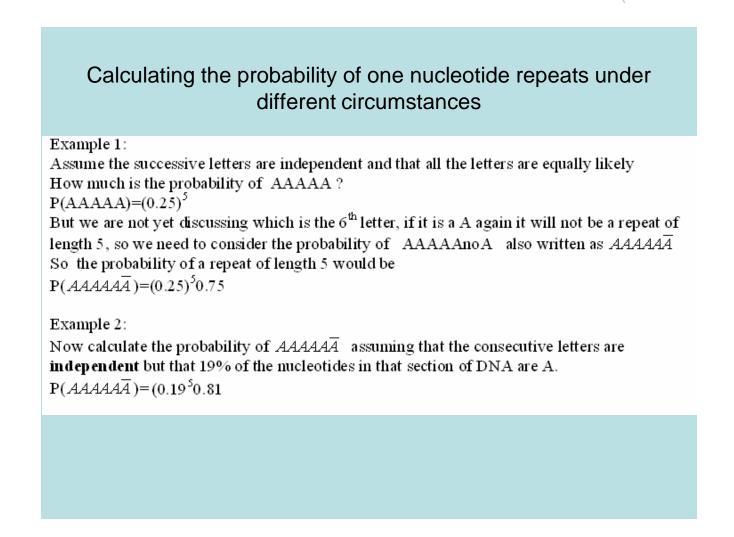
Are the consecutive letters independent? There are 100719 nucleotides in that sequence I counted how many times CC,CG,CA,CT, GC,GG,GA,GT,AC,AG,AA,AT, TC,TG,TA or TA happened. The frequencies of each pair are displayed in the table below: A 4048 7706 4557 2856 19167 G 5901 11388 8935 5020 31244 C 7098 4154 10615 7868 29735 T 2120 7997 5627 4828 20572 Total 100718 pairs of letters E.Seier- Exploring DNA data

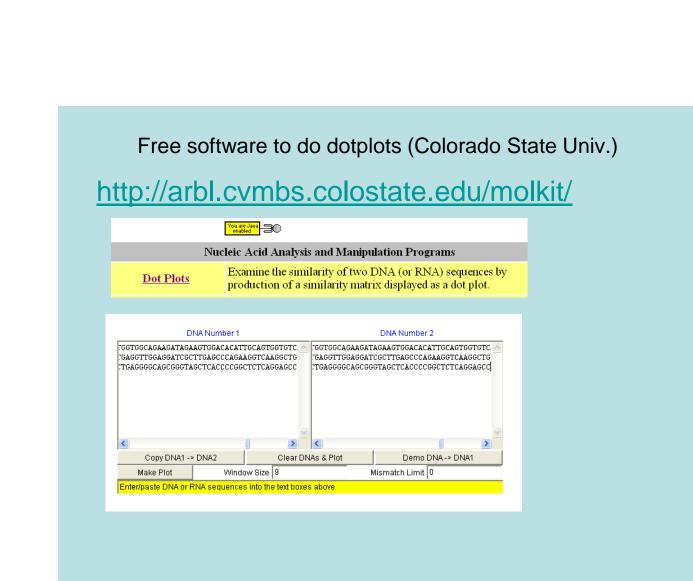


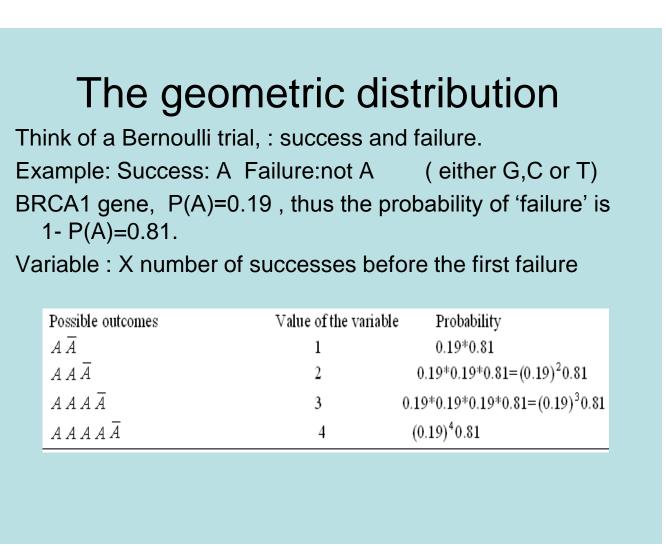
Repeats

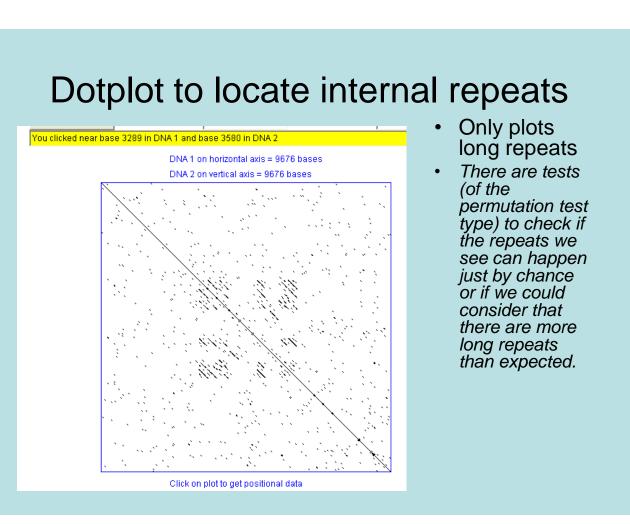


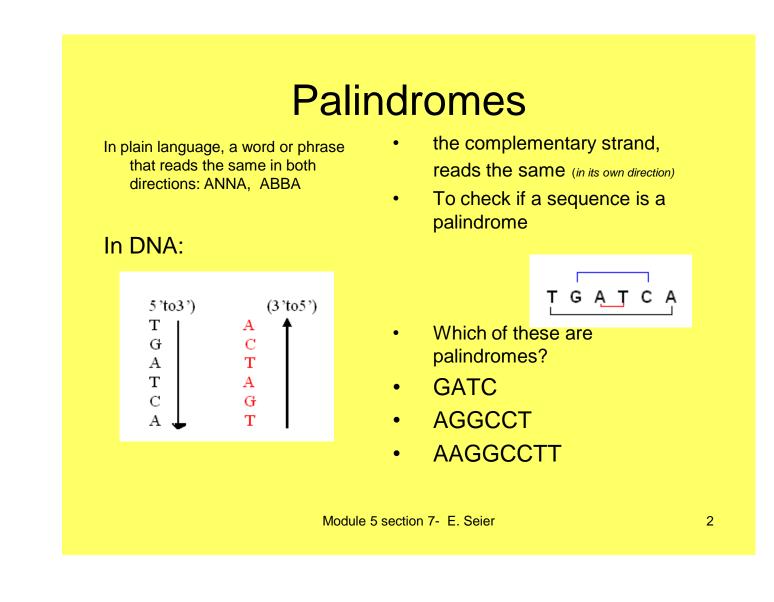


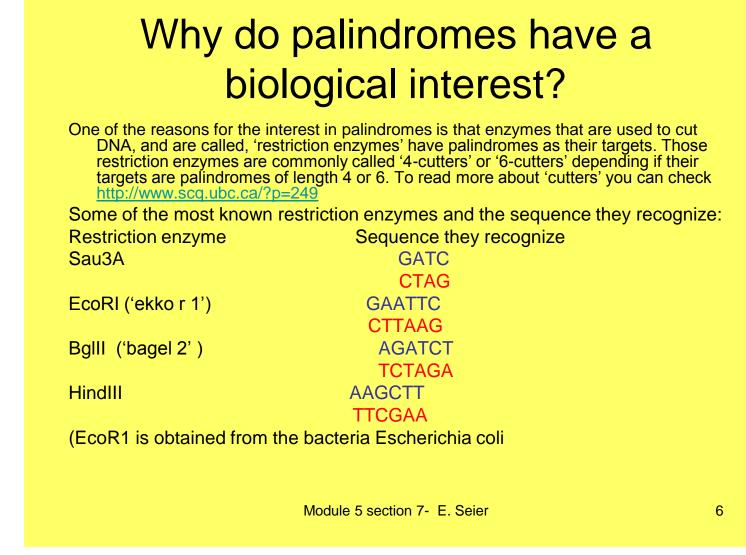












Palindromes

