

Biological words: (in genome)

1. $R=1$ (base composition): $\{A, T, G, C\}$
 $fr(A), fr(B), fr(G), fr(T) \Rightarrow$ Any one parameter is unk.

$\therefore fr(A) = fr(T), \quad \& \quad fr(C) = \frac{(1 - fr(A+T))}{2}$
 $= fr(G).$

<u>Org.</u>	<u>% G+C</u>	<u>Genome size (Mb)</u>
<u>Eubacteria:</u>		
E-coli	50.7	4.693
<u>Archaeobacteria</u>		
Pyrococcus abyssi	44.6	1.765
<u>Eukaryotes</u>		
C-elegans (Nematodes)	36	97
Arabidopsis (Thalassia (a flower's plant))	35	125
Homo Sapiens	41	3080.

In prokaryotic genomes, in particular, there is an excess G over C on the leading strands, (strands whose 5' to 3' direction corresponds to the direction of replication fork movement). \Rightarrow Known as "GC skew" $= \frac{\#G - \#C}{\#G + \#C}$

2. $R=2 \quad \{AA, AT, \dots, CC\}$

χ^2 - Test for verifying the observation with the model (or expected value). (Here "dinucleotide of freq").

Let $O =$ # dinucleotide (say "GC") in seq. of length n .

$E =$ expected no. of $\binom{(n-1) P_G P_C}{(n-1) \text{ pairs}}$

$\therefore \chi^2 = \frac{(O - E)^2}{E}$

$\chi^2 \Rightarrow$ small if $O \approx E$ else large.

Alg.: Let r_1 & r_2 be two nucleotides whose adjacent occurrence to be observed.

(a) Let $c = \begin{cases} 1 + 2Pr_1 - 3Pr_1^2 & r_1 = r_2 \\ 1 - 3Pr_1Pr_2 & r_1 \neq r_2 \end{cases}$

(b) If $\chi^2/c > 3.84^*$, conclude that the iid model (identically independent distribution) is not a good fit.
 [Statistics]
 from theory of

[Assignment to students]
 (* Accepted with 0.95 confidence level or ~~0.05~~ 5% level of significance)

$k = 3$ (Codons)

Codon Adaptation Index (for an amino acid):

$$CAI = \prod_{k=1}^L (p_k/q_k)^{1/L}$$

L = length of the amino acid.

p_k = Prob. of the k th ^{k th codons occurrence of the} amino acid (among highly expressed genes).

q_k = max. prob. of a codon of that amino acid (among the same set of highly expressed genes).

$$\log(CAI) = \frac{1}{L} \sum_{k=1}^L \log(p_k/q_k)$$

$0 < CAI \leq 1$. In E-Coli, a sample of 500 - protein coding genes displayed CAI values in the range of 0.2 to 0.85.

CAI vs +vely correlated with mRNA expression level.

k-word problem (a different approach)

Obj.: To find a particular 'k-word' occurs more often in a set of N strings (e.g. promoter regions) of length m (say).

Statistical Test:

$N_w =$ No. of strings where the word "w" is observed.

Compute $P_w =$ Prob. of occurrence of w ~~at~~ at least once in a seq. of length m .

[simulate by randomly generating string of length m for a large no. of times (say 5000) using the prob. of letters (A, T, C, G) in a genome (computed from the seq.)].

\therefore The following ~~is~~ ^{the} ~~used~~ randomness model given N strings:

$$P(k) = \text{Prob. of the string occurring } k \text{ times} \\ = \binom{N}{k} p_w^k (1-p_w)^{N-k}$$

$$E(k) = N \cdot p_w \quad \& \quad \text{Var.}(k) = N \cdot p_w (1-p_w)$$

Then the statistic

$$Z_w = \frac{N_w - N \cdot p_w}{\sqrt{N p_w (1-p_w)}} \approx N(0, 1)$$

Compute p-value ^{via} $P(Z_w > Z_w)$ for each word.

If $P(Z_w) < 0.001$, do not consider the string w as unusual.