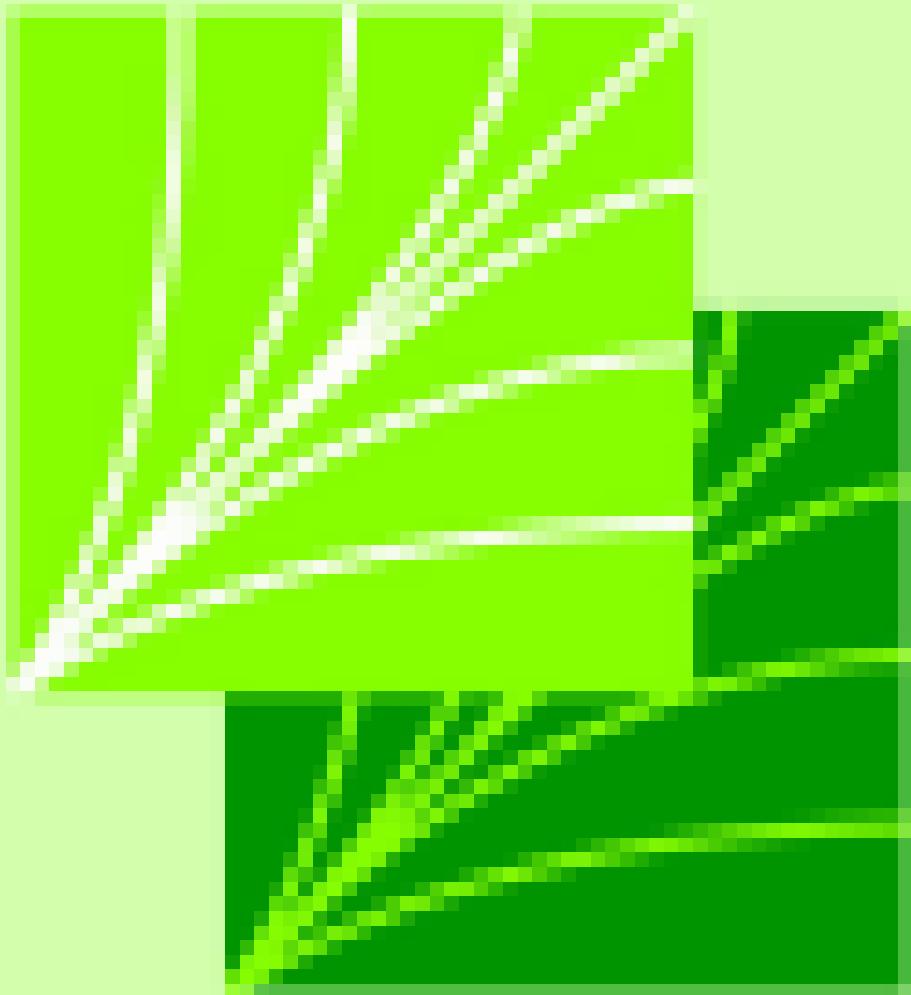


Mathematics of
Genomic Analysis

Mathematics of Genomic Analysis



This page intentionally left blank

MATHEMATICS OF GENOME ANALYSIS

The massive research effort known as the Human Genome Project is an attempt to reveal the sequence of the three billion nucleotides that make up the human genome and to identify individual genes within this sequence. Although the basic effort is of course a biological one, the description and classification of sequences also naturally lend themselves to mathematical and statistical modeling.

This short, but basic, on the mathematics of genome analysis presents a brief description of several ways in which mathematics and statistics are being used in genome analysis and sequencing. It will be of interest not only to students but also to professional mathematicians curious about the subject.

Jerome K. Pavao is Professor of Physics and Mathematics at the Courant Institute of Mathematical Sciences and Department of Physics at New York University, where he has taught since 1986. He has held visiting positions at Middelton Hospital Medical School in London, Columbia University, Rutgers University, Princeton University, Rockefeller University, Yukawa Institute in Kyoto, Tohoku University, Norwegian Institute of Technology, Max Planck Institute in Tübingen, Catholic University in Rio de Janeiro, Ecole Polytechnique de Lille, Soviet Academy of Sciences in Moscow, Liverpool, Kirov, and Leningrad, University of Paris, Nankai University, and Tsinghua University in China. He has received the Peleg (New York Academy of Science), Pattern Recognition Society and Hildebrand (American Chemical Society) Chemical Physics awards.

Cambridge Studies in Mathematical Biology

Editor

C. CHALMERS

University of Sheffield, UK

F. C. HOPPENSTEADT

Akiozo State University, Tempe, AZ, USA

L. A. MILNE

Wileyana Institute of Science, Tel Aviv, Israel

1. Brian Charlesworth, *Evolution in age-structured populations* (Cambridge)
2. Stephen Coarman, *Mathematics of learning and living*
3. C. Chalmer and E. A. Thompson, *Geotropophiles and parasitic nematodes*
4. Frank C. Hoppensteadt, *Mathematical methods of population biology*
5. G. Clark and B. S. Everitt, *An introduction to mathematical genetics*
6. Frank C. Hoppensteadt, *An introduction to the mathematics of neurons* (Cambridge)
7. Jane Cronin, *Mathematical aspects of Wright's Shapley model theory*
8. Henry C. Winkler, *Formulations in theoretical neurobiology*
William J. Adams, editor. *Neurobiology and dendrite structures*
9. William J. Adams, editor. *Neurobiology and stochastic processes*
10. H. Waddington, *Bioepigenetic systems*
11. Anthony G. Pakes and R. A. Miller, *Mathematical ecology of plant species competition*
12. Eric Bonham, *Modelling biological populations in space and time*
13. Lee A. Segel, *Biophysical kinetics*
14. Hal L. Smith and Paul Waltman, *The theory of the chemostat*
15. Brian Charlesworth, *Evolution in age-structured populations* (Cambridge)
16. D. J. Daley and D. G. Kendall, *Epidemic modelling - an introduction*
17. J. Maynard Smith, *An introduction to mathematical physiology and bioLOGY* (Cambridge)
18. Andrew R. Barron, *Mathematics of genome analysis*

MATHEMATICS OF GENOME ANALYSIS

JEROME K. PIERCE

New York University



ISSN 0022-215X (print version) or ISSN 1365-2753 (electronic version)
The Pitt Building, Trumpington Street, Cambridge, United Kingdom

CAMBRIDGE UNIVERSITY PRESS

The Edinburgh Building, Cambridge CB2 2RU, UK
32 West 45th Street, New York, NY 10036-4211, USA
677 Williamstown Road, Port Melbourne, VIC 3207, Australia
Ruiz de Alarcón 13, 28014 Madrid, Spain
Dux Road, The Woodlands, Cape Town 7400, South Africa

http://journals.cambridge.org

© Cambridge University Press 2004

First published in printed format 2004

ISBN 0-521-80855-3, eBook (ebook)

ISBN 0-521-58817-1, hardback

ISBN 0-521-58818-9, paperback

Contents

Page ix.	Page
	Preface
1	Decomposing DNA
1.1	DNA Sequences
1.2	Restriction Fragments
1.3	Clone Liberation
	Assignment 1
2	Recomposing DNA
2.1	Progenitor Assembly
2.2	Anchoring
2.3	Restriction Fragment Length Polymorphisms (RFLP) Analysis
	Assignment 2
2.4	Pooling
	Assignment 3
2.5	Replics
3	Sequence Statistics
3.1	Local Properties of DNA
3.2	Long-Range Properties of DNA
3.2.1	Longest Repeat
	Assignment 4
3.2.2	Dispersion Correlations
3.2.3	Markov-Level Criteria
3.2.4	Batch-Level Criteria
3.2.5	Statistical Models
	Assignment 5
3.3	Other Measures of Significance

3.3.1 Survival Analysis	73
3.3.2 Estrogen Criteria	76
4 Sequence Comparisons	81
4.1 Basic Matching	81
4.1.1 Mutual Evaluation Model	82
4.1.2 Independence Model	85
4.1.3 Direct Approach Evaluation	87
4.1.4 Entropy-Value Technique	89
Assignment 6	90
4.2 Matching with Implications	91
4.2.1 Score Distribution	92
4.2.2 Penalty-Free Limit	97
4.2.3 Effect of Total Penalty	99
4.2.4 Score Acquisition	100
4.3 Multisequence Comparisons	101
4.3.1 Locating a Common Pattern	103
4.3.2 Assessing Significance	109
Assignment 7	117
4.3.3 Category Analysis	118
4.3.4 Adaptive Techniques	121
5 Spatial Structure and Dynamics of DNA	128
5.1 Thermal Behavior	129
5.2 Dynamics	131
5.3 Effect of Heterogeneity	137
Assignment 8	137
Bibliography	139
Index	137

Preface

"What is life?" is a perennial question of transcendent importance that can be addressed at a bewildering set of levels. A quantitative scientist, met with such a question, will tend to adopt a reductionist attitude and first seek the discernible units of the system under study. These are, to be sure, molecules, but it has become clear only in recent decades that the "Rosetta" molecules share the primary structure of linear sequences – in accord with a temporal sequence of construction – subsequent to which chemical bonding as well as exchange can both extend the sequence meandering fully in real space and create a much larger population of molecular species. At the level of sequences, this characterization is, not surprisingly, an oversimplification because, overwhelmingly, the construction process of a life form proceeds via the linear sequences of DNA, then of RNA, then of protein, on the way to an explosion of types of molecular species. The smaller population of life subsequences is certainly RNA, which is the principal focus of our study, but not – from an informational viewpoint – to the exclusion of the proteins that serve as ubiquitous enzymes, as well as messengers and structural elements; the fascinating story of RNA, will be referred to only briefly.

That the molecules we have to deal with are fundamentally describable as ordered linear sequences is a great blessing to the quantitatively minded. Methods of statistical physics are particularly adept at treating such entities, and information science – the implied content of the bulk of our considerations – is also most comfortable with these objects. This hardly translates to triviality, as a moment's reflection on the structure of human language will make evident.

In the hyperactive field that "informatics" has become, the term evolves very rapidly, and "traditional" may refer to activities two or three years old. I first presented the bulk of this material to a highly heterogeneous class in 1993, again with modification in 1995, and once more, further modified, in 1997. The always new set forth the mathematical framework in which the burgeoning activity takes place, and, although hardly impervious to the passage of time,

This approach imparts a certain amount of stability to an intrinsically unstable divergent structure. I do, of course, take advantage of this nominal stability, leaving it to the reader to consult the numerous technical journals, as well as with due caution the increasing flood of semitechnical articles that document the important emerging facets of the overall generative field.

It is a pleasure to acknowledge the help of Connie Fugate and Daisy Calderon-Mejia in converting a largely illegible manuscript to, it is hoped, readable form, of Professor Gino Puccetti for translating, or noticing, the non-English words which the original manuscript abounded, and of numerous students who not only maintained intelligent faces, but also gently pointed out instances of confusion in the original features.

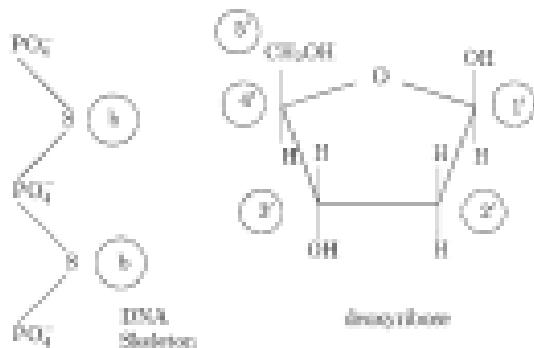
I

Decomposing DNA

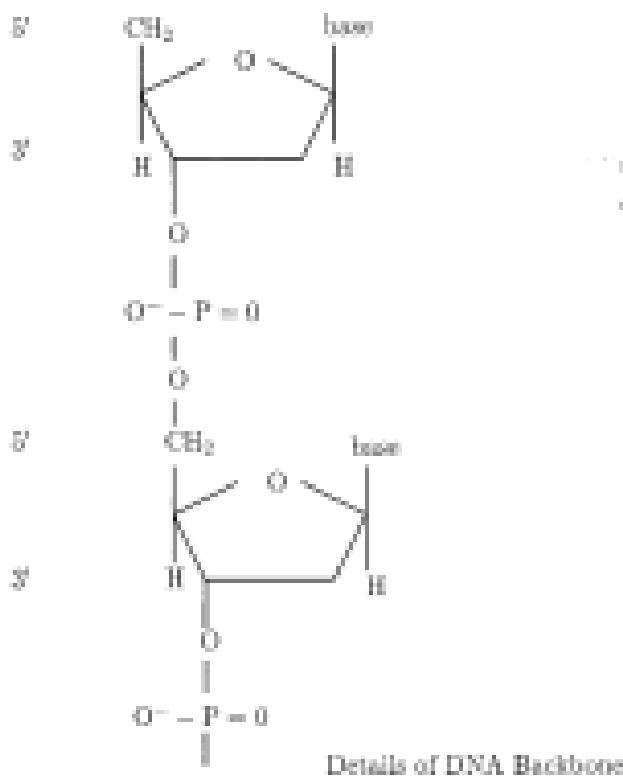
1.1. DNA Sequences

The realization that the genetic blueprint of a living organism is recorded in its DNA molecules developed over more than a century – slowly on the scale of the lifetime of the individual, but instantaneously on the scale of societal development. Divining the fashion in which this information is used by the organism is an enormous challenge that promises to dominate the life sciences for the foreseeable future. A crucial preliminary is, of course, that of actually compiling the sequence that defines the DNA of a given organism, and a fair amount of effort is devoted here to examples of how this has been and is being accomplished. We focus on nuclear DNA, ignoring the minuscule mitochondrial DNA.

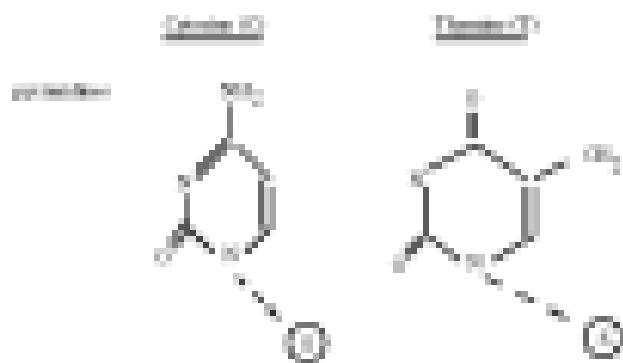
To start, let us introduce the major actor in the current show of life, the DNA chain, a very long polymer with a high degree of commonality – 99.8%, to within rearrangement of sections – among members of a given species [see Alberts et al. (1989) for an encyclopedic account of the biology, Cooper (1992) for a brief version, Miura (1986), and Gindkin (1992) for brief mathematical overviews]. The backbone of the DNA polymer is an alternating chain of phosphate (PO_4^-) and sugar (S) groups. The sugar is deoxyribose (an unmarked vertex in its diagrammatic representation always



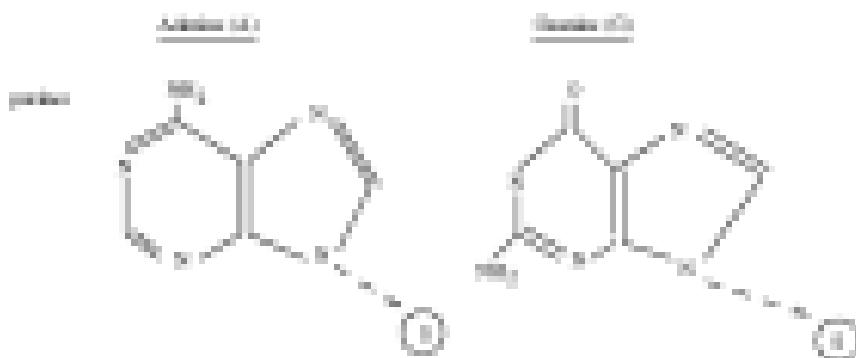
signifies a carbon atom) with standard identification of the five carbons as shown. Successive sugars are joined by a phosphate group (phosphoric acid, $H_2PO_4^-$, in which we can imagine that two hydrogens have combined with $3'$ and $5'OH$ s groups of the sugar, with the elimination of water, whereas one hydrogen has disappeared to create a negative ion); the whole chain then has a characteristic $3'-5'$ orientation (left to right in typical diagrams, corresponding to the direction of "reading," also upstream to downstream). However, the crucial components are the side chains or bases



(attached to $C1'$ of the sugar, again with elimination of water) of four types. Two of these are pyrimidines, built on a six-member ring of four carbons and two nitrogens (single and double bonds are indicated, carbons are implicit at line junctions). Note: Pyrimidine, cytosine, and thymine all have the letter y .



Two are the most biologically, built on joined five- and six-membered rings (adenine, with empirical formula $\text{H}_6\text{C}_5\text{N}_5$, and thymine the thieno-furan name pentahydrogen cyanate, of possible evolutionary significance).



DNA chains are normally present as pairs, in the famous Watson-Crick double-helix conformation, enhancing their molecular integrity. The two strands are bound through pairs of bases, pyrimidines to purines, by means of hydrogen bonds (...,), and chemical fitting requires that A must pair with T, G with C; thus each chain uniquely determines its partner. The DNA "alphabet" consists of only the four letters A, T, G, and C, but the full text is very long indeed, some 3×10^9 base pairs in the human. Roughly 1% of our DNA four-letter information is allocated to genes, "words" that translate into the proteins that, among other activities, create the enzymatic machinery that drives biochemistry, as well as instructional elements, the remaining unknown - perhaps mechanical - function.

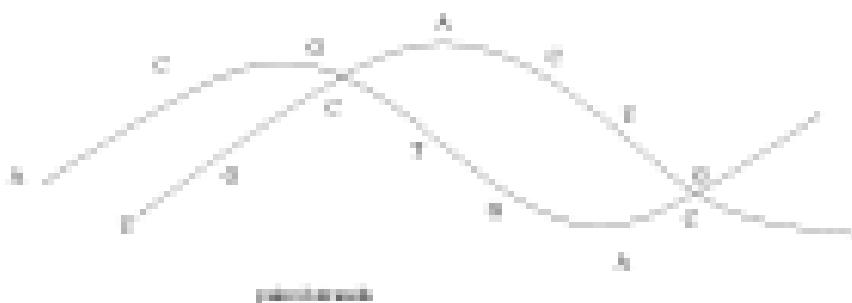


Double-chain DNA is typically represented in linear fashion, e.g.,

$$\begin{array}{ccccccccc} S & = & A & - & C & - & G & - & T \\ & & \vdots & & \vdots & & \vdots & & \vdots \\ Y & = & T & - & G & - & C & - & A \end{array}$$

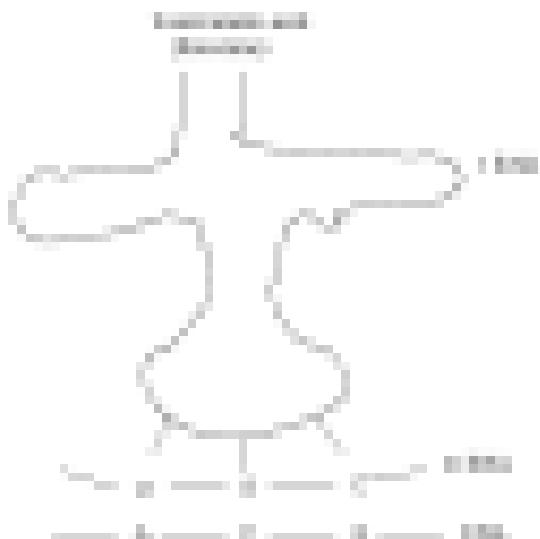
$$= C = G = T = A = C = G = T$$

(although the unique base pairing means that any single 2'-2' chain suffices), but because of the offset between 3' and 5' positions, the spatial structure is that of a spiral ribbon.



Thus the small portions of DNA – the genes – that code for proteins are not present in every part of the cell, but, especially in the non-gene-coding nucleotides, are interrupted by intervening (and often highly repetitive) DNA. The coding fragments – or exons – are also flanked by instructional subsequences, so that a small gene might look like: (5') upstream enhancer, promoter, start site, exon, intron, exon, poly-A site, stop site, downstream enhancer (3'). However, the vast remaining "junk DNA" – also richified by fairly complex repeats (e.g., 300 base pairs, 1.1, very long, microsatellites, very short) – adds greatly diverse mechanical properties. Indeed, e.g., the supercoiled structure grafted onto the double helix, is of unknown function, and may be only an evolutionary relic.

The single step inhibitor D906 → protein response overall stalled. Removal of the inhibitor allows the same four-step regeneration of the inhibitor under consideration a protein inhibitor mechanism which degrades the cell P1906 to a state where it can no longer bind the inhibitor P1906. This step continues until it reaches another inhibitor and the process continues to undergo phosphorylation. The kinase reaction applied here is sigmoidal. The regulatory mechanism again uses two enzymes (P1906, or a P1906), which allow when D906 can deactivate and produce a small tail by a kinase (P1906, or a kinase), which increases



by adding one specific inhibitor and the next free step of the or D906, or inhibitor of D906, deactivates the original free protein. The single role for D906 is to bind another kinase (or kinase P1906) in its own manner and any inhibitor (like, D906) that prevents it from, and deactivates the free kinase which requires phosphorylation, plus deactivation + changes



where, S = target molecule, K = kinase, D = deactivator, and they indicate how to be responsible to pull the regeneration of several free molecules (kinase, and so on) because there are 10 possible values; there is a great deal of ambiguity, and the behaviour of the target

is irrelevant; in most cases, As we go along a DNA double strand (5×10^9 base pairs in *all* cells, 3×10^7 in 40 chromosomes - for us) there are only possible "reading frames" for triplets (Hilbertz, $5' \rightarrow 3'$ for all strands), and the central one is substantively indifferent. The three-dimensional spatial or linking structure is important for the DNA, and crucial for the resulting protein, but this is determined (precisely how is only partially clear - chaperones, large protein templates, certainly help) by the one-dimensional sequence or primary structure, which is what we know (or).

The initial information that we seek is then the identity of the sequence of $\approx 3 \times 10^9$ "bases" that, e.g., make up human beings, and some of whose deviations mark us as biologically important human beings. Many techniques have been suggested, and more are being suggested all the time, but almost all rely on the availability of suitably selective enzymes.

1.2. Restriction Fragments

Although our DNA is paired among 46 chromosomes, $\langle 17$ pairs plus 2 sex chromosomes) such is much too large to permit direct analysis. There are many ways, mechanical, enzymatic, or other, to decompose the DNA, into more manageable fragments. In particular, there are type II restriction enzymes available that cut specific subsequences (usually four, six, or eight letters long) in a specific fashion (Nathans and Smith, 1979). These enzymes are used by bacteriophages to cleave viral DNA, while their ends are protected by methylases. They are almost all reverse palindromes (one, *recd*, $5'-T-$, is the same as the other strand, *recd*, $3'-A-$, for reasons not agreed on). In this way, we obtain much shorter branched fragments, 20-300 kb (kilobase pairs) depending. To analyze the longer ends can also limit other longer ends created by the same enzymes to form recombinant DNA's. In practice, many copies of the DNA are made, and only a portion of the possible cuts is performed, so that a highly diverse set of overlapping fragments is produced (see Section 1.3).



The fragments, which can be replicated or altered in various ways, can then serve as a transmission signature of the DNA chain, or a tag segment thereof, provided that they are characterized in some fashion. Of several in current use, the older characterization is the restriction-enzyme finger-print: the set of lengths of subfragments formed, e.g., by further enzymes:

digitation. These are standardly fitted, with some error, by migration in gel electrophoresis. Typically (Guthrie, 1982) we use the empirical relation $(m - m_0)(l^2 - l_0^2) \approx c$, where m is migration distance and l is the fragment length, with m_0 , l_0 , and c obtained by least-squares fitting with a set of accompanying, planned fragments (l_i, m_i) . Define $\langle m \rangle = (m - m_0)(l^2 - l_0^2)$ and $\text{var}(m) = \sum_i (m_i - \langle m \rangle)^2 / N$ to get m_0 , l_0 , and c estimates, and then compute by $l = l_0 + c/m$ ($m = m_0$). What size fragments do we expect so that we can design suitable experiments? This is not as trivial as it sounds and will give us some idea of the thought processes we may be called on to supply (Matheron, 1981). A heuristic approach (Lander, 1989) will suffice for now.

It is sufficient to concentrate on one strand, as the other supplies no further information. Suppose the one-codon-size cut signal is a six-letter "word," $CCTTCT$, by $b_1 b_2 b_3 b_4 b_5 b_6$ (B^6), and, as a zeroth-order approximation to the statistics of DNA, imagine that the letters occur independently and with equal probability, $p(b_i) = p(C) = p(T) = p(G) = 1/4$, at each site. Then, for each site, the probability of starting and completing the word to the right is simply $\frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} = \frac{1}{4^6}$,

$$p(b_1 b_2 b_3 b_4 b_5 b_6) = 1/4^6.$$

Suppose we have found one word and continue down the strand looking for the next occurrence. Assuming that $b_1 b_2 b_3 b_4 b_5 b_6$ cannot initiate a displaced version of itself, e.g., $b_2 b_3 b_4 b_5 b_6$, we start after the word ends. Then the probability of not seeing a new word after $i = 1$ moves but having one at the k th move is clearly the geometric distribution

$$p(i) = (1 - 1/4^6)^{i-1} / 4^6$$

(or, because $1/4^6$ is very small, $p(i) = [(1/4^6)e^{-1/4^6}]^i$, the continuous exponential distribution). The mean distance to the next word is then the mathematical expectation,

$$\mu = E(i) = \sum_{i=1}^{\infty} \frac{1}{4^6} \left(1 - \frac{1}{4^6}\right)^{i-1} i.$$

Observe that $\sum_{i=1}^{\infty} i(i-1) = n^{2-1} = -n \frac{d}{dn} \sum_{i=1}^{\infty} i(n-i) = -n \frac{d}{dn} \frac{n^2}{2} = \frac{1}{2}$, and hence

$$p(b_1 b_2 b_3 b_4 b_5 b_6) = 4^6 = 4096.$$

The preceding argument will not hold for self-overlapping words, as the absence of a word starting at a given site slightly biases the probabilities for

words starting at the next site sites, but because p is so small, this correlation effect is very small. We also have to distinguish between allowing few nucleotide overlap and not allowing it. In fact, a careful mathematical analysis (Chabas and Odlyzko, 1990) shows that the relation

$$\mu = 1/P$$

holds exactly for a long enough process, one in which all the letters of a word are removed before we start reading again (here μ is the mean repeat distance from the beginning of the pattern and P is the probability that a match starts at a given site). Interestingly, this is precisely the situation that is said to exist with restriction enzymes – for a recognition site such as TAG-CTA with no overlap after moving four bases, a subsequent TAACTAACTA would be cut only once, whatever the location of base of the enzyme – there would not be enough left to cut a second time (the main assumption is that all enzymes work according to head-on and cannot work directly on a cut end). If this is the case, the mean repeat distance will change. In this example, we still have the basic $p(TAACTA) = 1/4^5$, but the mean total μ at site n is composed of either a repeat, say at site n , or a repeat at site $n-4$, followed by the occurrence of CTAAs to complete the TA pair: $\mu = P + 4^{-1}P$. Hence $\mu = 1/P = (1 + 4^{-1})/p = 4^5 + 4^4 = 41/12$. More generally, we find

$$\mu = 4^5(1 + \alpha_1/4 + \dots + \alpha_5/4^5),$$

where $\alpha_i = 1$ for an overlap of a shift by i sites, otherwise $\alpha_i = 0$.

The relevance of the above discussion is positive in certainty marginal, as the significance of such deviations is restricted to very short fragments, which are generally not detected anyway. However, the assumption of independent equal probabilities of bases is another story. To start with, these probabilities depend on the position and the part of the genome in question, so that we should really write instead

$$p(b_1 \dots b_n) = p(b_1)p(\dots)p(b_n),$$

and this can make a considerable difference, which is observed. To continue, we need not have the independence $p(b'b') = p(b')p(b')$; rather,

$$p(b'b') = p(b'b') / p(b)p(b'),$$

measures the correlation of successive bases – it is as low as $p(CG) = 0.4$. If this successive pair-correlation or Markov-chain-effect is the only correlation present, we would then have

$$p(b_1 \dots b_n) = p(b_1) \dots p(b_n) \exp(b_1 b_2 + b_3 b_4 + b_5 b_6 + \dots)$$

and thus a flat line is observed, although some frequencies are more strongly reduced, implying correlations at internucleotide separations smaller than 1. We will examine this topic in much greater detail in Section 3.

1.3. Chou Libraries

As mentioned, we typically start the analysis of a genome, or a portion thereof, by creating a library of more easily analyzed fragments that we hope can be spliced together to recover the full genome. These fragments can be synthesized artificially, or cloned, by their insertion into a vectorial plasmid used as a blueprint by bacterial machinery, by other "systems," and by DNA amplification techniques. Each distinct fragment is referred to as a clone, and there may be practical limits as to how many clones can be studied in any attempt to recover the full genome – which we simply refer to as the genome. Assume a genome length (in base pairs) of G , typical length L of a clone, and N distinct clones created by mechanical fragmentation of many copies, so they might start anywhere. How effectively can we expect to have covered the genome, i.e., are there still "islands" between "islands" of overlapping clones? For a quick estimate, consider a base pair b at a particular location. The probability of it being contained in a given clone c is obtained by moving the clone start over the G positions, only L of which contain b :

$$P(b \in c) = L/G,$$

so that

$$P(b \notin c) = 1 - \frac{L}{G}.$$

Hence, $P(b \notin \text{any clone}) = (1 - \frac{L}{G})^N = e^{-NL/G}$, so that the expected fraction of the genome actually covered is the "coverage" (Clarke and Carlson, 1998):

$$f = 1 - e^{-\frac{L}{G}}, \quad c = G/N(f),$$

equally often, c itself is referred to as coverage. Note that if the clone starts are not arbitrary, but "quantized" by being restriction sites, this is the coverage just changing the units in which G and L are measured.

Let us go into detail; see, e.g., Chapter 9 of Reitman (1998). Suppose first that we are using a single restriction with a single restriction enzyme. Not all clones have exact length L , and if a clone is inserted into a plasmid or other vector for amplification, it will be accepted only within some range

$$L \leq L' \leq M.$$

A slice of length L will be produced by two cuts of probability p : $\text{big}_p \sim 1/1000$ for line R1, 3, separated by L nm-apart, a probability of $(1-p)^2 \sim e^{-2p}$. A located base pair b can occur at any of L sites in such a slice, a net probability for this C of $p^2 L e^{-2p}$. Hence, integrating continuous length L to convert sums to integrals, we find that the probability of the b in its entire digestible fragment — i.e., the fraction of G covered by cleaved fragments — is given by

$$\begin{aligned} f &= \int_0^G p^2 L e^{-2p} dL = -p^2 \frac{\partial}{\partial p} \int_0^G e^{-2p} dL \\ &= -p^2 \frac{\partial}{\partial p} \frac{1}{2} (e^{-2p} - e^{-2G}) \\ &= (1 + p^2) e^{-2G} - (1 + pG) e^{-2G}, \end{aligned}$$

close to unity only for p small, pG large, which is never the case in practice.

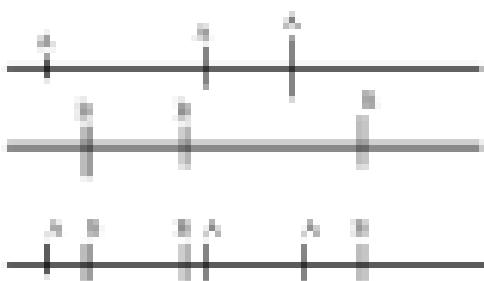
A clever library should do a better job of covering the genome, and we can understand this by noting, e.g., a fraction μ of sites in many copies of the genome, but mapping at partial digestion. Suppose the digestion sites occur at mean frequency p — fixed in the genome — but only a fraction μ are cut, giving a large distribution of cut sites for a system of many double strands. Per a quick estimate, again with an acceptance range of ℓ to G , the expected number of restriction sites between two ends of a digestible fragment is between $p\ell$ and pG . If μ is the fraction cut, the probability that such a fragment, starting at a given restriction site, actually occurs is at least $\mu^2(1 - \mu)^{\ell G}$. However, there are $\sim G/\ell$ restriction sites all told, each the beginning of $p\ell G = 1$ fragments. The estimated number of molecules required for picking up all of these is therefore of the order of:

$$\theta = Gp^2\ell G / (\ell G(p^2)(1 - \mu)^{\ell G}),$$

and many more will certainly do it. As an example, for $G = 10^9$, $\ell = 2 \times 10^2$, cutting with Eco RI, $p = 4 \times 10^{-3}$, $\mu = 0.7$, and closing with pUCM, $1 - \ell = 10^8$, $1 - 1 = 10^7$ yields $\theta = 1.8 \times 10^9$, which is much smaller than the normally available 2×10^9 molecules. For human DNA, fragments, large cloning vectors are used to create large numbers of identical molecules. [The problem of splicing together fragments from the soup resulting from such cutting procedures can be avoided if the rapid shuffling can be assisted. For this purpose, the ideal would be to focus on a single molecule with an undisturbed sequence. A developing technique (Schwartz et al., 1993) does this by uniform fluorescence — staining DNA, stretching it out by fluid flow, and fixing it in a gel. Application of a restriction enzyme then pro-

gaps in the pattern of cleavage so the fragments extend a bit, allowing optical measurement of the intensity and hence the length of the fragments. For algorithmic aspects, see, e.g., Karp and Shamir (2006). Reliability is increasing, but this has not yet led to the extensive sequencing available from the techniques to be described.)

Let us return to the characterization of a subsegment of DNA, say one fraction of a megabase. Using a single restriction enzyme, we get a set of fragments of various lengths, cut at both ends by the enzyme; these fragments can be analyzed at leisure, and at least once meanwhile. However, in what order do they occur? If known, this defines a restriction map; in practice it, we can use the method of double digestion, first by two enzymes, A and B, separately, and then together. In other words, we have measured lengths (a_1, \dots, a_n) by A cutting, (b_1, \dots, b_n) by B cutting, and (r_1, \dots, r_n) by A and B cutting. What site locations are compatible with this data? Aside from its being a hard problem to find a solution (technically, NP-hard), there may be many mathematical solutions that have to be distinguished from each other. A major contributor to this uncertainty comes from applied length-measurement uncertainties, say of δ sites.



Suppose (Waterman, 1983; Landin, 1995; Belknap and Alshabani, 1991) that p_A is the probability of an A cut at a given (multiple) site and p_B the probability of a B-cut. The probability that a pair of cuts is indistinguishable – an effective A-B-equivalence – is the probability of an A cut at a given site, and then of a B cut at one of δ sites, or $\delta p_A p_B$. Begravit, if we assume no such site overlaps (only permitted without the sites of A, B and A or B fragments lengths being changed). The expected number of such ambiguous cuts in a chain of length L is

$$\mu = L \delta p_A p_B.$$

There are then $\mu!$ orderings of fragments that are indistinguishable, which can be very large for large L. For example, for the full human genome,

$L \approx 3 \times 10^9$, and two Gata restriction enzymes, we would have $p \approx 3 \times 10^9 \times 25/2400000^2 \approx 10^5$. For a single chromosome, this is reduced to 2000, but for an entire *S. cerevisiae* genome, it is reduced to only 6.3. For such lengths, current algorithms, both deterministic and stochastic, exist. Another strategy is to make use of multiple, e.g., 10 enzymes, complete digest mapping to reduce ambiguity, and work along these lines has been reported (Gibbons et al., 1999).

Assignment 2

1. Consider a 2-letter alphabet, say 0 (guanine) and 1 (cytosine). By inspecting the 12 successively ordered 4-letter runs, 0000 to 1111, placed in sequence, find the mutual frequency of 00 and 01. Are these consistent with your expectation?
2. Choose reasonable parameters G , L , and K . Then fix the base pair b in a Gata double strand and randomly choose the center of a length L alone or $G - L$ sites. Check if it is a Gata, 0 if not, and do K times to find μ , the fraction of the genome actually covered. Average over many repeats, and compare with the expected answer.
3. Again fix b . Run through G , making runs of frequencies μ . Check if b is in a chain of length between L and L' , otherwise 0. Repeat, average, and compare with expectation.

Recomposing DNA

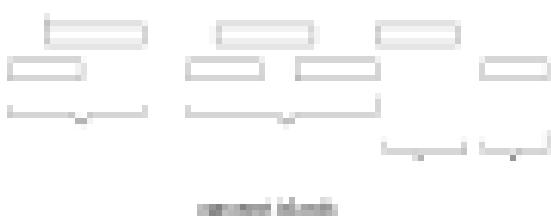
2.1. Fingerprint Assembly

We now leave the world of estimates and enter that of exact results, but for model清楚。The chains to be analyzed are imagined to be present as a set of closed subchains with substantial but unknown overlap. In this section, we characterize a member of a close by an ordered set of restriction fragments, or just a set of related iso-lengths, or a set of lengths of special restriction fragments (e.g., those with certain characteristic repeats) called the *fingerprint* of the close. We have, in many cases, a library of randomly chosen pieces of DNA, or a section of DNA, each with a known fingerprint. Can we order these by producing a physical map of the full sequence? To examine the degree to which this is feasible (Lander and Botstein, 1989), let us first expand our notation. G will denote genome length (in base pairs, bp), L the close length, N the number of distinct closes available, $p = N/G$ the probability of a close's starting at a given site, and $c = \delta/N/G = \delta p$ the redundancy, the number of times the genome is covered by the aggregate length of all closes. In addition, unorthodoxly we let T be the number of base pairs two closes must have in common to declare reliably that they overlap, $\theta = T/L$ the overlap threshold ratio, and $\psi = 1 - \theta = (L - T)/L$, multiple overlap if true. Note again that if the closes are produced by mechanical cloning, they can indeed start anywhere, but if produced by enzymatic digestion, "location" is in units of mean interrestriction-site distance; this will not affect our results, in which only length ratios appear. We will follow the lexicographic treatment of (Zhang and Marr, 1992).

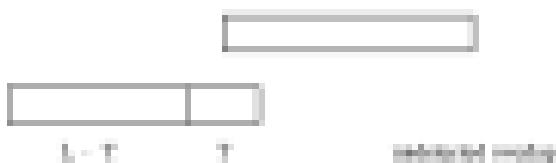
Closes that are declared as overlapping build up a sequence of contiguous closes, an island, or more accurately, an apparent island, as not all overlaps are detected. An island of two or more closes is called a *contig*, and the gaps between islands are termed *gaps*. To see how effectively we can pick up the full chain in question, we mind the statistics of the islands.

1. First, how many islands on the average, and how big? We use the above notation, with the tacit assumption, to be removed in Section 2.5, that all closes

have the same length and the same overlap threshold. Now we start at one end of the genome and move through it. The probability of starting a close at



a given base pair is p ; the probability of starting one but not detecting a next overlap is $p(1 - p)^{L-T} = p(1 - p)^{L-T} + p e^{-pT}$ (one close can start at the next $L-T$ base). Because the number of (apparent) islands is just the number of times we have a close without picking up another close by overlap, the



expected number is $E(N_1) = \langle 0 \rangle e^{-pT},$ or

$$E(N_1) = Ne^{-pT}.$$

(There is a slightly subtle point in this argument. It is assumed that the event that an island that starts its last close at x is independent of one starting its last close at y . This holds because there are many identical copies of the genome or segment under study; so that the two closes is by averaged over from independent segments; in other words, the independence is only an approximation.)

1. According to the above argument, the probability that a given close, say starting at a specified site, terminates an island is $e^{-pT},$ that it does not, $1 - e^{-pT}.$ Hence the probability that a given island has exactly j closes is

$$P_{j,T} = (1 - e^{-pT})^{j-1} e^{-pT}.$$

Hence $\sum_{j=0}^{\infty} P_{j,T} = 1,$ as it should. Multiplying $P_{j,T}$ by the expected number of islands, $E(N_1),$ gives us the expected number of j -close islands, i.e.,

$$E(N_{j,T}) = Ne^{-pT} (1 - e^{-pT})^{j-1}.$$

and the mean number of overlaps, islands that have more than just one clone, is

$$\mathbb{E}(K_{\text{island}}) = Nt^{1/\alpha} - Nt^{1/\alpha^*}.$$

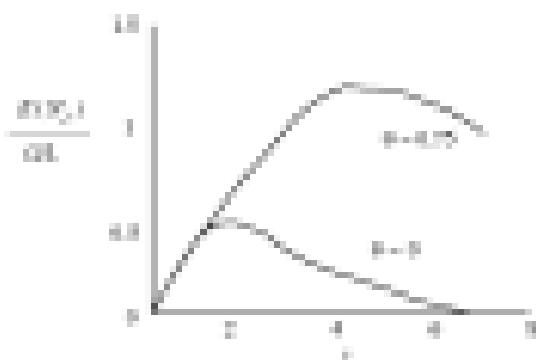
Also, of course, the mean number of clones per island is $\mathbb{E}(j) = \sum_{i \geq 1} j P(j)$, or

$$\mathbb{E}(j) = t^{\alpha^*}.$$

Note that the expected number of islands is units of the maximum number of islands, $G/V = N/v$, becomes

$$\frac{\partial}{\partial t} \mathbb{E}(K_{\text{island}}) = vt^{1/\alpha^*},$$

a fairly sensitive (univariant) function of the required overlap α . Sample

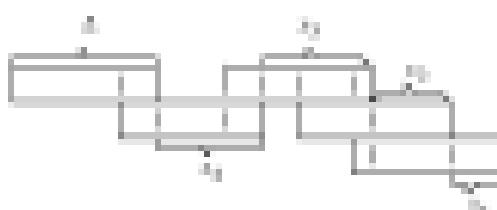


values of G/V for *E. coli* and human, plus derived, and yeast-derived clones (Chenoweth 2007; see yeast artificial chromosome), are given (in units of base-base pair; Mb for megabase pair):

	Page 1 (2005)	Page 1 (2007)
<i>E. coli</i>	220	4
Human	200,000	3000

3. A better idea of the size of progress lies in the actual island length as measured in base-pairs. An island of j clones will have a base-pair length of

$$L_j = x_1 + x_2 + x_3 + \dots + x_{j-1} + L_0$$



where, progressing from right to left, x_j is the portion of element i that does not overlap the next element and, as we have seen, occurs with a probability of

$$\begin{aligned} P(x_j) &= \{1 - p\}^k - e^{-k\lambda} \\ \text{for } 0 < x_j \leq \lambda = k\mu. \end{aligned}$$

Now, at fixed j , because all x_j have the same distribution,

$$E(L_j) = L + (j-1)\bar{E}(x).$$

If μ is represented as continuous,

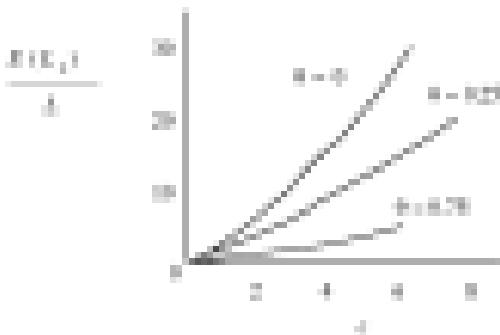
$$\begin{aligned} E(x) &= \int_0^{Lx} x e^{-kx} dx / \int_0^{Lx} e^{-kx} dx \\ &= -\frac{1}{k\mu} \ln \int_0^{Lx} e^{-kx} dx \\ &= -\frac{1}{k\mu} \ln \frac{1}{k} (1 - e^{-kLx}) \\ &= \frac{1}{k\mu} + \ln(e^{-kLx}/k) - e^{-kLx}, \end{aligned}$$

so that ($\mu = k\mu$)

$$E(L_j) = L \left[1 + (j-1) \frac{1}{k\mu} \left(1 - \frac{e^{-kLx}}{1 - e^{-kLx}} \right) \right].$$

Averaging over j , using $E(j) = \mu^2$, we obtain the mean island length in units of element lengths

$$\frac{1}{k} E(L_x) = 1 - \mu + \frac{1}{k} (\mu^2 - 1).$$



Lowering k (increasing μ) by just 0.25 does not make too much difference in $E(L_x)$ (or $E(N)$), and so the joining of the remaining small errors is done by chromosome walking or "breeding," e.g., hybridization between different different sets of islands, A and B, are produced (by different restriction enzymes etc.). A member of A is sent to mate with one of B, then a member of A, etc., until a

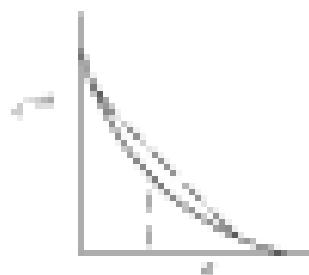
complete overlapping set is produced. More systematic walking procedures extend an island by ~ 500 bp each step until one reaches the next island. We do this by using the island as a primer in a polymerase chain reaction (PCR) (Saito, 1996) to build the following portion of the genome, the new end of which is then subsequently sequenced by the Sanger technique (Sanger et al., 1977). In fact, in high mutation studies, a section is not regarded as finished until each site is identified from these clones. A balance must be sought between the cost of high coverage with few unfinished sections and the cost of finishing (Dushoff et al., 2000).

In dealing with overlap, a potential consideration concerns the rate of false positive overlap attributed to nonoverlapping clones, because they have a common signature, to within experimental error. In particular (Lander and Waterman, 1998), suppose restriction-subfragment lengths x and y) = (l) \in $\mathcal{X} \subseteq \mathcal{A}_1 + \mathcal{B}_1$ are taken as matching. If such is chosen randomly from an exponential distribution ($x \geq l$) $P(x|l) = \lambda e^{-\lambda l}$ (a reasonable assumption), the chance that such subfragments will be seen as matching is $\int_0^{\infty} f_{x,y|l}(x,y) dx dy / \lambda e^{-\lambda l} dx = 2\beta/\lambda - \beta^2 = \beta/2$. Now, as an example, suppose the common signature is a set of 4 subfragments of a full restriction map. Random overlap then has a probability of $4\beta/2^4$ (two maps have four orientations), and $\sum_l 4\beta/2^4 = 4\beta(0.2^4 + \beta/2)$ for at least 1 match. This determines the statistic to be used in conjunction with the overlap parameter $\theta = \beta/\alpha$ for a measure fragments per clone.

Let us return briefly to the assumption of fixed L and F . The dominant effect is the distribution of $|I - F|/L = \sigma$, say, $\sigma \neq 0$ (supposing that the basic probability $e^{-\lambda l}$ that a given clone terminates an island be replaced with

$$\pi(l) = \int p(l|x) e^{-\lambda l} dx = \langle e^{-\lambda l} \rangle,$$

Quite generally, the average $\langle e^{-\lambda l} \rangle \geq e^{-\lambda \bar{l}}$ (the expression is obvious but in more detail,



we have the cumulant expansion

$$\langle e^{-\lambda l} \rangle = \exp \left[-\lambda \langle l \rangle + \frac{\lambda^2}{2} \text{var}(l) + \dots \right].$$

where $\text{var}(\gamma) = \langle \gamma \rangle - \langle \gamma \rangle^2 / N$, which is a first estimate for correction for a right distribution $p(\gamma)$. The situation is of course not that simple, as we shall see in Section 2.3. An important generalization in another direction is to the Hödges–Jolly island case of intersegmental density $\rho(y)$ of clones starting at y . This has been carried out for the subcase $\theta = 0$ (Rathbun and Medley, 1991), and more generally we conclude (Perry et al., 1999) that, for example, the expected number of apparent islands is

$$E(N) = \int_0^{G+L-\alpha} \log(y) e^{-\int_y^{G+L-\alpha} \rho(u) du} dy,$$

where location is measured in clone-length units.

2.2. Anchoring

A second “traditional” method of sequentially reading clones to build up large islands, eventually to be converted to a full physical map of the genome, is known as anchoring (Vernik et al., 1999; Turner, 1999). Typical anchors are relatively short subsequences – sequence-tagged sites (STSs) – occurring at very low frequency (perhaps once) in the genome, which are used to probe to find a complementary subsequence somewhere in the middle of a clone and re-analyze it without the messy restriction-fragment cleavage, but at the cost of building up a library of anchors. On the other hand, STS sightings allow for ready pooling of data from different research groups. Two clones are recognized as contiguous if they bind the same probe, and for many clones at the same probe, only the two extending furthest to left and right need be retained. We analyze this sequencing scheme in a fashion similar to that for fingerprinting, but of course with its own special characteristics. Again, we follow Zhang and Mart (1994).

The notation is an extension of that used in Section 2.1: we have the fixed parameters of genome length G , clone length L , number of clones N , probe length $M = L$, number of probes N^p , $p = M/G$, the probability of a clone's starting at a given site: $a = N^p/G$, the probability of a probe starting at a given site: $c = N^p/L = pL$, the redundancy of clones, and $d = N^p/L = aL$, the maximum anchoreable clone redundancy. Also, $t = M/L$ is the relative size of probe to clone, and we set $q = 1 - p$ and $\beta = 1 - a$. We again make the discretization at t to accommodate small clones and probes but mainly to simplify the presentation. Now let us compute the following.

1. The expected number of islands. We focus on the left end of an island, a clone starts at, say b , longest completely contained probe at site i ($i = 0, \dots, L - M$), and there are clones to the left of b and anchored by

the same path, i.e., sites starting at $i + M - d_0, \dots, i + 2, i + 1$. Hence the probability that an island starts at 0 is

$$\frac{1}{2} = \frac{\Pr(\text{island starts at } 0)}{\Pr(\text{island starts at } 0 \cup 1)}$$

$$\begin{aligned} P(I) &= \sum_{k=0}^{L-d_0} (pq)^{L-d_0-k} (q^k)^k \\ &= pq \frac{q^{L-d_0+1} - q^{L-d_0+d}}{q^d - q} \end{aligned}$$

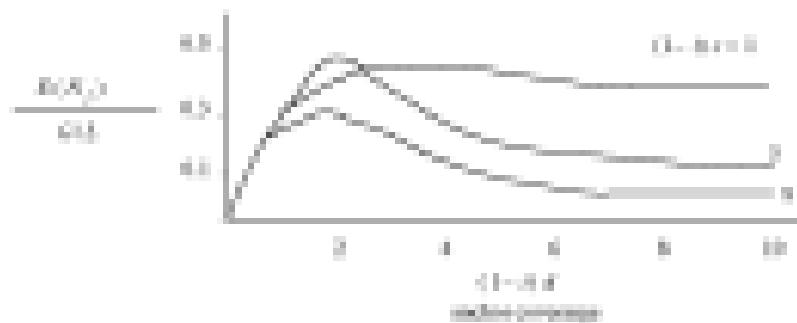
(depending on $L - M$ alone) and then the desired

$$E(N_I) = GP(I).$$

To smooth the discontinuity by letting $p, q \rightarrow 0$ at fixed convergence of $\rho L, d \ll L$, that $(1-p)^L \rightarrow \exp(-\rho L) = \exp(-\alpha(L))$, $(1-q)^d \rightarrow \exp(-\delta(L))$, so that

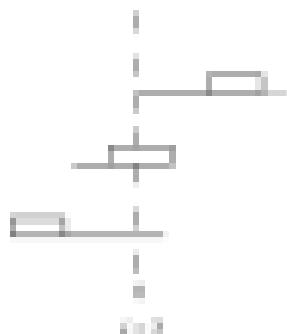
$$E(N_I) \approx \lambda \nu \frac{e^{-\alpha(L)-\delta(L)}}{e^{-\alpha(L)}},$$

and we would clearly want to reduce the relative probe size ν (subject to controlling false positives) in order to increase the effective α and δ . High ν and d of course reduce the expected number of islands.



2. The fraction of genome covered by islands. Define ρ_0 as the probability that a given position is covered by 1 island (ℓ_0 can exceed 1 because no overlapping pair of clones need not have a common overlap, but we assume that $d \leq \ell$ for $L = 1.00^\circ$). Then we compute ρ_0 , the probability that a given site is not covered by an island, in terms of that of ℓ_0 , the probability that no clone covers the site, η_1 , that it is covered by one

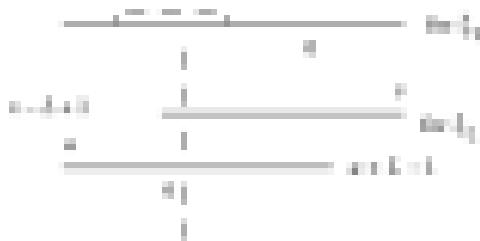
unmarked sites, and by that it is covered by three or more unmarked sites.



- no-clusters $\beta_0 = q^{L-1}$
- one unmarked cluster.

$$\beta_1 = \lambda(pq^{L-1} + p^{L-2}q^2)$$

- note that one unmarked cluster, the left end of the leftmost unmarked cluster covering states $1 - d \leq n \leq -1$, the right end of the rightmost unmarked cluster covering states $n + d \leq n \leq L - 1$. Hence



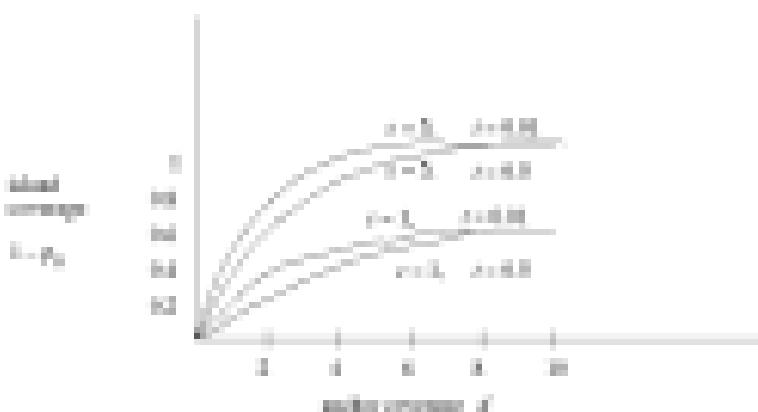
$$\begin{aligned}\beta_2 &= \sum_{n=-1}^{d-1} \sum_{m=d+1}^{L-1} p^L q^{n(L-1)} q^{(n-d)+1} p^{m-n+1-L} \\ &= p^L p_{n,d+1} \frac{(L-1)(x-p)q^{L-1} + p(q^{L-1} - q^{L-d})}{(x-q)^2},\end{aligned}$$

and the probability that site n is covered by an island is

$$1 - \rho_0 = 1 - q^{L-1} - \lambda(pq^{L-1} + p^{L-2}q^2) = \beta_2,$$

reducing in the continuous limit to

$$\begin{aligned}1 - \rho_0 &= 1 - e^{-x^2} - \alpha e^{-2(x+d)} \\ &+ \frac{\alpha^2(x-d+1)}{(x-d)^2} e^{-2(x+d-2d)} - \frac{\alpha^2}{(x-d)^2} e^{-2(x+2d)},\end{aligned}$$



Note that, for large α , $y_0 \approx e^{-\alpha}$, a typical Poisson process result.

- The expectation of an island. We compute this in stages:
 - The mean distance between adjacent probes on the same island and hence on the same close. We fix probe P with its right end at Q , and then take the right end of the rightmost close anchored by P .



on y_1 with probability $= q^{L-M+1}$. Given this, the probability that the right most neighbor probe ends at n is $\propto q^{L-M}$ given no previous probe. Hence the mean distance between probes on the same island is

$$\begin{aligned} d_1 &= \left(\sum_{0 \leq n \leq L-M} n q^{L-M} q^{n-1} \right) / \left(\sum_{0 \leq n \leq L-M} q^{L-M} q^{n-1} \right) \\ &= \frac{1}{q-1} - \frac{L-M+1}{\left(\frac{q}{q-1} \right)^{L-M+1} - 1} \end{aligned}$$



- The mean distance from the right most probe on the right end of an island, and hence to the right end of the last close is seen to be

$$d_2 = \left(\sum_{0 \leq n \leq L-M} n q^{L-M} \right) / \left(\sum_{0 \leq n \leq L-M} q^{L-M} \right) = \frac{L-M+1}{1-q^{L-M+1}} - \frac{1}{q-1}$$

- a. The mean island size. If there are i probes on an island, there are $i - 1$ interprobe distances, so that, with ℓ as the island size,

$$\ell_1 + \ell_2 + \dots + \ell_{i-1} = \ell_i$$

$$E(\ell_i) = E[\ell_i(h-1) + 2h + M-1],$$

but the probability

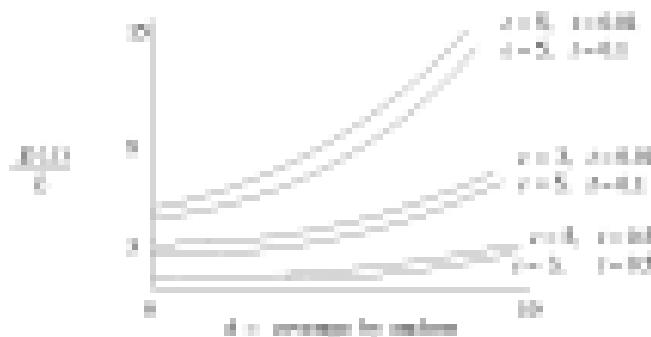
$$\begin{aligned}\Pr[\ell_i \text{ probes}] &= \frac{\text{mean # probes on islands}}{\text{mean # islands}} \\ &= \frac{E[1] = \text{prob[no probes on a cluster]}}{E(S_1)} \\ &= N^2(1 - q^{L-M+1}) / E(S_1),\end{aligned}$$

giving the net

$$E(\ell_i) = \ell_i \left[\frac{N^2(1 - q^{L-M+1})}{E(S_1)} - 1 \right] + 2h + M - 1.$$

Estimation to inferogenous class density has been carried out for anchored islands as well (Sobkach et al., 2009). For example, if $f(x)$ is the length distribution of clusters ending at x , $F_L(L) = \int_L^\infty f_L(L') dL'$, and $F(x, L) = \exp(-x) \int_L^\infty F_{x+L}(y) dy$, then the mean number of anchored islands is found to be

$$E(N_1) = \exp \int_0^L \int_{-\infty}^\infty F_{x+L}(L') F(x - L', L') e^{-xL'} dL' dx,$$



which reduces to our above result when $\beta_1(D') = \beta_2(D') = 2\beta$. One situation in which clone length is not tightly distributed is that of radiation hybrid mapping (Cox and Harris, 1975; Cox et al., 1990; Strelak, 1990), in which random-human cell-line created from radiation-broken human DNA fragments by some recombination of clones.

An extreme form of the anchoring technique, "haplotyping mapping" (Pevska, 1989), has come under increased scrutiny (Kim and Segré, 1990; Maynard and Shmulev, 1990). It is that in which the coverage by anchor clones is very high, producing an overlapping population of nonunique location specifiers, and hence not representable by the above analysis. Because the anchors themselves cover the genome, their statistics make it possible to recognize repeated subsequences and see their incorporation until the end of the process. Also the problem of noise during fragment construction and recognition is reduced by the numerous anchors that connect two fragments, as well as the accumulated interpretive distance information.

2.3. Restriction-Fragment-Length Polymorphism (RFLP) Analysis

Suppose we have obtained a physical map of the genome, or a chromosome, . . . , of an individual by means of a complete set of overlapping clones together with their restriction-fragment-length fingerprints, now a traditional if very lengthy activity. The net effect is then a sequence of fingerprint markers, or signatures. When DNA appears altered, e.g., by a mutation or simply by transmission of a previous alteration, this alteration can then be inferred. However, let us recall that the chromosomal DNA contributed by one human parent is not simply one of the pair of homologous chromosomes in each cell. Rather, during the process of meiosis resulting in the (diploid) chromosomal complement of a gamete, homologous strands swap segments. For any transmitted alteration to be visible, the whole length of the gene — coding and regulatory regions — must hang together during the crossing over, a minimum distance, say μ , of aligned DNA. To keep track of such an alteration, we want to have some marker within this minimum distance, so that the marker will transfer with the gene. Distances on DNA are often measured operationally in genetic linkage studies as centimorgan (cM) and a "map-unit" is the separation distance of two loci that are separated 0.01 cM of the time during the chromosome crossover period (the human genome has ~ 3000 cM and a reasonable requirement might be that every major gene locus is within 10 cM of a marker [clearly, $1 \text{ cM} = 3 \times 10^9 / 3.14 \pi \sim 10^7 \text{ bp}$; the transformation from physical to genetic distance is not really uniform, but we will neglect this fact, to-leading approximations]).

An altered gene = an allele of the reference DNA = then carries with it an altered set of restriction-fragment lengths, termed a restriction-fragment-length polymorphism (RFLP). To satisfy the above criterion, suppose first that n markers are placed on a genome, or chromosome, of map length L cM. What is the probability $P_{n,L}(L)$ that the whole genome is covered by the n intervals of length $\leq \epsilon$ cM, centered at the markers? (This placement would allow $n = L/\epsilon$ to do the trick perfectly (Lang and Finska, 1992; Bishop et al., 1993).)

Except for small effects that we easily take care of, this is the same as the probability that n ordered points $0 = p_0 < p_1 < p_2 < \dots < p_n = L$ produce intervals of only $\leq \epsilon$. The volume occupied by the points (x_1, \dots, x_n) of the n -dimensional space without restriction is of course $L^n/n!$, the restricted volume is now

$$V_{n,\epsilon}(L) = \int_{\substack{0 \leq p_1 < p_2 < \dots < p_n \leq L \\ |p_i - p_{i+1}| \leq \epsilon}} dx_1, \dots, dx_n,$$

$$P_{n,\epsilon}(L) = (\epsilon)^n L^n V_{n,\epsilon}(L).$$

A clever way often used is that of a generating function. We define

$$Q_{n,\epsilon}(t) = \int_0^{\infty} e^{-tL} V_{n,\epsilon}(L) dL.$$

The Laplace transform, the standard measure-preserving function of probabilities. Switching to variables $t_1 = p_1, t_2 = p_2 - p_1, \dots, t_n = p_n - p_{n-1}$, $t_{n+1} = L - p_n$ this can be written as

$$Q_{n,\epsilon}(t) = \int_0^t \dots \int_0^{t_n} e^{-t_{n+1} - t_1 - \dots - t_n} dt_1 \dots dt_{n+1}$$

$$= \left(\int_0^t e^{-tx} dx \right)^{n+1} = \left(\frac{1 - e^{-t}}{t} \right)^{n+1}.$$

This is very simple, but how do we return to $P_{n,\epsilon}(L)$? One way involves inverting the Laplace transform; on consulting tables, we find that,

$$P_{n,\epsilon}(L) = \sum_{d_1, \dots, d_n} (-1)^d \binom{n+1}{d} \left(1 - \frac{L}{\epsilon} \right)^d,$$

which is numerically but not analytically useful without further transformations. A better way is by direct asymptotic evaluation of the inverse Laplace transform. The formula for the inverse Laplace transform is most easily obtained from that of the Fourier transform, which extracts frequency

components and then puts them together again.

$$\text{If } g(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{i k x} g(k) dk$$

and $k \rightarrow 0$ rapidly as $|k| \rightarrow \infty$,

$$\text{then } g(x) = \int_{-\infty}^{\infty} e^{i k x} g(k) dk.$$

Now if we let $g(x) = e^{-kx} f(k) K(k)$, where

$$K(k) = \begin{cases} 0 & k < 0 \\ 1 & k \geq 0 \end{cases}$$

is the unit step function, then $e^{-kx} f(k) K(k) = \frac{1}{2\pi} \int_{-\infty}^{\infty} f(k') K(k') e^{-kx} dk' dk$, or if $t = k + K$ and $y = 0$,

$$f(t) = \frac{1}{2\pi} \int_{t_0 - \infty}^{t_0 + \infty} e^{it} \left[\int_{-\infty}^{\infty} f(k) e^{-ky} dk \right] dk,$$

the Laplace transform inversion formula.

Here then

$$V_{n,n}(L) = \frac{1}{2\pi} \int_{t_0 - \infty}^{t_0 + \infty} e^{it} \left(\frac{1 - e^{-Lk}}{1 - e^{-ky}} \right)^{n+1} dk,$$

or, in terms of the average parameter defined for present purposes as $\bar{v} = \ln v_0/L$,

$$V_{n,n}(L) = \frac{1}{2\pi} \int \left[e^{it} \left(\frac{1 - e^{-Lk}}{1 - e^{-ky}} \right) \right]^n \frac{1 - e^{-Lk}}{1 - e^{-ky}} dk.$$

In general,

$$\lim_{n \rightarrow \infty} \left| \int v(t)^n \ln(v) dk \right|^{\frac{1}{n}} = \max(v(t)).$$

In the present case, $v(t)$ is stationary at a real value of t , which is minimum by the real criterion, but maximum in the imaginary direction. It is given by

$$\frac{v}{\bar{v}} = \frac{1}{2} + \frac{\pi e^{-Lk}}{1 - e^{-ky}} = 0,$$

so that, if v_0 is large, then $(v_0/\bar{v}) - 1 = -(v_0/\bar{v}^2 - 1) \rightarrow 0$ as $t \rightarrow \bar{v}/k$, and (by means of Stirling's approximation in the form $(v)^{1/n} \sim v/v$)

$$\begin{aligned} P_{n,n}(L)^{1/n} &\sim \frac{\pi e^{-Lk}}{2} V_{n,n}(L)^{1/n} \\ &\sim \frac{C}{n} e^{-Lk} \sqrt{\frac{1 - t^2}{v(t)}} = 1 - e^{-Lk}. \end{aligned}$$

Hence:

$$P_{n,\ell}(t) = \exp(-(\pi t)^{\alpha/\ell}).$$

Thus to have a reasonable probability of full coverage, we need $\pi \approx \sigma^2$. In particular, for the full human genome, with $\pi \approx 20$ cM for X-linked recombination, $L_1\pi \approx 2000/20 = 100$, another fact $\approx \pi/100$, leading to the requirement $\pi \approx 1$ or $\pi \approx 1120$. Substantially higher values of π will get the probability very close to 1.

In further detail, we may study the mean and the standard deviation of the proportion C_n of the genome covered by n randomly placed intervals of length $\ell_1, \ell_2, \dots, \ell_n > x$ (Measuring together with $\sum \ell_i = L$). It can be shown (Kobza, 1944), and is certainly reasonable, that the expectation of the uncovered part is given by

$$E(1 - C_n) = p_{1,n},$$

the probability that a randomly placed location is not covered by n randomly placed intervals, uniform over the whole genome. Similarly,

$$E((1 - C_n)^2) = p_{2,n},$$

the probability that two randomly placed sites are not covered by n randomly placed intervals. Let us look at $p_{1,n}$. There are two possibilities (1) a random site falls inside an interval, at least x cM from endpoint; this has probability (interval length/length) $= (\ell - x)/\ell = 1 - x/\ell$. The probability that this site does not overlap any of these intervals of course $(1 - x/\ell)^n$; or (2) a site falls $y < x/2$ cM from a chromosome end with probability x/ℓ , and is not covered by n intervals with probability $(1 - (x + x/2)/\ell)^n$. Summing, y from 0 to $x/2$ and adding the three contributions, we obtain

$$\begin{aligned} 1 - E(C_n) = p_{1,n} &= \left(1 - \frac{x}{L}\right) \left(1 - \frac{x}{\ell}\right)^n \\ &\quad + \frac{2x}{x+1} \left[\left(1 - \frac{x}{2\ell}\right)^{n+1} - \left(1 - \frac{x}{\ell}\right)^{n+1} \right]. \end{aligned}$$

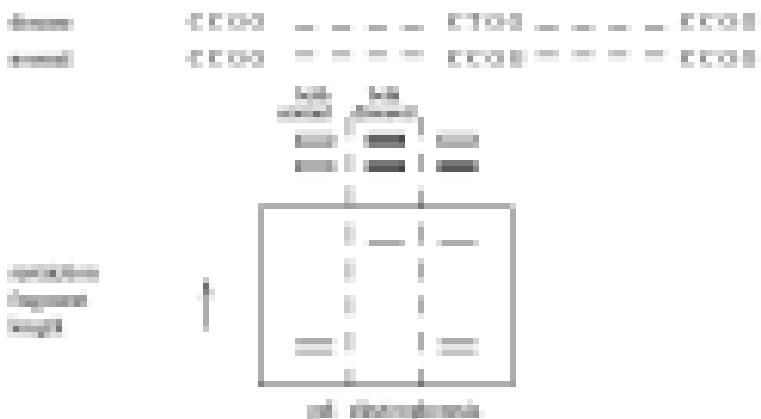
or, because $(1 - x/\ell)^n \approx e^{-x/\ell} \approx e^{-x^2/L^2}$ in the continuous limit,

$$E(C_n) \approx 1 - \left(1 - \frac{x}{L}\right) e^{-x^2/L^2} - \frac{2x}{x+1} (e^{-x^2/2L^2} - e^{-x^2/L^2}),$$

and for $x \rightarrow \sigma^2$, this is again dominated by the e^{-x^2/L^2} term.

With the coverage under control, we can take advantage of RFLPs in both medical and forensic directions. For the moment, we just note that RFLPs can be generated both by variation in the distribution of restriction sites, by a

very important subclass, the variation in number of tandem repeats (VNTRs), and more usually by repeats of the variation of microsatellites with no genetic implications. The tandem repeats are genes that were initially repeated by recombination of pieces and then accumulated a hyperdistribution by crossing over; their multiple nature allows them to be preferentially selected by hybridization. A typical situation involving the transmission of an altered gene (here, one C is replaced by a T) is shown (Koszinai et al., 1990).



A CTCG1 cutter acting on the two-strand genome of these chromosomes (white for normal, gray for mutated) produces three characteristic restriction-fragment signatures. If the disease gene is recessive, their (two copies on the two strands are necessary) progeny, both diseased, will have restriction digests that might go as



and so the genetic risk is readily assessed.

Assignment 3

1. In inappropriate assembly, suppose that the close density ($c(j)$) value increases linearly from 0 at each end of the genome to a maximum at the center.

- Compare $A(P_i)$ with that obtained by using the average $\bar{c}(P)$ and explain the result.
1. In the anchoring technique, find the probability that no gap between island-starts is four pairs.
 2. The problem of superglue in RPLP can be solved as a saddle-point estimation. Plot the coefficient in front of the exponential.

3.4. Pooling

A major use of a sequenced-clone library is to answer the following question: Where on the DNA, or chromosome, or large sequenced fragment (each comprising a set of which the library contains fully overlapping subsets) is a given protein, or portion thereof produced, or more generally the same for a set of associated proteins? For this purpose, we can, e.g., translate back to a relevant piece of mRNA, to use as a probe whose binding says that we have located the position of the “target” in question. One obvious procedure is to check every clone for binding by the probe, but this is slow, tedious, and error-prone. Instead, we do this by pooling (Balding and Tolley, 1997; Piatou et al., 1999). There are two general types of pooling:

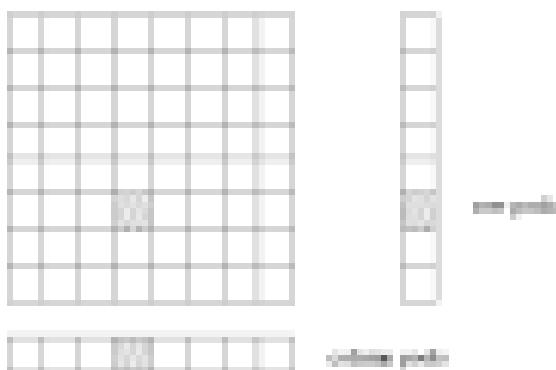
1. Adaptive pooling, most typically performed often on the traditional technique of detecting a (theoretical) counterfeit coin. We mix half the clones



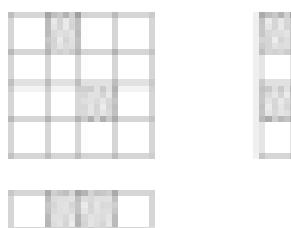
into one pool, half into another, and check each half for binding. We divide the clones of the “positive” pool (containing just one) into two halves, check antipool, etc. The desired clone, of the N in the library, is hence located in 2^k bins, where $2^k = N$, or $2^k = \ln N/\ln 2 = \log_2 N$. However, this requires the rotation of records, has many sequential operations, must be redone for each subsequent k to be detected, and is also error-prone.

2. **Phenotypic row-pool phasing.** An experimental design is chosen in which each pool contains predesignated clones, the whole set of pools is observed, and data are collected. No further experiments are carried out.

Example: Row and column designs. $N = r \cdot c$ clones are placed at the cells of an r -row c -column grid. Rows and columns are pooled separately. If the

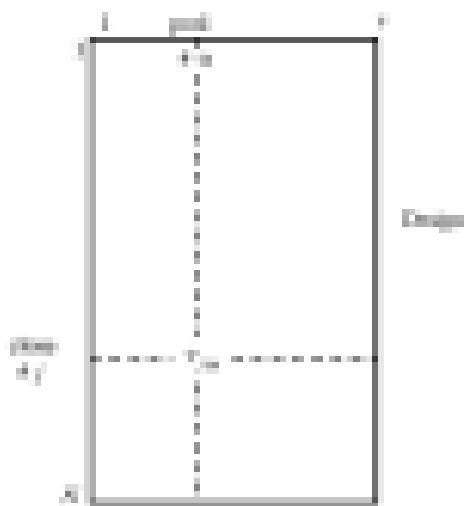


"word" is rare enough and the clone overlap is small enough, one column and one row will show positive, and no-the word-containing clone is located. Optimally to minimize the required number of pools, $r = c = N^{1/2}$, so there are $\sigma = 2N^{1/2}$ pools. We can in principle use a D -dimensional hypercube design with $(D - 1)$ -dimensional state pools, so $\sigma = D N^{(D-1)/D}$, minimized $D = \ln N$; then $\sigma = \sigma \ln N$ is quite small, but the setup is too complicated. Of course, when the word appears in several clones two cannot have less than



this coverage, neither overlapping row nor column red positions because only the sum of rows and the sum of columns are known. We resolve this ambiguity by applying the procedure to subsets of clones or by then examining separately each unresolved clone, partially but minimally adaptive.

What is an optimal "block-based" design to reduce the pool number at a given level of association (number of contacted clients)? This should presumably be an "on-the-average" strategy, as the identity of the word-carrying client is to be regarded as unknown. We examine this issue mainly by replacing probability arguments by explicit Boolean operations, reducing the



ad-hoc and strengthening the modularistic components of our discussion. We start by generalizing square incidence matrix to define our experimental designs. Here

$$\pi_{ij} := \begin{cases} 1 & \text{if client } j \in \text{pool } i \\ 0 & \text{if } j \notin \text{pool } i \end{cases}$$

Which client are positive, i.e., contains the word being probed, is not known to us, but the vast majority of pools will have most of the few positive clients. These negative pools comprise the full set of residual or definitely negative clients. When the negative pools and residual negative clients are studied, we get a new small matrix $\hat{\pi}$ that still contains all the positive information. The simplest design criterion, and the one we adopt here, is to minimize the number of clients left over, as they are the "noise" in which the positive clients are embedded. (A presumably better criterion, but harder to design for, recognises that some positive pools in $\hat{\pi}$ may contain just one remaining client; such a client is definite or reached positive, and it makes sense to try to maximize their number — but we will not do so here.)

The way that probability enters the picture at this stage is that if the small target word occurs only once in the full fragment, but the library coverage is c , it would be expected to occur c times in the library. The best we can do is introduce a class designator:

$$\eta_i = \begin{cases} 1 & \text{if done } i \text{ in } (+) \\ 0 & \text{if done } i \text{ in } (-) \end{cases}$$

$$i = 1, \dots, M.$$

With the knowledge that c of the η_i are 1 on the average, the rest being 0. Thus the average of any of the η_i is c/M , and they are independent of each other. In other words, with $\langle \cdot \rangle$ denoting expectation,

$$\langle \eta_i \rangle = \eta_i/M = c/M.$$

The pair of sets (ν_{ij}, η_i) of known unknown quantities, 0 or 1, completely specifies the setup. Let us develop the required derived quantities. To start with, we define

$$\Omega_{ij} = \begin{cases} 1 & \text{if no class } i \text{ is part in } (+) \text{ or is negative} \\ 0 & \text{if at least one class } i \text{ is in } (+) \text{ or is positive} \end{cases}.$$

This is clearly given by the expression

$$\Omega_{ij} = \prod_{k=1}^M (1 - \nu_{kj} \eta_{ik}),$$

which is 1 iff and only iff every $\nu_{kj} \eta_{ik} = 0$, i.e., either $\nu_{kj} = 0$, i not in k or $\eta_{ik} = 1$ but $\nu_{kj} = 0$, j is in i but i is negative).

Next, a class i' will be passed or used red negative if it is a member of some negative pool, i.e., if $\nu_{ij} = 1$ and $\Omega_{ij} = 1$ for at least one j .

i is marked negative if

$$A_{ii'} = \Omega_{ij} \nu_{ij} = 1 \text{ for at least one } j.$$

This is equivalent to saying that

$$\prod_{j=1}^M (1 - \nu_{ij}) = 0.$$

However, if

$$P_i = 1 - \prod_{j=1}^M (1 - \nu_{ij}),$$

Now

$$P_i = \begin{cases} 1 & \text{if } i \text{ is resolved negative} \\ 0 & \text{if } i \text{ is not resolved negative} \end{cases}$$

Consequently, the probability that i is resolved negative,

$$\Pr[\text{E}_i \text{ is resolved negative}] = 1 - \left(\prod_{j=1}^n (1 - A_{ij}) \right),$$

As there is no a priori way of distinguishing the characteristics of two different clauses, we might as well focus on clause #1:

$$\Pr[\text{E}_1 \text{ is resolved } (+)] = 1 - \left(\prod_{j=1}^n (1 - A_{1j}) \right)$$

$$\text{where } A_{1j} = Q_j \cdot e_{1j} \quad \text{and } Q_j = \prod_{i=1}^m (1 - b_{ij} q_{ij}).$$

A bit more conveniently, $\Pr[\text{E}_1 \text{ is } (+)] = p$ means that $\Pr[\text{E}_1 \text{ is } (-)] = 1 - p$; if we subtract the resolved $(+)$ clauses from the merely $(-)$ clauses, we get the actual but unresolved $(-)$ clauses.

$$\Pr[\text{E}_1 \text{ is unresolved negative}] = \left(\prod_{j=1}^n (1 - A_{1j}) \right) = p.$$

Carrying out this averaging requires expanding the product $\prod_{j=1}^n (1 - A_{1j})$ and corresponds to the “inclusion–exclusion” theorem of probability; see, e.g., Feller (1971). The expansion consists, to within sign, of terms A_{1j} , pair terms $A_{1j_1} A_{1j_2}$, triplets $A_{1j_1} A_{1j_2} A_{1j_3}, \dots$, with the condition that no two indices are equal, and that permuting the order $A_{1j_1} A_{1j_2} A_{1j_3} \leftrightarrow A_{1j_3} A_{1j_1} A_{1j_2}$ does not create distinct terms. We can take care of the latter by allowing everything, but dividing by $l!$, the number of permutations of l distinct letters. In other words,

$$\begin{aligned} \prod_{j=1}^n (1 - A_{1j}) &= 1 - \sum_{i=1}^l A_{1j_i} + \frac{1}{2!} \sum_{\{i,j\} \subset [l]} A_{1j_i} A_{1j_j} \\ &\quad - \frac{1}{3!} \sum_{\{i,j,k\} \subset [l]} A_{1j_i} A_{1j_j} A_{1j_k} + \dots, \end{aligned}$$

and we conclude that

$$\Pr[\text{E}_1 \text{ is unresolved negative}] =$$

$$1 - p - \sum_{i=1}^l (A_{1j_i}) + \frac{1}{2!} \sum_{\{i,j\} \subset [l]} (A_{1j_i} A_{1j_j}) - \frac{1}{3!} \sum_{\{i,j,k\} \subset [l]} (A_{1j_i} A_{1j_j} A_{1j_k}) \dots$$

For the evaluation, we see first that, because $r_{ij}^2 = r_{ij}$, then $A_0 = r_{00} \prod_{i=1}^N (1 - r_i/r_{00}) = r_{00} (1 - r_1/r_{00}) (1 - r_2/r_{00}) \cdots (1 - r_N/r_{00})$, or

$$A_0 = (1 - r_1/r_{00}) r_{00} \prod_{i=2}^N (1 - r_i/r_{00}).$$

However, $\langle r_0 \rangle = \cdots = \langle r_i \rangle = p$, so:

$$\langle A_0 \rangle = (1 - p)r_{00} \prod_{i=1}^N (1 - pr_{00}).$$

Thus, using $r_0^2 = r_0, \dots, r_i^2 = r_i$, we have

$$\begin{aligned} A_0 \langle A_0 \rangle &= (1 - r_1^2/r_0) r_{00} r_{00} \prod_{i=2}^N (1 - r_i/r_{00})(1 - r_i/r_0) \\ &= (1 - r_1/r_0) r_{00} r_{00} \prod_{i=2}^N (1 - r_i(1 - (1 - r_{00})(1 - r_0))) \end{aligned}$$

and similarly

$$\begin{aligned} A_0 \langle A_0 \rangle A_1 &= (1 - r_1/r_0) r_{00} r_{00} r_{00} \prod_{i=3}^N (1 - r_i(1 - (1 - r_{00})(1 - r_0)(1 - r_1))), \\ &\vdots \end{aligned}$$

etc., and we conclude that, in general,

$$\langle A_0, A_1, A_2, \dots \rangle$$

$$= (1 - p/r_{00}, r_{00}/r_1, r_1/r_2, \dots) \prod_{i=1}^N (1 - p + pr(1 - r_{00})(1 - r_0)(1 - r_1)\cdots)$$

We therefore have the relatively simple result that, at fixed design,

$$\frac{1}{1-p} \Pr\{1\text{ is measured in }q\}$$

$$\begin{aligned} &= 1 - \sum_{n_0} \left\{ r_{00} \prod_{i=1}^N (1 - p + pr(1 - r_{00})) \right\} \\ &+ \frac{1}{p} \sum_{n_1} \left\{ r_{00} r_{00} \prod_{i=2}^N (1 - p + pr(1 - r_{00})(1 - r_0)) \right\} \\ &- \frac{1}{p^2} \sum_{n_2} \left\{ r_{00} r_{00} r_{00} \prod_{i=3}^N (1 - p + pr(1 - r_{00})(1 - r_0)(1 - r_1)) \right\} \\ &+ \dots \end{aligned} \tag{2.1}$$

How then in the world do we use Eq. (2.1) to, e.g., plan a design that minimizes the number of points required for achieving a given resolution? Any complicated algorithm would defeat the whole objective. A neat way is to hedge our bets and choose the design in some random fashion. The simplest is random binomial. Just choose a parameter b so that each draw occurs on the average in b points. Hence $\sum_{i=1}^n m_i = b$ for any i , or because the m_i for different i are independent,

$$\langle m_i \rangle = b/n, \quad \text{independently.}$$

The averaging of Eq. (2.1) is then next to trivial. The sum $\sum_{A_1, A_2, \dots, A_M}$, (A_1, \dots, A_M) contains $n(r-1)(r-2)\dots(r-M-1)$ terms, all of which have the same value $n(r-1) - (r-1)(r-2)\dots(r-M) = n!/((r-n)!)$, the binomial coefficient $\binom{n}{r}$, and every v is replaced by d/v in the average. We thus have

$$\begin{aligned} & \frac{1}{1-p} \Pr[1 \text{ is unmeasured negative}] \\ &= 1 - \binom{r}{1} \binom{M}{r} \left[1 - p + p \left(1 - \frac{d}{v} \right) \right]^{M-1} \\ &+ \binom{r}{2} \binom{M}{r}^2 \left[1 - p + p \left(1 - \frac{d}{v} \right)^2 \right]^{M-1} \\ &- \binom{r}{3} \binom{M}{r}^3 \left[1 - p + p \left(1 - \frac{d}{v} \right)^3 \right]^{M-1} + \dots \end{aligned}$$

Again, M is large, but p is small, and $p(1 - (1 - \frac{d}{v}))$ even smaller, so we can use $(1 - x)^{M-1} \rightarrow e^{-px}$. Finally because $1 - p = \Pr[1 \text{ is negative}]$, we can write

$$\begin{aligned} P_{\text{un}} &\equiv \Pr[1 \text{ is measured } (-)] \Pr[1 \text{ is } (-)] \\ &= \frac{1}{1-p} \Pr[1 \text{ is unmeasured } (-)] \\ &= \sum_{i=0}^r (-1)^i \binom{M}{i} \binom{M}{r}^i \exp[-v \left(1 - \left(1 - \frac{d}{v} \right)^i \right)], \quad (2.2) \end{aligned}$$

where $Np = c$ is just the coverage.

This is a nice explicit formula, but it has r terms, and they almost sum up if the exponential were constant, we would have $\sum_{i=0}^r (-1)^i \binom{M}{i} = (1 - 1)^M = 0$. Therefore numerical evaluation and tabulation are problems. We postpone the numerical evaluation to at the conclusion of analytical processing, so that, at the very least, we can get a better feeling for parameters.

dependence. A number of analytic tricks are available. Here is one. We introduce the asymmetric difference operator

$$\delta_0^k f(t) = k f(t+1) - f(t) \approx (k, k-1) f(t),$$

and hence

$$(\delta_0^k)^T f(t) = (k, k-1)^T f(t) = \sum_{j=0}^{k-1} (-1)^{j+k} \binom{k}{j} k^j f(t+j).$$

so that in fact

$$\begin{aligned} R_{\text{max}} &= e^{kt} (-1)^k (\delta_0^{k+n})^T e^{(n-k)t} h(t) \\ &= e^{kt} (-\delta_0^{kn})^T \sum_j c_j(t) \left(1 - \frac{t}{n}\right)^m \Big|_{t=0}. \end{aligned}$$

However,

$$\begin{aligned} -\delta_0^{kn} \left(1 - \frac{t}{n}\right)^m &= \left(1 - \frac{t}{n}\right)^m - \frac{t}{n} \left(1 - \frac{t}{n}\right)^{m+1} \\ &= \left[1 - \frac{t}{n} \left(1 - \frac{t}{n}\right)\right] \left(1 - \frac{t}{n}\right)^m. \end{aligned}$$

so

$$\begin{aligned} R_{\text{max}} &= e^{kt} \sum_j \frac{c_j'}{n!} \left[1 - \frac{t}{n} \left(1 - \frac{t}{n}\right)\right]^m \\ &= \sum_j \left(1 - \frac{t}{n} e^{kt}\right)^m \frac{c_j'}{n!} e^{kt} \sim \sum_j e^{kt m} \frac{c_j'}{n!} e^{kt}, \quad (2.2) \end{aligned}$$

where

$$c_j' := -\ln \left(1 - \frac{t}{n}\right).$$

Now all terms are positive, so computation is quick and accurate.

In practice, a stratified random-influence design, in which we impose a fixed number of points ordered by each class, $\sum_i n_{ij} = k$ for each j , has experimental advantages. It can be treated similarly, but the analysis is a bit more involved.

Assignment 3

1. Choose reasonable values of a , b , and n in Eq. (2.2) and carry out the evaluation at a few levels of precision to see how much computational accuracy is needed.

2. Because Eq. (2.2) has only positive terms, analytical approximations are straightforward. They can be found and compared with problem 1.

2.4. Reptiles

There are, as we have seen, various ways of joining clones and building up islands for the purpose of constructing maps of the genome, and many are being suggested all the time. A dominant role in the increasing common whole genome "shotgun sequencing" is now being played by the use of end-characterized clones, i.e., those in which a few base-pair ends at each end are identified, base after base, by use of the old Sanger technique (Chapman et al., 1997). Then overlap is recognized if the leading end of a "new" clone overlaps, by at least T base pairs, the trailing end of the last clone of the currently developed island. Rather than develop a new clever technique of analysis such time, it might pay to create a descriptive machinery that is tailored to the general concept and therefore allows closure questions to be answered more readily. This is not a novel concept, but let us carry it a bit further than is usually done (Peters and Peters, 1999). We start here with only the basic fingerprint assembly of fixed clone length L and overlap criterion T , but in a fashion that extension to a distribution of clone size and overlap-threshold can be carried out with relative ease.

If genome ends effects are unimportant, we can imagine that the island-building process starts with a clone whose left end is at site 1, the right end at its length L . Now we build up an island by putting down a clone of step k , starting at site k , with probability $p_{k+1}(N/G)$. If it overlaps the current ($k-1$) island sufficiently, a new island with right end at $R_k = k + L - 1$ will be produced. If it does not, the island terminates at the previous value R_{k-1} . If no clone is found starting at k , with probability $q = 1 - p$, then $R_k = R_{k-1}$.

The first crucial criterion is that of sufficient overlap, T . In the basic case we are considering, this simply requires that a clone be found starting at site k such that

$$R_{k-1} \geq k + T - 1.$$

The complementary criterion is that of the island terminating; it will do so



at step $k - 1$ as soon as $R_{k-1} = k + T - 2$, for then a clone starting at k can have an overlap of at most $T - 1$. Hence, at step k ,

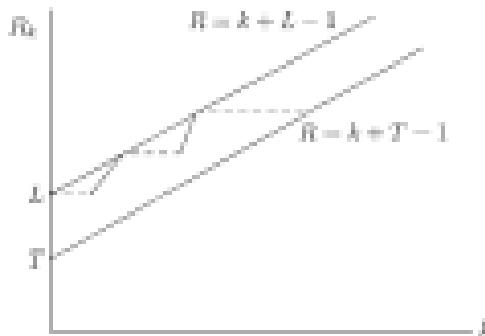
$$R_k = k + T - 1 \text{ : stop.}$$

Pictorially, then, the island develops as a random walk in which at each step, $\Delta k = 1$, either a clone is added, so that R_k jumps to $k + L - 1$, or no clone is added, so that $R_k = R_{k-1}$, but the island stops if this results in $R_k = k + T - 1$.

What we want to do now is to find

$$P_k(R),$$

which is the probability that our developing island at step k has length R .



Because an island is completed when $R_k = k + T - 1$, the expected length of an island will be just

$$E(L_I) = \sum_{k=1}^{\infty} (k + T - 1) \cdot P_k(k + T - 1).$$

$P_k(R)$ is to be found iteratively. If the island has "arrived" at the point (k, R) , it can have done so in two ways. Either (1) $R = k + L - 1$, corresponding to an overlapping clone appearing at step k , and this has probability p providing that R_{k-1} was anything between its minimum $(k-1) + T$ and its maximum $(k-1) + (L-1)$,

$$P_k(k + L - 1) = p \sum_{R'=k-1+T}^{k-1+L} P_{k-1}(R'),$$

or (2) $R = R_{k-1}$ because no clone appears at step k , and of course R_{k-1} cannot have been smaller than $(k-1) + T$ (or the island would already have

terminated).

$$P_2(R) = q \cdot P_{2-1}(R) \quad \text{if } R \geq d-1+T, \quad R \leq d-1+L.$$

Finally, the whole process starts at:

$$P_2(R) = \delta_{R, L}.$$

The Kronecker delta $\delta_{i,j} = 1$ if $i = j$, otherwise 0.

The "walk" that we are analysing is now with a reflecting barrier along the slope $R = d + L - 1$, an absorbing barrier along $R = d + T - 1$. Analysis is much easier if the barriers are fixed, so we set:

$$P_2(R+j) = Q_2(j),$$

and this has:

$$\begin{aligned} Q_2(L-1) &= p \sum_{i=0}^{d-1} Q_{2-i}(1), \\ Q_2(j) &= q \cdot Q_{2-1}(j+1) \quad \text{for } j \geq T-1, \quad j \leq L-2, \\ Q_2(1) &= \delta_{j,L-1}, \end{aligned}$$

subject to which

$$Q_2(L_j) = \sum_{i=1}^{\infty} (t+T-1) \cdot Q_2(T-1).$$

For solving this, the method of choice is usually a generating function. We set:

$$\begin{aligned} Q(x, j) &= \sum_{i=1}^{\infty} x^i Q_2(i) \\ &= x \cdot \delta_{j,L-1} + \sum_{i=1}^{\infty} x^i Q_2(i). \end{aligned}$$

Now we multiply the recursive relations by x^j and sum over it from 2 to ∞ . We see that:

$$Q(x, L-1) = x + p \sum_{i=0}^{d-1} x \cdot Q(x, i),$$

$$Q(x, j) = q \cdot x \cdot Q(x, j+1) \quad \text{for } L-2 \leq j \leq T-1,$$

$$Q(L_j) = t \cdot Q(x, T-1) \delta_{j,L-1} + (T-1) \cdot Q(x, T-1).$$

Working up from $Q(x, T-1)$ in the second relation, we have $Q(x, j) = Q(x, T-1) q^{j-T+1}/x^j$ for $j \leq T-1$, and substituting this into the first

relation, we obtain

$$\begin{aligned}\frac{Q(x, T=1)}{(q/p)^{k+T}} &= x + px Q(x, T=0) \sum_{t=0}^{k+T} \frac{1}{(q/p)^{k+t+T}} \\ &= x + \frac{p}{q} Q(x, T=0) \left[\left(\frac{1}{q/p}\right)^{k+T} - 1 \right] / \left(\frac{1}{q/p} - 1\right),\end{aligned}$$

so that

$$\begin{aligned}Q(x, T=1) &= [x(1 - qx(p/q)^{k+T})/(1 - qx - px(1 - qx)^{k+T})] \\ &= [x(1 - qx(p/q)^{k+T})/(1 - x + px(p/q)^{k+T})].\end{aligned}$$

We see that $Q(1, T=1) = 1$, which just says that the probability that the biased coin comes up heads is 1. What we need then is $Q(0, T=1)$, and a brief calculation shows that

$$Q(0, T=1) = \frac{1}{p} [p^{-(k+T)} - q],$$

from which we conclude that

$$E(L_1) = T + \frac{1}{p} [p^{-(k+T)} - 1],$$

or, because $p^{-(k+T)} = (1 - p)^{k+T} \approx e^{-pkT}$ for small p , that

$$E(L_1) = T + \frac{1}{p} [e^{-pkT} - 1],$$

which we have already seen.

What deeper questions can we ask? The simplest might be about the dispersion, or standard deviation, of L_1 : after all, a mean of 1.0 is nothing guarantee if 10% of the population are 1.0 and the rest are 0;

$$\begin{aligned}v(L_1) &= [E(L_1^2) - E(L_1)^2]^{1/2} \\ &= \left[\sum_T T^2 Q_0(T=1) - \left[\sum_T T Q_0(T=1) \right]^2 \right]^{1/2},\end{aligned}$$

(the constant part of L_1 does not contribute)

$$\begin{aligned}&= [120 + Q_1'(1) - Q_1''(1)/2]^{1/2} \\ &= [10 \ln Q_1'(1) + 10 \ln Q_1''(1)]^{1/2},\end{aligned}$$

where the quantity $T - 1$ is invited. After some minor algebra,

$$\begin{aligned} \mathbb{E}(L) &= \left[\left(T + \frac{q^T}{p^T} - \frac{2L-1}{pq^{T-1}} + \frac{1}{pq^{T-1}} \right) + \left(T - \frac{q}{p} + \frac{1}{pq^{T-1}} \right) \right]^{1/2} \\ &= \left[\frac{q^T}{p^T} - \frac{q}{p} - \frac{2L-1}{pq^{T-1}} + \frac{1}{pq^{T-1}} \right]^{1/2} \\ &\rightarrow \frac{1}{p} e^{\mu L - \nu} \left[1 - \left(L - \frac{1}{2} \right) \mu e^{-\mu L - \nu} \dots \right]. \end{aligned}$$

In fact very close to $\mathbb{E}(L)$ itself, as in an exponential, or basic survival time, distribution – which this is very close to.

It is a bit more complicated to restrict the islands to contigs, which are the objects we'd really pick up. This means counting those islands that have not hit the upper border, i.e., that have not encountered at least one weight- p jump. For this purpose, we simply refrain from using the fact that $p + q = 1$ in our first expression for $Q(x, T - 1)$:

$$Q_0(x, T - 1) = q^{T-1} x^{L-1} (1 - qx)(1 - (p + q)x + px) q^{L-1} x^{L-1}.$$

The coefficient of x^L then represents the weight of the $(k+1)$ -close islands. Subtracting out $Q_0(x, T - 1)$ gives the resulting generating function:

$$\begin{aligned} Q_{\text{cont}}(x, T - 1) &= Q(x, T - 1) - Q_0(x, T - 1) \\ &= q^{T-1} x^{L-1} px^L (1 - qx)^{L-1} (1 - x + px) q^{L-1}, \end{aligned}$$

which is no longer normalised: $Q_{\text{cont}}(1, T - 1) = 1 - q^{L-1}$, so that

$$Q_{\text{cont}}(x, T - 1) = \frac{q^{T-1}}{1 - q^{L-1}} px^{L+1-T} \frac{1 - qx^{L-T}}{1 - x + px} q^{L-1}.$$

Hence

$$Q_{\text{cont}}(1, T - 1) = 1 + \frac{1}{p} q^{L-1} - (L - T) \frac{q^{L-1}}{1 - q^{L-1}},$$

$$\mathbb{E}(L_{\text{cont}}) = \frac{1}{p} q^{-L+T} + \frac{T - L q^{L-T}}{1 - q^{L-T}}.$$

The advantage of this route of derivation is that generalisation to a distribution of close lengths, $f(L)$, and a distribution of overlap detection thresholds, $w(T)$, is easy to carry out. Although a complete closed-form analysis is available in only special cases, means and variances of required characteristic func-

be obtained quite generally. For example, for the mean island length, we find

$$E(L_A) = \sum_{k=1}^{\infty} k P(k) / \prod_{j=1}^{A-1} \pi(j),$$

where

$$P(k) = \sum_{n=1}^{\infty} f(n), \quad q(n) = 1 - p(n)/P(n).$$

In the continuum limit, in which the unit length, say T , corresponds to very many base pairs, writing

$$f(t) = \lim_{n \rightarrow \infty} f(n) T^n / T, \quad W(t) = w(t/T), \quad \rho = T/p,$$

we see that this reduces to

$$w(t) = \int \rho(u) \exp\left(-c \int_u^t \rho(u') w(u') du'\right) du,$$

showing that the generalization in Section 2.1 from the fixed island length, $P(k) = \delta(k - K)$, fixed threshold, $w(t) = \delta(t - T)$, case works nicely.

3

Sequence Statistics

Once we can assume that long stretches of DNA are completely sequenced, it is possible to analyze the information contained therein. In the "language" of DNA, there are many "verbal" biases – dialects, accents, pauses – that have evolved for broad reasons of physical and biochemical accessibility and function. They contribute to the "default state," the random or null hypothesis with respect to which additional information must be assessed. Ignorance of this background bias not only ignores available information but also poses as noise against which a signal must compete. In brief synthesis, as we have noted, we know that DNA is not homogeneous, but contains regions that code for specific proteins as well as those with associated regulatory functions. Each of the former is divided into exons, which place bases at a time) are transcribed and translated into amino acids, domains, which are spliced out during this process, and untranslated subregions. The latter are often confined to a flanking region of this coding section, all resulting in a set of "junk DNA," which may or may not be functional.

3.1. Local Properties of DNA

To start with, the bases are not even at equal frequencies. For example, for the human mitochondrial (16,000 bp), we have typically (see, e.g., Weis (1993), Chap. VII)

A	C	G	T
0.31	0.19	0.23	0.37

For different regions of the human fetal globin gene, the pair (A_1, A_2) , the distributions are given in the following table.

	Total Length	A	C	G	T
U random (1)	1000	0.25	0.25	0.25	0.25
U random (2)	1000	0.25	0.25	0.25	0.25
Uniform (1)	1000	0.25	0.25	0.25	0.25
Uniform (2)	1000	0.25	0.25	0.25	0.25
Independent (1)	1000	0.25	0.25	0.25	0.25

Clearly, this bias should be taken into account, and we should bunch together, for statistical purposes, only subsequences of similar character. However, even if this is not done (assuming 100,000 bp from the vertebrate sequence as an example), an additional local structure immediately appears. If p_i is the relative frequency of base i , p_{ij} of the pair $i \cdots j \cdots i$, then the pair correlation ratio

$$p_{ii}/(p_i p_j)$$

is not at all unity, as it would be for independent placement, but rather is as in the following table:

	A	C	G	T	Second base
A	1.15	0.84	1.36	0.89	
C	1.15	0.98	0.41	1.28	
G	1.00	0.99	1.11	0.92	
T	0.95	0.99	1.29	1.07	

Observe the very small CG frequency, which presumably is due to genetic "mismatches."

Bunching together sequences of different single frequencies will maximize the pair correlation, e.g., $C_1 C_2 + C_2 C_1 = (C_1 C_2)^2$ does not have the form $N_1 N_2$ when the full sequence consists of two internally uncorrelated sequences; hence the tacit assumption of homogeneity is implicit in assessing correlations.

A Fisher test for pair correlations is to stick to a single functional entity, here a dinucleotide gene state, organized again according to the number of successive pairs of mutations of the 16 possible types.

	A	C	G	T	Total
A	22	36	22	32	97
C	32	51	14	47	145
G	22	36	26	39	123
T	3	37	40	34	114
Total	87	144	117	109	457

Note that the N_{ij} in this contingency table, the actual number of times $\{j\}$ is found, do not quite satisfy $N_{ij} = N_{ji}$, etc., where $N_{ij} = \sum_k N_{ijk}$, $N_{ji} = \sum_k N_{kij}$, because the left term of the leftmost pair of the chain cannot be the right term of any pair, etc.

The probability law for independent trials (probabilities in the p^2 -rule) is the p^2 -rule (Fisher, 1949). This goes as follows. Suppose the result of a measurement is indicated by α , with N_α the number of times α occurs, and of course $\sum_\alpha N_\alpha = N$, for a type of result, α , the total number of trials. If the probability of α is p_α , then, on the assumption of independence, the set (N_α) occurs with a probability

$$\left(N! \prod_{\alpha} \frac{1}{N_\alpha!} p_\alpha^{N_\alpha} \right) \prod_{\alpha} p_\alpha^{N_\alpha}.$$

If N is large, we let $N_\alpha = N x_\alpha$, which converts the above to

$$P(N_\alpha) = \left[N! \prod_{\alpha} \frac{1}{N_\alpha!} (p_\alpha N)^{N_\alpha} \right] \prod_{\alpha} p_\alpha^{N_\alpha},$$

and we then use Stirling's formula $n! \sim \sqrt{2\pi n} (n/e)^n$ to obtain

$$P(N_\alpha) \sim \frac{N!}{(\prod_\alpha N_\alpha!)^2 N^N} \left[\prod_{\alpha} \left(\frac{N x_\alpha}{N_\alpha} \right)^{N_\alpha} \right]^2 / \left(\prod_{\alpha} x_\alpha \right)^{N^2}.$$

Now $x_\alpha \ln(p_\alpha x_\alpha/N_\alpha)$ has a single maximum at $x_\alpha = p_\alpha$,

$$x_\alpha \ln(p_\alpha x_\alpha/N_\alpha) = p_\alpha - \frac{1}{2p_\alpha} D_\alpha = p_\alpha^2 - \dots,$$

so that the N^2 power in $P(N_\alpha)$ has a very sharp maximum, and

$$P(N_\alpha) \sim \frac{N! (N^N e^N)}{\prod_\alpha (2\pi N_\alpha p_\alpha)^{N_\alpha}} e^{-N/2 \sum_\alpha (N_\alpha - p_\alpha N)^2/p_\alpha},$$

depending on only the quantity

$$\begin{aligned} x^2 &= N \sum_{\alpha=1}^m \frac{(N_\alpha - p_\alpha N)^2}{p_\alpha} = \sum_{\alpha=1}^m \frac{(N_\alpha - p_\alpha N)^2}{p_\alpha N} \\ &= \sum_{\alpha=1}^m \frac{(N_\alpha - \langle N_\alpha \rangle)^2}{\langle N_\alpha \rangle}, \end{aligned}$$

where we have used (and will use) the notation $\langle N_\alpha \rangle$ interchangeably with

$\Delta(N_1)$: Introducing the new variable

$$\beta_{ij} = \delta^{1/2} \frac{p_i - p_j}{\mu_i^{1/2}},$$

we see that the distribution is independent of the p_{ij} . For large N , we may regard the β_{ij} as continuous variables with a charge $\Delta N_1 = 1$, and have $d\alpha_{ij} = N^{1/2} \beta_{ij}^{-1/2} d\beta_{ij} = (\Delta p_{ij})^{1/2} d\beta_{ij} = (\Delta p_{ij})^{1/2}$. If the (p_{ij}) were independent, the probability for the (β_{ij}) would then be determined by $p(p_{ij}) dp_{ij} = P(K_1) dK_1$, or

$$P(\beta_{ij}) = N_1 \partial N_1 / \pi^N \left(\frac{1}{2\pi} \right)^{N/2} e^{-N/2 \beta_{ij}^2},$$

where

$$y^2 = \sum_{ij} \beta_{ij}^2.$$

The reason for the large prefactor $N_1 (\partial N_1)^N = \sqrt{2\pi N}$ will appear in a moment.

Now what is the distribution of y^2 ? The y^2 's are not free in n -dimensional y space because they are restricted to the hyperplane $\sum_i K_i = N$, or

$$\sum_{ij} \beta_{ij}^{1/2} p_{ij} = 0,$$

and it is the surface "thickness" of only $\sim N^{-1/2}$ on the restricted surface of the (β_{ij}) that is responsible for the above $\sqrt{2\pi N}$. However, we do not have to worry about normalization; it can be supplied in the last step. In general, suppose there are r linear homogeneous relations that the $K_i = (N_{ij})$, and hence the p_{ij} , have to satisfy. Then the space is r -dimensional dimensional

$$y = \mathbf{y} - \mathbf{z},$$

the number of degrees of freedom, but the intersection of the spherically symmetric unnormalized distribution with the hyperplane through the origin of r -dimensional y remains spherically symmetric in the r -dimensional space. Call the new coordinate in this space the (Z_{ij}) . Then we make the distribution of

$$y^2 = \sum_{ij} Y_{ij}^2$$

where

$$P(T_i) = \prod_{j=1}^i \left[\frac{1}{\sqrt{2\pi}} e^{-t_j^2/2\sigma^2} \right]$$

the constant is now obtained by the normalisation condition $\int \cdots \int P(T_i) dT_i = 1$. Because

$$\begin{aligned} \langle e^{-t_1 t_2} \rangle &= \prod_{j=1}^2 \int \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(t_j + t_2)^2/2\sigma^2} dT_j \\ &= (1 - 2e^{-\sigma^2}) \end{aligned}$$

defines $\int_0^\infty p(x^2) e^{-x^2/2\sigma^2} dx^2$, we readily find, on consulting inverse Laplace transform tables, that the x^2 distribution for n degrees of freedom is

$$p_n(x^2) = (2^{-n/2})^{-n} e^{-x^2/2\sigma^2} (x^2)^{n/2} \left(\frac{n}{2} - 1 \right)!,$$

from which it follows that, for any沉iable v , $p_n(x^2)$ is also sharply distributed, with $E(x^2) = v$.

Now back to our example. We want to test the independent joint hypothesis, i.e., whether (using the distinction between a left member of a pair and a right member, because of clockwise) the observed (N_{ij}) are consistent with the p_{ij} , computed as

$$p_{ij} = p_i p_j,$$

We do not actually know the p_{ij} and p_{ij} , we have to estimate them by means of $\hat{p}_{ij} := 1/W \sum_j N_{ij} = N_{ij}/W$, $\hat{p}_{ij} = N_{ij}/W$, and then the question is whether

$$(N_{ij}) \sim (W_i W_j)/W$$

is a valid estimate, modified of course by known restrictions. However, there are many restrictions. Because $\sum_i (N_{ij} - \frac{W_i W_j}{W}) = \sum_i (N_{ij} - \frac{W_i W_j}{W}) = 0$, as is readily verified, according to $(i = 1, \dots, n) = 1 - 1 = n$ independent restrictions, we compute

$$\chi^2 = \sum_{i,j} \left(N_{ij} - \frac{W_i W_j}{W} \right)^2 \int (N_{ij} W_j/W) - 5n \chi^2$$

In the example quoted, compared with the $E(\chi^2) = 2n - T = 2$ degrees of freedom. This shows that the assumption $p_{ij} = p_i p_j$ is extremely inconsistent with the data.

However, is the strictly pair-considering Markov chain assumption good enough? Maybe there are really three, or four, or more successive influences

base correlations. Let us use the notation:

$$\delta_x(j) = \begin{cases} 1, & \text{base } j \text{ is at site } x \\ 0, &quad & \text{base } i \neq j \text{ is at site } x \end{cases}$$

Then the multilite correlations are defined as

$$p_1 = \langle \delta_{x_1}(j_1) \rangle,$$

$$p_2 = \langle \delta_{x_1}(j_1) \delta_{x_2}(j_2) \rangle,$$

$$p_3 = \langle \delta_{x_1}(j_1) \delta_{x_2}(j_2) \delta_{x_3}(j_3) \rangle, \dots,$$

where $\langle \dots \rangle$ means taking the average of the multiplet over all possible chain positions and averaging. The corresponding conditional correlations are then defined as

$$P_{111} = p_1/p_0, \quad P_{1111} = p_2/p_0, \dots,$$

for independence, a Markov chain of order 0, $P_{111} = p_1$, or $P_{1111} = p_2/p_1 p_0, \dots$

for a Markov chain of order 1, $p_{1111} = p_{111}/p_0 = \frac{p_2}{p_0 p_1}$,

for a Markov chain of order K , $p_{111\dots 11} = p_{11\dots 11}/p_0$ ($K \leq L$), or

$$P_{11\dots 11} = \frac{P_{11\dots 11} P_{11\dots 11} P_{11\dots 11}}{P_{11\dots 11} P_{11\dots 11} \dots},$$

Because $p_{11\dots 11} = \sum_{k=0}^L P_{11\dots 11}$, the Bernoulli chain is fully specified by $4^{L+1} - 1$ independent parameters, and the constraint $p_{11\dots 11} \leq 1 \times 4^L$ (as $\sum_k p_k = 1$). Clearly we need a lot of data – i.e., very long chains – to determine these parameters, assuming homogeneity of the chain.

Perhaps we can at least get the intrinsic correlation length in the sense of the order of K of the underlying Markov chain (Katz, 1993), one technique is (Ding, 1993; Tavaré and Gelfand, 1994) the Bayesian information criterion (BIC). Consider the subsequence x_{L+1}, \dots, x_{L+K} of the test sequence as the (x, E') piece of data, $\text{data}(x, E')$, and choose some subset of the $\{(x, E')\}$. Given this data, then, according to the familiar simplified probability theorem, the a posteriori probability that a given E' is the correct order is

$$P(E' | (\text{data}(x, E'))) = P(\text{data}(x, E') | E') P(E') / P(\text{data}(x, E')).$$

However we want knowledge E' to maximize the likelihood $P(\text{data}(x, E') | E) = \prod_{i=L+1}^L P(\text{data}(x_i, E') | E)$, evaluated as

$$\prod_{i=L+1}^L P(x_{L+i} | x_{L+1}, \dots, x_L, E)^{2^{(L-i)+1}},$$

where $\pi(j_1, \dots, j_{K+1})$ is the number of times the subsequence $j_1 \dots j_{K+1}$ is encountered in the data. In other words, we want to maximize

$$\ln P(\text{data}|j_1, K) | K = \sum_{\substack{j_1, \dots, j_{K+1} \\ \text{in data}}} \pi(j_1, \dots, j_{K+1}) \ln \frac{\pi(j_1, \dots, j_{K+1})}{\sum_j \pi(j_1, \dots, j_K)}.$$

However, there are $2 < 4^K$ parameters that are implicitly determined by the process, such as the result of $\pi = [\sum_{j_1, \dots, j_K} \pi(j_1, \dots, j_K, j)]$ places of data. This suggests a parameter to control penalty function, and R_{max} is defined as the K that maximizes

$$R(K) = \sum_{\substack{j_1, \dots, j_K \\ \text{in data}}} \pi(j_1, \dots, j_K) \ln \frac{\pi(j_1, \dots, j_K)}{\sum_j \pi(j_1, \dots, j_K)} - \frac{1}{2} K < 4^K \ln n.$$

The penalty function is not unique, satisfying two weak information-theoretic criteria, but we can indeed show that

$$\lim_{n \rightarrow \infty} \Pr(R_{\text{max}} = R_{\text{true}}) = 1.$$

Example. The 48,500bp "head region" of phage λ. We have

K	0	1	2	3	4
$-2R(K)$	24,504	24,482	24,413	24,414	24,414

showing (weakly) that $K = 2$ is suggested, consistent at least with the triple structure of DNA as amino acid translation.

Once an estimate of K is made, we can return to estimate the parameters p_{j_1, \dots, j_K} . Because there are far too many for a reasonable statistical estimation, i.e., guaranteed homogeneous place of DNA, one way out is to model the parameter set (see, e.g., Rafferty 1987)

$$p_{j_1, \dots, j_K} = \sum_{i=1}^I \lambda_i q_{ij_i},$$

where

$$\sum_j \lambda_i = 1 \quad (\text{I = 3 parameters}),$$

$$\sum_{ij} q_{ij} = 1 \quad (\text{12 parameters}).$$

a total of only $I + 11$ parameters. These can then be found from the maximum-likelihood estimator, i.e.,

$$\hat{\lambda}_i = \sum_{\substack{j_1, \dots, j_K \\ \text{in data}}} \pi(j_1, \dots, j_K) \ln \sum_{j=1}^I \lambda_i q_{ij_i}$$

subject to

$$\sum_{k=1}^K \lambda_k = 1, \quad \sum_{k=1}^K p_{ik} = 1,$$

without the necessity of an additional penalty function.

3.2. Long-Range Properties of DNA

Regions of uniform A, C, G, and T distributions at low resolution are called isochors (Bernardi, 1999); they correspond to 10 base pair (10μ) to hundreds of thousands pairs. Because DNA consists a whole number of double strands, distributions are often specified by the G + C content of either or both strands, which shifts when the isochores switch. One way (Bernardi-Baldarelli, 1998) of detecting an isochores or other segmented structure is (Bernardi-Labeyrie et al., 2000) by means of the Jensen-Shannon entropy. If the sequence S is decomposed into segments S_i , i long and with base frequency p_{ij} , then the S -entropy per site is defined as $H(S_i) = -\sum_j p_{ij} \ln_2 p_{ij}$, and the total Jensen-Shannon entropy is

$$J(S, S_i) = \sum_i \lambda_i H(S_i) \{H(S)\} = H(S, S_i) \geq 0,$$

where

$$\lambda_i = \sum_{j=1}^n \lambda_{ij}.$$

For each decomposition of an intelligently refined set of likely decompositions, we find the probability P that a random decomposition of S has a value equal to or less than $J(S, S_i)$, and take the maximizing decomposition (S_i) as fixed. P is the optimal decomposition at confidence level P . Other statistical techniques have been suggested (Rauschky et al., 2000).

If this procedure is believed biologically, it must be repeated and taken into account. However, are there correlations within such long stretches, between successive ones, and how do we measure them? These are questions concerning the long-range properties of DNA.

3.2.2. Longest Repeat

Another general criterion of obvious biological relevance – typical of a number of tests used by Karlin and collaborators (see, e.g., Karlin et al., 1998) but not directly expressible in terms of correlations, is as follows. For a very long sequence, any subsequence may be expected to repeat, e.g., a seven-base

subsequence every $\ell' \sim 15,000$ bp, but longer ones may perhaps not repeat at all. A repeat of a longer subsequence might indicate biological function, evolutionary history, etc., and deserve significance. The most similar way of checking would be to look at the longest subsequence that repeats and test whether this is an unusual consequence of the length of the full sequence. For that purpose, a good thing to do would be to measure the mean and the variance of the longest repeat lengths—e.g., for simplicity over the assumption of independently distributed bases, and compare with observations. Therefore let us consider a chain of length n in base pairs. To search for subsequence repeats computationally, we may first count low-order repeats, and successively refine them by raising the iterated-order. We construct the self-comparison “dot matrix”

$$d_{\alpha\beta\gamma} = \begin{cases} 1, & \text{if } \alpha_1 \gamma_1 \dots \alpha_r \gamma_r = \beta_1 \dots \beta_r \\ 0, & \text{otherwise} \end{cases},$$

where base type is denoted by α_i , base location by $\alpha_1, \alpha_2, \dots$.

A	T	T	G	A	T
A					
T					
T					
G					
A					
T					

here, the by the matrix are represented by dots, 0s by no entry. A run of r dots down some upper or lower diagonal (i.e., $\alpha_1 \cdots \alpha_r$) then signifies a repeat of an ℓ -subsequence (All in the example shown). This matrix is symmetric. If we run half the diagonals — say the upper ones — one after the other, creating a string of size $= \frac{1}{2}n(n+1)$, we will include only a replicable fraction r/n of sites within r of the end of a diagonal that are recorded as possible members of r runs, but should not be. Also, although this is an matrix, represents only a piece of information, the probability of finding a repeat of the longest repeating subsequence is very small to start with, so we do not have to worry about correlations between run locations.

Hence the probability that the longest repeat is of length r is the same as the probability $P(r = \text{max})$ that the longest run of “successes,” or dots in a chain of length $t = \lfloor \frac{r}{2} \rfloor n^2$, is r . A success means the matching of two bases and hence occurs with probability p .

$$P = \sum_{r=1}^t p_r^2.$$

Now $P(r = \text{max})$ is born simply the probability that there is no $r + 1$ run minus the probability that there is no r run:

$$P(Q = \text{max}) = P(r + 1) - P(r).$$

Over a sequence of length t ($\lfloor \frac{r}{2} \rfloor n^2$). (Technically, $P(r = \text{max}) = P(r \text{ run and no } r + 1 \text{ run}) = 1 - P(\text{no } r \text{ run, or no } r + 1 \text{ run}) = 1 - P(r \text{ run}) - P(r + 1 \text{ run}) = P(\text{no } r \text{ run, but no } r + 1 \text{ run}) = 1 - P(r) - P(Q + 1) = 0$.) Let $q = 1 - p$. Then, for $P(r)$, we have typically x success runs of length $\leq r - 1$, separated by $y = x - 1$ failures, where $t = x + \sum_{j=1}^x d_j$:

$$\begin{array}{c} d_1 \quad d_2 \quad \dots \quad d_x \\ \hline \dots \quad 0 \quad 0 \quad \dots \quad 0 \end{array}$$

a probability of $q^{x-1} p^{\sum d_j}$. To sum up these probabilities at given t and r , we construct the generating function

$$\begin{aligned} \sum P(r) z^r &= \sum_{r=t}^{\infty} \sum_{\substack{x, y \geq 0 \\ x+y=r-1}} q^{x-1} p^{\sum d_j} q^y p^{t-x-y} \\ &= \sum_{x,y} q^{x-1} \left[\sum_{j=1}^{x-1} (pz)^j \right]^y \\ &= \sum_{x,y} q^{x-1} \left[\frac{1-(pz)^x}{1-pz} \right]^y \\ &= \left[\frac{1-(qz)^x}{1-pz} \right] \int \left\{ 1 - \frac{qz(1-(qz)^x)}{1-pz} \right\}^y \\ &= [1-(qz)^x] y! = z + g(qz)^x. \end{aligned}$$

This is a ratio of polynomials, $g(z) = z/(1-z)$ (the factor $1 - pz$ cancels out and the subscript indicates order); and hence can be partial fractioned as

$$\sum_{k=1}^r a_k P_k = b_k,$$

where $\alpha_n(\mu)/\mu$ satisfies $1 - \alpha_n + \alpha_n^2/\mu\alpha_n^2 = 0$ and

$$\begin{aligned} \alpha_n &= \lim_{t \rightarrow 0} Q_{n-1}(t)(1 - \alpha_n)/Q_0(t) = Q_{n-1}(1)/Q_0(1) \\ &= \frac{1 - q\mu\alpha_n^2}{-1 + q\mu(\mu\alpha_n^2)} = \frac{\alpha_n - \frac{\mu\alpha_n^2}{2}}{-\alpha_n + \mu(\alpha_n - 1)}. \end{aligned}$$

Now $\sum_{n=1}^{\infty} \alpha_n(\mu) = \alpha_1 = \sum_{n=1}^{\infty} \left(\sum_{k=1}^n (-\alpha_k(\mu)) t^k/k! \right)$, and we conclude that

$$P(x) = \sum_{n=1}^{\infty} \frac{\alpha_n}{x^n}.$$

For very large x , only the α_n , called x^* , of smallest absolute value will contribute, so

$$P(x) = \frac{x^{*} - \frac{\mu\alpha_1^2}{2}}{x^{*} - \mu(x^{*} - 1)} \frac{1}{x^{*+1}}.$$

For x^* , assuming sufficiently large x , we simply iterate $x^* = 1 + qx^2/(qx^2 - 1) + \dots$. Hence, to leading order, $P(x) = \exp(-\lambda qx^2)$, or

$$\begin{aligned} P(x) &= \exp(-\lambda qx^2) = e^{-\lambda qx^2} = e^{-\lambda q^2} \\ &= e^{-\lambda q^2} (e^{\lambda q^2 x^2} - 1) \\ &= \exp\left[-\lambda q^2 x^2 + \ln(e^{\lambda q^2 x^2} - 1)\right]. \end{aligned}$$

This has a sharp maximum w.r.t. x^2 , located at

$$\begin{aligned} -\lambda q^2 \ln p(\mu^2) + \lambda q^2 (\ln p(\mu^2) q^2 x^2)^{1/2}/(q^2 x^2 - 1) \\ = \lambda q^2 (\ln p(\mu^2) (-1 + q)) - e^{-\lambda q^2 x^2}] = 0, \end{aligned}$$

or $(1/\mu)^2 = \lambda q^2/(\ln 1/\mu)$. The derivative vanishes at the maximizing value,

$$F := \ln(p(\mu^2)/\ln 1/\mu)/(\ln 1/\mu),$$

and for the second derivative we have

$$\begin{aligned} \frac{\partial^2}{\partial x^2} [\ln(p(\mu^2))(-1 + q/(1 - e^{-\lambda q^2 x^2}))] \\ = \lambda q^2 \ln p(\mu^2) \left[-\ln p(\mu^2) + q \ln p(\mu^2) (1 - e^{-\lambda q^2 x^2}) \right. \\ \left. - q^2 \ln p(\mu^2) e^{-\lambda q^2 x^2} / \{ (1 - e^{-\lambda q^2 x^2})^2 x^2 \} \right] \\ = -\ell^2 q^2 \ln p(\mu^2) q^2 e^{-\lambda q^2 x^2} / q^2 = -\ell^2 q^2 \ln p(\mu^2) q^{2+\ell} \\ = -p(\mu^2) (\ln p(\mu^2))^2. \end{aligned}$$

Replacing the exponent of $P(r = \text{max})$ yields

$$P(r = \text{max}) = C \exp\left(-\frac{1}{2} \frac{R}{q^2} \left(\ln \frac{R}{p}\right)^2\right) r = R^2,$$

where C is a normalization constant, and we end up with

$$\begin{aligned} E(P_{\text{max}}) &= 2 \frac{\ln n}{\ln(1/p)} + \frac{\ln(1/p)^2 \ln(1/p)}{\ln(1/p)} + \dots, \\ \sigma^2(P_{\text{max}}) &= 2p^2/(n^2 \ln(1/p)^2). \end{aligned}$$

See also Karlin and Ost (1982) and Miel et al. (1998). A brief derivation in a more general context will be given in Section 4.1.8.

As an immediate application, we can look at the first DNA column as a mixture of SV40 and λ phage:

	Observed	Expected	χ^2
SV40	32	32	0
λ phage	31	31	0

The χ^2 -test is obviously significant.

Assignment 4

- Suppose we collapse the DNA information into patterns (0) and (1) (i.e., binary) (1). Show that the χ^2 -test for independence has $r - 1$ degrees of freedom, motivated by the fact that all four of the observed ($20, 20, 20, 20$) are the same.
- How would you generalize the long-range result to account with the order-1 Markov chain as pattern defined?

3.2.2. Displaced Correlations

Let us examine the long-range structure more systematically. To get a feeling for the quantities of interest, suppose that the we are not far from the independent placement of bases (with the base type now denoted by $i, i+1, \dots, j$). The most general question we might ask would be about the nature of the A_i subsignature distribution, which, in view of the approximate independence,

we write as

$$\begin{aligned} E[\delta_{i+1}(t_1)\delta_{i+2}(t_2) \cdots \delta_{i+j}(t_j)] &= \lambda p_i \prod_{k=1}^{j-1} [p_k + (\mu_{i+k}(t_k) - p_k)] \\ &= \left(\prod_{k=1}^j p_k \right) \lambda p_i \left[1 + \sum_{k=1}^{j-1} \frac{\delta_{i+k}(t_k) - p_k}{p_k} \right. \\ &\quad \left. + \sum_{1 \leq i < j} \frac{(\mu_{i+k}(t_k) - p_k)(\mu_{i+l}(t_l) - p_l)}{p_k p_l} + \dots \right] \\ &= \left(\prod_{k=1}^j p_k \right) \left[1 + \sum_{1 \leq i < j} \frac{p_{i+1}(1-p_i) - p_k p_l}{p_k p_l} + \dots \right]. \end{aligned}$$

Here we have taken advantage of the definition $p_k = \lambda p_i \delta_{i+k}(t_k)$, as well as of the translation invariance – *i.e.* within and *between* – of the averaging, which allows us to write

$$\lambda p_i \delta_{i+1}(t_1) \delta_{i+2}(t_2) = p_i \delta_i(t_1) \quad \text{for } i < j.$$

Thus it is sufficient at this level of analysis to know the covariances,

$$\begin{aligned} C_{ii'}(t) &= \text{Cov}(\delta_i(t), \delta_{i'}(t)) \\ &= p_i \delta_i(t) - p_i p_{i'}. \end{aligned}$$

the mean product of base type i , and type i' , fluctuating when t also spans. In fact, we now phrase a set of significant questions directly in terms of the Covariance Function approximations.

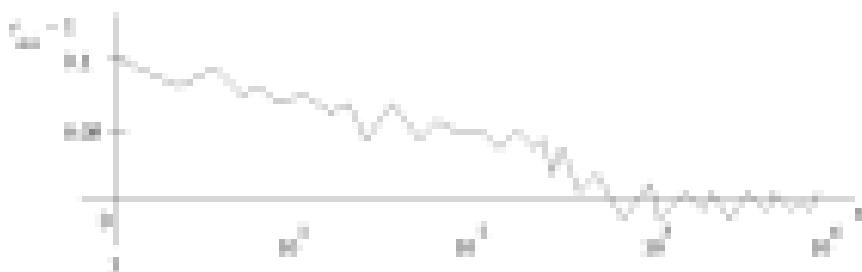
3.2.3. Nucleotide-level details

Imagine then a long fragment of a single strand of DNA. We expect that for the whole fragment, $\mu_A/\mu_C \cdot \mu_G/\mu_T$ will be different, although it is true that if we average over fragments of two-stranded complementary DNA, we should find that $\mu_A = \mu_T$ and $\mu_G = \mu_C$. In fact, these equalities tend to hold even for a long-enough piece of a single strand, the so-called pseudo-symmetry condition, *e.g.*, consider (Zheng and Marx, 2004a) the completely known 210-kb single strand of the yeast *Saccharomyces cerevisiae* Chromosome III (pseudo-random access), in which we find $\mu_A = 0.31$, $\mu_T = 0.30$, $\mu_G = 0.35$, and $\mu_C = 0.20$. We might sacrifice this symmetry, for example, to anomalous repressive events that result in frequent inversions of segment inversions.

In general detail, we want to study base correlations, the correlations $C_{ij}(R)$, or the corresponding correlation coefficients

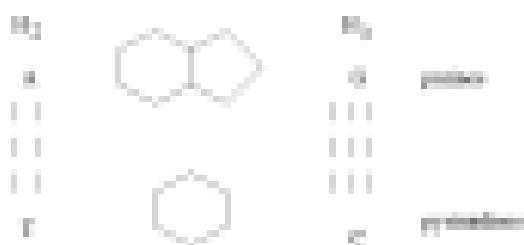
$$r_{ij}(R) = \rho_{ij}(R)/\rho_{ii}\rho_{jj}$$

The first question is that of persistence, the way that $r_{ij}(R)$ decays on the



average to its random value of 0. All four bases in this dimension behave similarly in this respect, with a range of ~ 140 nm (base-pairs). r_{AA} and r_{CC} are very similar, whereas r_{GG} is similar to r_{TT} , with higher amplitude fluctuation. It is not possible at this level to distinguish among detailed decay forms for $r_{ij} = 1$, such as e^{-R} and e^{-R^2} .

Going over displaced-base correlations $r_{ij}(R)$, it is useful to distinguish



on the basis of pairing type, two or three hydrogen bonds, and heterocyclic type, purine or pyrimidine. A first result is that $r_{AA(AA)}$ and $r_{TT(TT)}$ have negative correlations ($\rho = 1 \approx 0$) for $R \geq 6$, with a range of ~ 130 nm. For $r_{GG(GG)}$ and $r_{CC(CC)}$, with pairs not complementary in the sense of pairing type either (i.e., A and C, G and T), we start with smaller negative correlations, which then reverse. Finally, for a complementary pair, there is small negative correlation, up to ~ 150 nm. Although the wide fluctuations do not let us say much more, there is no significant regularity in the data on the average.

$$r_{ij}(R) = r_{ji}(R)$$

(where the asterisk indicates complementary base), a much more detailed kind of symmetry that would be exactly true for a pair of complementary strands.

We obtain some reduction in fluctuations by further amalgamating the bases. One such amalgamation figure of merit is the mutual information function of Li (1994),

$$M(x) = \sum_{i,j} p_{ij}(x) \ln [(p_{ij}(x)/p_i p_j)]$$

which vanishes if and only if $p_{ij}(x)$ factorizes as $p_i p_j$. In the human sequence cluster (H3k36M) VII of 12360 at (-0.7% intra, 0.11% inter), this gives an estimated correlation range of ~ 900 nt, consistent with the general



observation that introns tend to have long correlations. Compared with the uncorrected sequence of the same base ratio, but randomly placed bases, the result is quite sharp. Another statistical quantity used (West, 1993) that is symmetric over the bases

$$C(x) = \text{tr } C_0(x) = \sum_{i,j} [p_{ij}(x) - p_i^2]$$

has quite similar behavior, as does (Li et al., 1994)

$$C(x) = \sum_{i,j} C_{ij}(x) \delta_{i,j}$$

where the vector a_i is one of the four vertices of a tetrahedron at distance 1 from the origin (clearly $a_i \cdot a_j = 1$, $a_i \cdot a_k = -\frac{1}{2}$ for $i \neq j$).

Expansion-Mutation Model. There is no particular reason for total symmetry, and so we may choose to only distinguish between dual pairs,

$$C_0(x) = \frac{1}{2} [C_{AA}(x) + C_{CC}(x)] \neq 0$$

where $C_{AA+CC} = C_{AA} + C_{CC} + C_{AC} + C_{CA}$, or tetrahedral type

$$C_1(x) = \frac{1}{2} [C_{AA+CC}(x) + C_{AC+CA}(x)]$$

Obviously, because $\lambda_1(A) + \lambda_2(A) = 1 - \lambda_1(B) = \lambda_1(T)$ at each site, it is clear that $C_{\text{cov},\text{corr}}(k) = C_{\text{cov},\text{corr}}(n)$. It is similarly clear that, in terms of the variances, e.g.,

$$V_{1,2,\dots,k}(n) = C_{1,2,\dots,k,1,k+1,\dots,n}(n),$$

we have

$$C_k(n) = \frac{1}{k} V_{k+1,\dots,n}(n) = T_{k+1}(n).$$

In other words, we are reduced to a sequence of two symbols, say $B_i(t) = \lambda_1(G) = 0$ or 1, and ask for the process of extended persistence. A number of models have been suggested, in particular for what is claimed (Peng et al., 1992) to be a fractal or power-law decay $n^{-\alpha}$ dominated by noncoding DNA. One of these (Li and Kaneko, 1992) is evolutionary expansion with error. In its most primitive version, this is a simple cellular automaton (Lindenmayer, 1968) model with stochastic evolution. For a string of 0s and 1s at some point in time, a normal step is a doubling of an element, 0 \rightarrow 00 or 1 \rightarrow 11, but a substitution of the form (p) this is replaced with an error, 0 \rightarrow 1 or 1 \rightarrow 0. Thus there is a linear expansion rate per site of $\delta = 2(1 - p) + p = 2 - p$ at each time step. A heuristic consequence is this. Let $p_{\text{err}}^{(N)}(n)$ be the probability that symbol 0 is repeated by n after N time steps. Then assume uncorrelated law from 0 and 1 are 0 and 1)

$$p_{\text{err}}^{(N+1)}(n) = \sum_{n' \neq n} T_{n,n',n'} p_{\text{err}}^{(N)}(n').$$

For $N = \infty$, we will have $p_{\text{err}}(p) = \sum_n T_{n,n,n} p_n$, so that the covariance $C_{\text{err}}^{(N)}(n) = p_{\text{err}}^{(N)}(n) - p_{\text{err}}(p)$ satisfies the same equation. However, the transmission of n is still to be dominated by the same expansion rate at the spectrum $n \sim \delta n^{\alpha}$, thus we can write instead (assuming no further n dependence)

$$C_{\text{err}}^{(N+1)}(n) = \sum_{n' \neq n} T_{n,n',n'} C_{\text{err}}^{(N)}(n').$$

If $\delta > 1$ is the maximum eigenvalue of the 2×2 matrix T and π is the corresponding eigenvector, then asymptotically we must have

$$C^{(N+1)}(n) = 2 C^{(N)}(n/\delta),$$

where C is the π component of C . This has the obvious stationary solution

$$C(n) \propto n^{(\ln \delta)/2},$$

the desired power law. This may indeed be valid in prebiotic expansion and in larger entities that are repeated, with mutations.

Simple Sequence Repeats. In the preceding subsection, the symbols b_1 , b_2 , and b_3 might refer to small subsequences that are duplicated. Hence we might not be surprised to find many repeats, such as (GFT)ⁿ with n as large as 20, although the probable number in the two-stranded human genome under random equivalent placement of bases would be only $\sim 2 \times 10^{-1} \times 10^{29} = 2 \times 10^2$. Actually such microsatellites with $n > 10$ have $\sim 10^2$ bases in the human genome; they are highly polymorphic in length and provide a diagnostic signature of an individual much used for forensic purposes. What sort of length distribution would we expect?

Suppose [Perez, 1998] that a run of n_0 repeats of length k , with $n = n_0 k$, is required for viability and that this appears by mutation from a "twinkling" sequence. The latter must be correct at $n - 1$ positions and wrong at 1 of the n positions, and so will occur with probability $p_0 = (3/4)^{n_0}(1/4)^{k-1}$. In the resulting sequence (FTGF)ⁿ will be equivalent to (FTGF)ⁿ, and also (or the other strand) to (GFT)ⁿ or (GCF)ⁿ – an equivalence class of $N_{eq} = 4$ (whereas (FTGF)ⁿ would have $N_{eq} = 2$). Also, a non-optimum entry of the last base of a repeat would signal a different repeat and not be counted (probability 3/4 that the desired repeat is recognized), whereas a synonymous repeat of the repeat end (probability 1/4^k) would be rated as a new value of n . Hence the actual probability to be used is

$$p'_1 = \left(\frac{3}{4}\right)^{n_0} = N_{eq} \left(\frac{1}{4}\right)^{n_0-1} \left[1 - \left(\frac{1}{4}\right)^k\right].$$

For G base pairs, there are G possible starting positions for the sequence (ignoring end effects), so if β is the mean base substitution rate – hence β/k for a strand substitution – then we have a count of n_0 -fold repeats given by

$$\bar{n} = \frac{1}{3} G p'_1.$$

Thus, for (GFT)ⁿ, knowing that $G = 8 \times 10^{29}$ bp and $\beta = 5 \times 10^{-10}$ bp/year, we would have $\bar{n} = 0.001$ (quasistationary) for $n_0 = 2$, 0.002 for $n_0 = 3$, etc.

Given the source of n_0 -fold repeats, we need a model for the dynamics. With $n = n_0 + k - 1$, so that $k = 1$ initially, imagine a birth-death process of k Amt base per microsatellite, i.e., probability $P(k)$ that it has k $\rightarrow k + 1$ and also $k\downarrow$ for $k \rightarrow k - 1$. Thus, if $P(k, t, t')$ is the probability of having k repeats at time t , given that we started with $k = 1$, counting the disappearance from the k pool, the survival from $k - 1$, and the arrivals from $k + 1$, we clearly

have

$$\begin{aligned}\frac{d}{dt} P(1, k, t) = & -2\lambda k P(1, k, t) + \lambda(k-1)P(1, k-1, t) \\ & + \lambda k + 1)P(1, k+1, t)\end{aligned}$$

for $k = 1, 2, \dots$. On the other hand, if k decreases to 0, the repeat is not feasible and the process stops,

$$\frac{d}{dt} P(1, 0, t) = 2\lambda P(1, 1, t),$$

and of course we have the initial condition

$$P(1, k, 0) = A_{k,1},$$

We solve the difference equation in standard fashion by writing up the generating functions

$$P_{1,t}(t) = \sum_{k=0}^{\infty} P(1, k, t)z^k,$$

multiplying the difference equation by z^k , summing over $k = 1, 2, \dots$ and adding the $k = 0$ part gives us that

$$\frac{d}{dt} P(1, t) = \lambda z - 1)^2 P_{1,t}(t)$$

$$P(1, 0) = 1.$$

This is readily solved as

$$P(1, t) = 1 + \frac{z-1}{(1-\lambda z-1)^2},$$

yielding the distribution

$$P(1, k, t) = \frac{(z\lambda)^{k-1}}{(1+\lambda z)^{k+1}} \quad \text{for } k \geq 1$$

[so that $P(1, 0, t) = \lambda z / (1 + \lambda z)$. With a steady-state S of seeds from $t = 0$ to $t = T$, the total number of P is

$$P(1, k, T) = S \int_0^T P(1, k, t) dt = \frac{S}{\lambda z} \left(\frac{z\lambda^T - 1}{1 + \lambda z} \right),$$

the full distribution. Note that the total number of repeat sequences is

$$\begin{aligned} N(T) &= \sum_{k=1}^{\infty} N_k(k, T) \sim \frac{S}{k} \int_0^{\infty} \frac{1}{k} e^{-kx/T} dx \\ &= \frac{T}{k} E_1\left(\frac{1}{kT}\right). \end{aligned}$$

In terms of the exponential integral E_1 , and has the asymptotic $\ln(T)$ form

$$N(T) \sim \frac{S}{k} (\ln(AT) + \gamma)$$

($\gamma = 0.577\ldots$ is Euler's constant), a very slow increase even in the absence of mutational deterioration.

Length Distributions. Repeated not by a simple, well-behaved mechanism, and need not be confined to preceding regions. For example, it has been suggested numerous occasions that not only are genes just a shuffling of a small set number of motifs (Gilbert, 1967; Strelakova et al., 1994) (but see Durbin et al., 1998), but also (see, e.g., Dewey, 1998) that exons are composed of relatively small numbers of ancestral units, of course evolved to some extent. One thing is certain: The three-letter codons are highly repeated, but are not equally likely, and this can very much affect the correlation structure within an exon. That is, if codon (r, s, t) occurs at relative frequency f_{rst} , then, neglecting any correlations between codons in a very long exon, we will have (Bloodgood & Chessa, 1997)

$$f_m^R(M) = \frac{1}{3} \sum_{\substack{(r,s,t) \\ \text{fixed}}} (f_{rrr} f_{rst} + f_{sss} f_{rst} + f_{ttt} f_{rst})$$

independently of b , and

$$f_m^R(M+1) = \frac{1}{3} \sum_{\substack{(r,s,t) \\ \text{fixed}}} (f_{rrr} f_{rst} + f_{sss} f_{rst} + f_{ttt} f_{rst}) = f_m^R(M-1).$$

Thus a symmetric probe such as the mutual information function $M(r)$ will simply have a fine structure with period 3.

For a distribution of finite mean lengths $\mu(p)$, mutations are more interesting. We first note that for weak dependence, i.e., small $C_{rs}(0) = p_{rs}(0) - p_r p_s$ the mutual information can be expanded as

$$M(r) = \frac{1}{2} \sum_{s,t} C_{rs}(0)^2 / p_r p_s.$$

Now suppose, as in prokaryotes, we have almost all exons. Then the observed correlations are still for a pair within an exon, but there are only $T - n$ possible

placements, compatible with $L = n = L$, for the full-length L sequence, so that

$$C_{1,0}(x) = \sum_{l=1}^{L-n} p(l|x_l - x)/L \cdot C_{1,0}^{(0)}(x).$$

The factor $(nL) = \sum_l p(l|x_l - x)/L$ produces an exponential distribution if $p(l)$ is exponential and a power law $\propto l^{-1-\beta}$ if $p(l) \propto l^{-\beta}$ above some cutoff, and so in fact the long-range correlation represented by $p(l)$ is a direct reflection of the main length distribution.

The length distribution of OBDP's (open reading frames = DNA stretch between a start-codon and a stop-codon in the same frame = coding sequence) in a number of organisms has been studied (Li, 1995), with the conclusion that these are always exponential, although the exponent may change sharply at a few $>10^3$ bp. Thus the main contribution to asymptotic correlations will only be exponential. On the other hand, it was noted (Almirante and Pirota, 1995) that there is empirical evidence that the length distribution of pulse and pyramidal clusters has the large/l form

$$p(l) \propto l^{-1-\alpha},$$

which suggests that the distribution be moderately stable distribution (Feller, 1968). If this is the case, then a noncoding region, modeled as a concatenation of P_a and P_b clusters, will have the same asymptotic length distribution, with the same consequences for rate of losses.

3.2.4. Windowed Correlations

To improve the statistics, we should use bigger batches of data. Most obvious is to accumulate information in windows of bases, $W(x)$, N bases starting at x , and average over x (Heckert and Timp, 1992). The windows, as expanded points, should be made much shorter than the fragment being measured and will then give meaningful results and statistics. We distinguish between correlations within a window and between windows.

Within-Window Correlations

For a window of size N , we tally n_x , the number of occurrences of x in a given window, and construct

$$\text{Cov}^{(N)}(y_x, R) = \langle y_x(R) - \langle y_x \rangle(R) \rangle \langle R \rangle.$$

In practice, the n_x are usually chosen so that the windows do not overlap, thereby avoiding next correlations between windows, but it turns out

surprisingly that none of the significant results alter if we let α run over all sites, and we shall do so. Of course, $(\theta_0) = P(\theta_0)$, but also, because $\theta_0 = \sum_{i=1}^{N-1} \theta_{i+1}$, we have $\text{Cov}^{(2)}(\theta_0, \theta_0) = \sum_{i=1}^{N-1} \sum_{j=i+1}^N \text{Cov}[\theta_{i+1}(t), \theta_{j+1}(t)]$, i.e., gathering together common values of $|i - j| = \alpha$,

$$\text{Cov}^{(2)}(\theta_0, \theta_0) = N(p_1 \theta_1 - p_2 \theta_2) + \sum_{\alpha=1}^{N-1} (N - \alpha) [C_\alpha(\alpha) + C_{\alpha}(\alpha)],$$

and, as an important special case, the variance

$$\text{Var}^{(2)}(\theta_0) = \text{Cov}^{(2)}(\theta_0, \theta_0) = Np_2(1 - p_2) + 2 \sum_{\alpha=1}^{N-1} (N - \alpha) C_\alpha(\alpha).$$

A convenient measure is the correlation coefficient of standardised maps (not to be confused with the correlation coefficient $r_{ij}(x)$ of Definition 2.2.2),

$$\text{corr}^{(2)}(\theta_0, \theta_0) = \text{Cov}^{(2)}(\theta_0, \theta_0) / [\text{Var}^{(2)}(\theta_0) \text{Var}^{(2)}(\theta_0)]^{1/2},$$

where it is readily verified that $-1 \leq \text{corr}^{(2)} \leq 1$. Note that if the base-phases were independent (stochastically), so that $C_\alpha(\alpha) = 0$, we would have

$$\text{Var}^{(2)}(\theta_0) = Np_2(1 - p_2),$$

$$\text{corr}^{(2)}(\theta_0, \theta_0) = -(p_1/(1 - p_1))^{1/2} (p_2/(1 - p_2))^{1/2} \quad \text{for } x \neq 1,$$

the latter being $-1/2$ when all $p_i = 1/2$.

Let us look a bit at $\text{Var}^{(2)}(\theta_0)$ [similar considerations apply to $\text{Cov}^{(2)}(\theta_0, \theta_0)$]. If the correlation is short range, e.g., $C_\alpha(\alpha) \sim \lambda e^{-\alpha \beta}$, as in Motter's order 1, then $\sum (N - \alpha) C_\alpha(\alpha) \sim \int_0^\infty (N - \alpha) \lambda e^{-\alpha \beta} d\alpha = \lambda e^{-\beta} (\lambda/\beta \alpha) \int_0^\infty e^{\alpha \beta} d\alpha = (\lambda/\beta)(\lambda - \lambda)(1 - e^{-\beta})/\beta^2$, so that asymptotically

$$\text{Var}^{(2)}(\theta_0) = N \left[p_2(1 - p_2) + 2 \frac{\lambda}{\beta} \right] + \dots$$

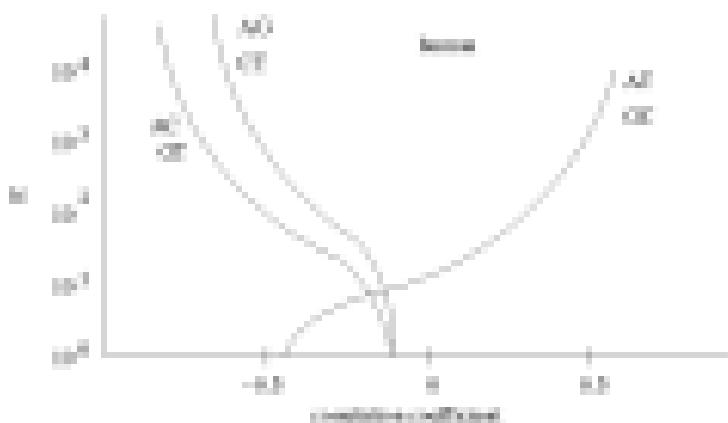
and similarly N . On the other hand, if the decay is strong, e.g., the power law $C_\alpha(\alpha) \sim \lambda \alpha^{-\gamma}$, where $0 < \gamma < 1$, then $\int_0^\infty (N - \alpha) \lambda \alpha^{-\gamma} d\alpha = \lambda/(1 - \gamma) (\lambda - \alpha)^{1-\gamma}$, and

$$\text{Var}^{(2)}(\theta_0) = \frac{2\lambda}{(1 - \gamma)(2 - \gamma)} (N^{1-\gamma} + p_2(1 - p_2) N + \dots)$$

is dominated by a power higher than N^1 .

What happens in actual genomes? Typically the extremes of bases (mainly with the rRNA, which is complementary to the processed RNA that codes for proteins, excluded as being specifically *E. coli* bacteria are studied). The first observation is that variances are much larger than theory. Secondly,

an even very reasonable (e.g., 10) Markov order. For example, at $N = 1000$, we have $\text{Var} \approx 18,000$ for the full human contained $A + T$ content, 2500 for



corresponding C rate, in a 1-kilobase pair window, compared with -0.02 for independence. There are also dramatic effects in correlations (Pelizzetti and Tang, 1992) that start at ~ -0.15 , as expected for small windows, but quickly depart to extreme values, and by 1.2 kb have arrived at:

$$\text{corr}(A, C) = -0.79, \quad \text{corr}(A, G) = -0.65, \quad \text{corr}(A, T) = 0.42, \\ \text{corr}(G, T) = -0.77, \quad \text{corr}(G, C) = -0.67, \quad \text{corr}(C, T) = 0.45.$$

The AT and GC correlations quickly become larger and positive, perhaps approaching the $\text{corr}(A, T) = 1 = \text{corr}(G, C)$ for full strand symmetry for which $A0 = GT$ and $G0 = GC$. This provides focusing on $A + T$ or $G + C$ content, leaving out fluctuations, without losing information.

There is as well the qualitative fact that closely

$$\text{corr}(B_0, B_0') = \text{corr}(B_0, B_0''),$$

where B, B' , and B'' cover all four base types. In any order – an equality that holds to $\sim 1\%$ when we look at the corresponding covariances. This is actually not too informative. At fixed N , $(B_0 + B_0') \approx -(B_0 + B_0'') + B_0'''$, so $\text{Var}(B_0 + B_0') = \text{Var}(B_0) + \text{Var}(B_0'') + 2 \text{Cov}(B_0, B_0'') + \text{Var}(B_0''') = \text{Var}(B_0) + 2 \text{Cov}(B_0, B_0'') + \text{Var}(B_0''')$, equivalent to

$$\text{Var}(B_0) = \text{Var}(B_0') = \text{Var}(B_0'') = \text{Var}(B_0'''),$$

which is hardly surprising.

Between-Window Correlations

Here the persistence effect noted in the previous subsection again appears. Define

$$\text{Corr}(\mathbf{x}(t), \mathbf{x}(t')) = \text{Cor}^{(B)}(\mathbf{x}(t) + \mathbf{B}, \mathbf{x}(t') + \mathbf{B}) / \sqrt{\text{Var}^{(B)}(\mathbf{x}(t) + \mathbf{B})}$$

at fixed window size N , say 1 kb. For two nucleotides of two windows at distance d apart. Then both human and E. coli DNA have very large ranges indeed. Note that, quite generally,



$$\begin{aligned}\text{Cor}^{(B)}(\mathbf{x}(t), \mathbf{x}(t')) &= \sum_{\alpha=0}^{N-1} \text{Cor}(\mathbf{x}_{t+\alpha}(t), \mathbf{x}_{t+\alpha}(t')) \\ &= \sum_{\alpha=0}^N (\alpha + N) C_{\alpha}(t + \alpha) + \sum_{\alpha=1}^N (\alpha - N) C_{\alpha}(t + \alpha) \\ &= B^2 C_0(d),\end{aligned}$$

where $C_0(d)$ is a triangularly weighted average of correlations at separations from $d - N$ to $d + N$ (of course, $d > N$).

Fluctual Moments

Higher within-window moments of \mathbf{x}_t , or of combined bases, e.g., $\mathbf{x}_t + \mathbf{x}_{t+1}$, have been considered as well (Oliverity and Nanyena Rao, 2000), and in particular the normalized factorial moments

$$F_q^{(B)}(q) = (\langle \mathbf{x}_t \rangle^q / \langle \mathbf{x}_t \rangle - q!) / \langle \mathbf{x}_t \rangle^q$$

have some desirable properties. For example, if \mathbf{x}_t is Poisson distributed, $\text{Prob}(\mathbf{x}_t = \lambda)^B / \langle \mathbf{x}_t \rangle^B \propto e^{-\lambda}$, then all $F_q^{(B)}(q) = 1$. For a number of real DNA sequences, both with and without biases, there appears to be a monotonic trend as a function of the window size N : for $N \approx N_c = 10^4 - 10^5$ bp, all $F_q^{(B)} < 1$, variance increases as N and fluctuations appear Gaussian; for $N > N_c$, $F_q^{(B)} > 1$, variance increases as $N^{1.5-1.7}$ and fluctuations appear to be non-Gaussian. The Gaussian property can be assessed (Villegas

et al., 1998) by means of the function

$$\eta(N) = 1 - 2 \frac{\langle D(N)^2 \rangle}{\langle D(N) \rangle^2},$$

where $D(N)$ is the window- N sum $D(N) = \sum_{i=1}^{N-1} s(i)s(i+N)$ ($s(i)$ is the base at i ; a pyrimidine is parity). For a Gaussian distribution, $\eta(N) = 0$, whereas a model of correlated base sequences embedded in a sea of uncorrelated bases reproduces observations of the preceding type quite well.

3.2.5. Statistical Models

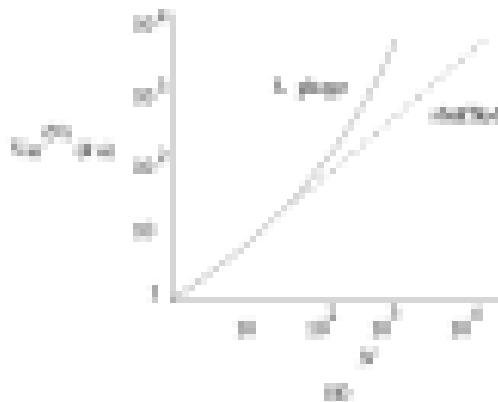
Percol Model

We examine the question of the genesis of long-range correlations, measured again by means of the window- N quantities

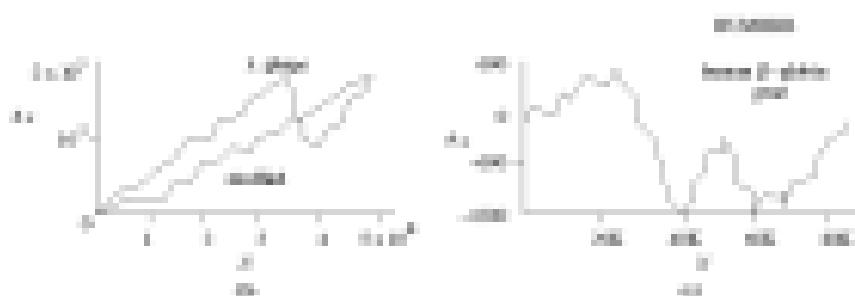
$$\delta C := BC + BT - BA - BG,$$

whose parity value is $n = +1$ for a pyrimidine and $n = -1$ for a purine (Karin and Bresler, 1993). As we have seen, $\langle \delta C \rangle_{\text{stat}} = 4 \langle \delta C \rangle_{\text{per}}$, and so $\langle \delta C^{(2)} \rangle_{\text{per}} = 4 \langle \delta C^{(2)} \rangle_{\text{stat}} = 4C$ as well, but let us switch the variable the instant, which, as N increases, creates a simple random walk, $\Delta(\mu)$ (μ is odd at each step), on the integers. We have, in the usual way (but noting that here $\pi^2 = 1$),

$$\langle \delta C^{(2)} \rangle_{\text{per}} = [N+1] \sum_{k=1}^{N-1} (N-k) C_{kk}(N).$$



A plot of the hexatriplet λ sequence shows an increase of lower than N^2 compared with the shuffled sequence, which indeed goes as N . In fact, a



simple count of the $\langle \lambda \rangle$ at a fixed starting position shows an expected linear increase for shuffled λ phage compared with piecewise linear sections for λ itself and even more pieces for a section of human genome, looking as if we have patches of different nucleic acid (or CT) composition, i.e., nucleotide heterogeneity. We can model this situation most simply by supposing two types of patch or subsection, one with $\langle x \rangle = u^2$, the other with $\langle x \rangle = v^2$. The probabilities that $x = +1$ (CT) and that $x = -1$ (GATC) are, respectively, $\frac{1}{2}(1 + u^2)$ and $\frac{1}{2}(1 - u^2)$ in these two regions. Now suppose that there are no correlations within each region, i.e., Bernoulli or Markov order 0. Then the joint probabilities for $x = +1$ or -1 at sites i and $i+1$, conditional on O and x being in u^2 regions (the first and the second arguments u^2 are distinct) are given by

$$\begin{bmatrix} \Pr(O | x^k, u^2) & \Pr(O | x^{k+1}, u^2) \\ \Pr(O | x^k, v^2) & \Pr(O | x^{k+1}, v^2) \end{bmatrix} = \begin{bmatrix} \log_2 \log_2 & \log_2 \log_2 \\ \log_2 \log_2 & \log_2 \log_2 \end{bmatrix}$$

What we need are the probabilities of O and x being in u^2 , v^2 regions, which depends on how many uniform intervals separate O and x .

There are many models for such interval distributions, but the simplest is to suppose that the switches from region to region, or patch to patch, occur independently and randomly at rate λ . Then $\Pr(O \text{ and } x \text{ both in same type region}) = \Pr(\text{zero switches}) = \Pr(\text{no switch}) + \Pr(\text{one switch}) + \dots + \infty(1 - \lambda)^2 + (\frac{\lambda}{2})(\lambda)^2(1 - \lambda)^{2-1} + \dots = \frac{1}{2}(1 - \lambda) + \lambda^2 + (1 - \lambda) - \lambda^2) = \frac{1}{2}(1 + 2\lambda)^2$, and similarly $\Pr(O \text{ and } x \text{ in different regions}) = \frac{1}{2}(1 - (1 - 2\lambda)^2)$. We conclude that

$$\begin{aligned} \pi_{uv}(x) &= \frac{1}{2}(1 + (1 - 2\lambda)^2) \left[\left(\frac{1 + u^2}{2} \right)^2 + \left(\frac{1 + v^2}{2} \right)^2 \right] \\ &\quad + \frac{1}{2}(1 - (1 - 2\lambda)^2) \left(\frac{1 + u^2}{2} \frac{1 + v^2}{2} \right) \\ &= \frac{1}{4} \left(1 + \frac{u^2 + v^2}{2} \right)^2 + \frac{1}{4} \left(\frac{u^2 - v^2}{2} \right)^2 (1 - 2\lambda)^2. \end{aligned}$$

so that $C_{11}(k) = p_1(k) - p_1(0)\delta_k = \frac{1}{2}(x^k - u^{-k})^2(1 - 2u)^2$. It follows that

$$\text{Var}^{(2)}(\theta_N) = N + \frac{1}{4}(x^N - u^{-N})^2 \sum_{k=1}^{[N/2]} (M - \delta_k)k = 2M^2.$$

In fact, for large patches, $\lambda \sim 0$, we have

$$\text{Var}^{(2)}(\theta_N) = N + \frac{1}{16}(x^N - u^{-N})^2/N \approx 1,$$

which, in turn, matches the observed N -dependence as well as the "fitted" $N + \mu N^{1-\alpha}$.

However, quantitative comparison is not reliable, and evidence has accumulated that the long-range power-law dependence of correlations is indeed the norm. In particular, Peng et al. (1994) (see also Bernaola-Galve et al., 1998), have examined in greater detail the extent to which a patch model – and thus certainly no patches – can mimic long-range power-law correlations. They did this by first dividing the sequence into subsequences of length $l = 10^2$, small enough with anticipated patches, finding an optimal linear fit to $\text{Var}(A)$ in each subregion, and subtracting it out (intercepting), thus eliminating any patch bias. The resulting $\theta_{\text{reg}}(N)$ hence fluctuates around zero, and it is to this that the window standard-variation analysis is applied. When this is done, the power-law violation is reversed as well as the threshold ℓ that measures the extent of patch contributions.

With confidence that the power-law variation is intrinsic, the analysis was formalized (Lo et al., 1996) in terms of the Hurst index for self-similar patterns. That is, if $E^{(m)}$ is a measured characteristic averaged over a block of length m , and the system is divided into nonoverlapping blocks instead by l , then the assertion that

$$\varepsilon^{(m)}(k) = [E^{(m)}_k - \bar{E}] [E^{(m)}_{(k+l)} - \bar{E}] = k^{-\beta}/\ell(k)$$

with slowly varying \bar{E} , independently of m (large m), defines a Hurst index, $H = 1 - \frac{1}{2}\beta$ for the self-similar sequence in question. Actually, a related assessment, that of the dimensionality of a walk constructed with A, T, G, and C as 2D steps on coordinate axes in the plane, had been carried out much earlier (Berthoumieu et al., 1992). The conclusion (arguably by use of the approximate equality of A and T steps and G and C steps) was that dimensionality could be defined and is anomalously small compared with that of a random sequence, consistent with power-law asymptotic correlations.

Hidden Markov Models

The Karlin–Stewart model is a special case of a class of stochastic models (Davidson, 1997; Baldi and Trifunovic, 1998; Dublin et al., 1998) that is being

and more and more to represent partially random data nonstationarity, e.g., in pattern recognition. We imagine a "hidden" set of main states $\{i_1 \in J\}$ and associated first-order Markov transition probabilities ($T_{ij} := P_{ij}(t \mid t')$). The process can be regarded either as starting with a state distribution P_1 , or,



equivalently, selecting a fixed starting state i_0 appended to the seq. S of possible states, with $T_{0i_1} = P_1 P_{0,i_1}$. Likewise, one may stop in "top band" at the n th stage or we can append an end state 0 and associated stopping probabilities T_{0i_n} , with the constraint $T_{01} = \delta_{0,1}$, thus producing a probabilistic stopping point.

To complete the description, there are the observed or "visible" states $\{y_i \in A\}$ that are independently chosen by means of the transition probabilities $\{Q_{ij} := P_j(y_i \mid i)\}$. Hence, at time n ,

$$P_v(i_1, \dots, i_n) := \sum_{y_1, \dots, y_n} Q_{i_1,y_1} P_{i_1,i_2} Q_{i_2,y_2} \dots \dots \dots T_{i_n,0} Q_{i_n,y_n} P_{0,y_1}$$

and we can now inquire about correlations such as

$$p_{xy}(k) := \sum_{i_1, \dots, i_k} P_v(i_1, \dots, i_k, y_{k+1}, \dots, y^k).$$

There are two extreme classes of transition matrices, the first being that of recurrent matrices, in which $T_{ii} > 0$ for all main states. For example, in the Karlin–Brenner model, with two hidden types,

$$T = \begin{bmatrix} 1-\lambda & \lambda \\ \lambda & 1-\lambda \end{bmatrix},$$

with start and stop included by extending this to

$$T = \begin{bmatrix} 0 & 0 & 0 & E \\ 0 & p(1-\lambda) & \lambda & 0 \\ 0 & \lambda & p(1-\lambda) & 0 \\ 0 & 1-p & 1-p & 0 \end{bmatrix}.$$

For the emission process, we would have

$$\hat{Q} = \begin{bmatrix} \{1 - w^2\} & \{1 - w^2\} \\ \{1 + w^2\} & \{1 + w^2\} \end{bmatrix}.$$

A second class is termed left-to-right. Here, we are looking for an underlying pattern, say x^1, x^2, \dots, x^n , and so we have main states x^1, x^2, \dots, x^n , and X (the anything), with transitions restricted to $X \rightarrow X$ or x^1 , but $x^1 \rightarrow x^2 \rightarrow \dots \rightarrow x^n \rightarrow X$. A typical problem associated with a such a hidden Markov model (HMM) is that of extracting the most likely hidden state sequence (x) given the output (y) , which we may accomplish, e.g., by first maximizing $\Pr(y_1, \dots, y_n | T, Q)$ to find optimal (T, Q) , and then maximizing $\Pr(y_1, \dots, y_n | x_1, \dots, x_n)$.

HMMs are nice examples of a broad class of models (3) in the figure below) in which coupled outputs are controlled by hidden coupled inputs.



In (b) in the figure the outputs are coupled only to the inputs, i.e., Bernoulli, and we have a highly reduced model only assignable at a sufficiently coarse level of resolution. The HMM, (c) in the figure, has further specification... At a first level, the outputs should presumably at least be coupled informationally to their neighbors, (d) in the figure, and this is the thought behind the next model we consider.

Markov Blanket Model

The hallmark of the Kerlin-Grenville model was a set of “hidden” instructions, namely a subdivision into distinct homogeneous regions that was itself statistically determined. We thereafter had a Bernoulli sampling with parameters dictated by an independent interval process, which is to say a special case of a Bernoulli process (strictly speaking, “process” infers a continuous duration of the transition event, valid here only if the base-to-successive-base interval is regarded as infinitesimal). To do a better job on short-range structure, we would probably want each homogeneous subregion to be at least a first-order Markov chain (“chain” means a discrete duration, and the interval between bases is certainly not small for low-order Markov). To do a better job on the long-range structure, we could replace the independent intervals or patches with a general first-order Markov process, making transitions between various

types of patterns (Klein, 1971) but, as indicated and more readily analyzed, we choose a model (Chandraratna, 1989; Hockett and Tang, 1992) in which a transition is attempted, with very low frequency, at every base, thereby succeeding on only a large scale. The model of Hockett and Tang that we now turn to also involves an $A + T$ assumption, consistent with the high degree of correlation between A and T .

Imagine that a first-order base-to-base Markov chain that depends on a hidden parameter w (or parameter w_0) that varies from step to step. The assumption is that w changes autoregressively by small increments, with transition probability $W(w' | w)$. Then at each site there is a combination of these a and parameter w —a Markov chain on a higher space—and we can proceed sequentially, starting, say with (r, w_0) , at site 0:

site 0: (r_1, w_0)

site 1: (r_1, w_1) with probability $P_{w_0}(w_1) W(w_1 | w_0)$

site 2: (r_2, w_2) with probability

$$\sum_{w'} P_{w_1}(w') \int W(w_1 | w') P_{w_0}(w) W(w_1 | w_0) dw_0$$

...

...

...

and at site M , the probability of (r, w) , given (r, w_0) initially, is

$$P_M^{(0)}(w | w_0) = \sum_{w'} P_{w_0}(w) \int W(w | w') P_{w_0}^{(M-1)}(w' | w_0) dw'$$

where $P_{w_0}^{(0)}(w | w_0) = \delta_{w_0, w}$.

Now we specialize to the two-state context. At time 0 , a dichotomous Markov transition matrix must have the form $\begin{bmatrix} 1-\alpha & \beta \\ \alpha & 1-\beta \end{bmatrix}$, so that:

$$\begin{bmatrix} P_{w_0,11}(w | w_0) \\ P_{w_0,12}(w | w_0) \end{bmatrix} = \begin{bmatrix} 1-\alpha(w) & \beta(w) \\ \alpha(w) & 1-\beta(w) \end{bmatrix} \int W(w | w') \begin{bmatrix} P_{w_0,11}^{(M-1)}(w' | w_0) \\ P_{w_0,12}^{(M-1)}(w' | w_0) \end{bmatrix} dw'.$$

and similarly for probabilities conditioned on G^C . If (r, w_0) initially has

probability $p(x_0) f(x_0)$. This becomes, on integration over x_0 ,

$$\begin{bmatrix} p_{\text{tot},\text{tot}}^{(0)}(w) \\ p_{\text{tot},\text{tot}}^{(0)}(w') \end{bmatrix} = \begin{bmatrix} 1 - \alpha(w) & \beta(w) \\ \alpha(w) & 1 - \beta(w) \end{bmatrix} \int W(w + w') \begin{bmatrix} p_{\text{tot},\text{tot}}^{(0)}(w') \\ p_{\text{tot},\text{tot}}^{(0)}(w') \end{bmatrix} dw'$$

where

$$\begin{bmatrix} p_{\text{tot},\text{tot}}^{(0)}(w) \\ p_{\text{tot},\text{tot}}^{(0)}(w') \end{bmatrix} = p(x_0) f(x_0) \begin{pmatrix} 1 \\ \alpha \end{pmatrix}.$$

Adding the two rows, we see that

$$p^{(0)}(w) = p_{\text{tot},\text{tot}}^{(0)}(w) + p_{\text{tot},\text{tot}}^{(0)}(w')$$

which is

$$p^{(0)}(w) = \int W(w + w') p^{(0)-1}(w') dw'.$$

$$p^{(0)}(w) = p(w) / f(w).$$

This allows us to eliminate $p_{\text{tot},\text{tot}}^{(0)}(w')$, obtaining, with $p_{\text{tot},\text{tot}}^{(0)}(w)$ abbreviated as $p^{(0)}(w)$,

$$p^{(0)}(w) = [1 - \alpha(w) - \beta(w)] \int W(w + w') p^{(0)-1}(w') dw' + \beta(w) p^{(0)}(w),$$

$$p^{(0)}(w) = p(w) / f(w).$$

It is convenient to choose the parameter α so the value of the probability p that the system is "missing" all $\lfloor \frac{w}{\gamma} \rfloor - \lfloor \frac{w'}{\gamma} \rfloor = \lfloor \frac{w-w'}{\gamma} \rfloor$ is $p = \beta(w + \beta) = \alpha$, and introduce $\beta' = w + \beta$, the total deviation from independence. Hence

$$\alpha(w) = (1 - \beta(w)) p(w), \quad \beta(w) = \log(p(w)).$$

Taking the initial $p(w) = \infty$, as well, we therefore have

$$p^{(0)}(w) = (1 - \beta(w)) \int W(w + w') p^{(0)-1}(w') dw' + \beta(w) p^{(0)}(w),$$

where

$$p^{(0)}(w) = \int W(w + w') p^{(0)-1}(w') dw'.$$

$$p^{(0)}(w) = p^{(0)}(w) = w f(w),$$

subsequent to which, of course,

$$\text{Per}_{\text{ext}}(B) := \int p^{(B)}(y) dy.$$

All that we need is the function $p(y)$. In practice, we can find this by taking k-th order windows from the sequence of interest, putting them in bins of size Δw , or $w = \frac{1}{2}\Delta w$ and fitting p to each bin. In the analysis of Pickard and Tang, it turns out that $0 < p \leq 1.1$, very close to independence. $\text{Per}_{\text{ext}} = \text{Per}_{\text{per}}$, and hence is the standard HMM. To the extent that $p = 1$ is valid, the distribution degenerates to

$$\text{Per}_{\text{ext}}(B) = \int w p^{(B)}(w) dw,$$

where

$$p^{(B)}(w) := \int W(w+w') p^{(B-w)}(w') dw',$$

$$p^{(B)}(w) = w p^B(w).$$

This motivates the question of the addition of w , as represented by $p^{(B)}(w)$. The iterated action of $W(w+w')$ on $p^{(B)}(w) = w p^B(w)$. Assuming that there is only a small change in w each time, we can set $p^{(B)}(w') = p^{(B)}(w) + (w' - w) p'(w)p^{(B)}(w) + (w' - w)^2 p''(w)p^{(B)}(w) + \dots$ in the transition equation. If the transitions are symmetric, $W(w+w') = W(w'+w)$, then $\int (w' - w) W(w+w') dw' = 0$; hence $\int W(w+w') dw' = \int W(w'+w) dw' = 1$, and we define $w^2(p) := \int (w' - w)^2 W(w+w') dw'$. The "dynamic" then takes on the Fokker-Planck form (Feller, 1950, Chap. XIV)

$$p^{(B)}(w) - p^{(B-w)}(w) = \frac{1}{2} w^2(p) \frac{\partial^2}{\partial w^2} p^{(B-w)}(w).$$

Finally, regarding the dependence on B as well-known, we can write this as

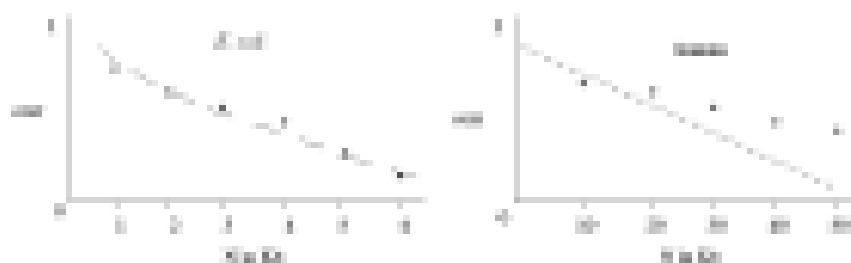
$$\frac{\partial}{\partial B} p(w, B) = \frac{1}{2} w^2(p) \frac{\partial^2}{\partial w^2} p(w, B),$$

$$p(w, 0) = w p(w).$$

a standard diffusion equation.

In Pickard and Tang (1992) the assumption is made that $w^2(p) = w^2$ is a constant on the allowed range of w and is obtained by fitting the large B

asymptotic solution to the bin distribution above. For the human genome, this corresponds to $\frac{1}{2} \leq n \leq \frac{3}{2}$, with $\alpha = 0.0115$, and for $K = 100$ to $0.4 \leq n \leq 0.6$, $\alpha = 0.0025$. Further, $f(n)$ was taken as uniform over the half n interval. The nondegenerate diffusion problem, with full $p_{ij}(n)$, was then solved by simulation, and the $\langle ET, AT \rangle W$ correlation coefficient computed. The results for $K = 100$ are very good, whereas that for human DNA, drops a bit too rapidly, perhaps because of the assumption of n -independent σ^2 or truncated uniform $f(n)$.



Assignment 3

1. What is the general class of base-symmetric figures of merit C_{sym} ? Compute the relevant coefficients in $L_P C_{\text{sym}}$.
2. If $\text{Var}^{(K)}(\mu_{ij})$ is known as an explicit function of K , find an expression for the required $C_{\text{var}}(K)$.
3. Show explicitly how to specialize the resulting Markov model to our version of the Karlin-Sternoi patch model.

3.3. Other Measures of Significance

3.3.1. Spectral Analysis

In Section 3.2 we looked at correlations in nominally homogeneous DNA, “noise.” Now we will start to focus in on substructures or meaningful heterogeneities in patterns. We will soon use the evolutionarily relevant technique of recognizing signal noise by its presence in more than one species of DNA, rRNA, or protein, or a long segment of DNA, but here we discuss very briefly how we might pick up substructures by looking at correlations. Again, we need to bunch data in some fashion to reduce noise-driven details, but primitive computations such as $\text{gc}(ET)$ or $\text{gc}(AT)$ make less sense – they would for example poorly single out things like $(CT)^*$ repeats.

The simplest form of autocorrelation consists of continuing to identify bases, and, for a fragment of length N , we define the Fourier transform

$$\alpha_k(t) = \frac{1}{N} \sum_{n=1}^N e^{j2\pi t n/N} A_n(t),$$

$$\text{integer } |t| < N/2.$$

For each t , corresponding to this is the power spectrum

$$\begin{aligned} P_{tt}(t)^2 &= N |\alpha_k(t)|^2 = \frac{1}{N} \sum_{n=1}^N e^{-j2\pi t n - j2\pi t N} A_n(t) \bar{A}_n(t) \\ &= \frac{1}{N} \sum_n A_n(t) + \frac{1}{N} \sum_{n=1}^{N-1} [e^{j2\pi t (N-n)} + e^{-j2\pi t (N-n)}] \sum_{m=n+1}^N A_m(t) \bar{A}_{m-N}(t), \end{aligned}$$

which can be amalgamated as $P_t^2 = \sum_n P_{tt}(t)^2$. Noting that $(1/N) \sum_n A_n(t) = p_t$, $(1/N - |t|) \sum_{n=1}^{N-1} A_n(t) \bar{A}_{N-n}(t) = C_{tt}(t)$ and $\sum_{n=1}^{N-1} (N - |t|) \cos(2\pi N t)$ by $|t| = 0$ (and $k \neq 0$ below this), we readily find that, for large N ,

$$P_t^2 = 1 - 2 \sum_n p_n^2 + 2 \sum_{n=1}^{N-1} \left[\cos(2\pi \frac{k}{N}) + \sum_n C_{nn}(t) \right], \quad k \neq 0.$$

On the one hand, we have $P_t^2 = 1 - 2 \sum_n p_n^2$ for random occurrences of bases and $P_t^2 \approx N^{-1/2}$ for $\sum_n C_{nn}(t) \sim N^{-1/2}$ (show this itself if partial correlations are indeed a real phenomenon).

On the other hand, P_t^2 is meant for picking up repeating substructures. Suppose that a repeated pattern P is described by $A_{p+x}(x = 1, \dots, p)$ for several values of x ; then, for each one, there is a contribution to $P_{tt}(t)$ of

$$\sum_{n=1}^N e^{j2\pi t (N-x)/N} A_{p+x}(t) = e^{j2\pi t N / N} \beta_t(P),$$

where

$$\beta_t(P) = \sum_{x=1}^p e^{j2\pi t x / N} p_x$$

is the power factor for the pattern P . If the pattern is effectively random, we expect $|\beta_t(P)|^2 \sim p$. Now if the $\{\ell_j\}$, $j = 1, \dots, q$, are the intervals between occurrences of P , so that

$$R_1 = \sum_{j=1}^q \ell_j,$$

The coefficient of $\delta_0(P)$ in $a_k(x)$ will be

$$\delta_0(x) = \frac{1}{N} \sum_{j=1}^N e^{2\pi i x \delta_j / N},$$

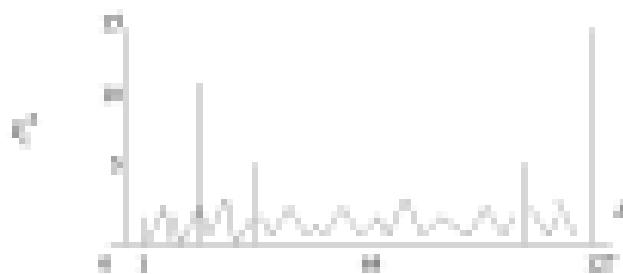
If the intervals are estimated as independent and identically distributed, with mean \bar{x} and variance s^2 (that δ_j is just a half-interval, as in A_{2M+1}), then $N = \sum_{j=1}^{2M+1} \delta_j = q\bar{x}$, and

$$\delta_0(x) \sim \frac{1}{N} \sum_{j=1}^N e^{2\pi i x \delta_j / q\bar{x}} \sim \frac{1}{N} \sum_{j=1}^N e^{2\pi i x \delta_j - 2\pi i x \delta_j^2 / q\bar{x}^2},$$

Hence if it is any multiple of q , $\delta = rj = r\bar{x}/q$, so that $e^{2\pi i x \delta_j} = 1$, we have

$$\begin{aligned}\delta_0(x) &\sim \frac{1}{N} \sum_{j=1}^N e^{-2\pi i x \delta_j^2 / q\bar{x}^2} \\ &= \frac{1}{N} (1 - e^{-2\pi i x \delta_j^2 / q\bar{x}^2}) / (2 \sin(\pi^2 x^2 s^2 / \bar{x}^2)).\end{aligned}$$

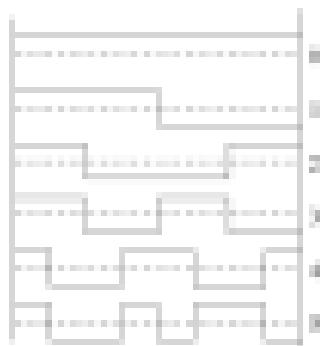
For a mutator $\sim T/\lambda$, this just gives $\delta_0(x) = 1/\bar{x}$, except that instead of $-p\bar{x}/T^2$ to δ_j^2 , which can be way above zero. As an example, the 128-kb 40-kb spacer region of Xbaege 58 DNA is shown. The strong peak at $x = 1/16$ is strong evidence of an eight-base repeated motif, and indeed the observed amplitude is several units of $p\bar{x}/T^2 \sim 4$. Of course, if $\sigma x/\bar{x}$ is not small,



$P(x)$ is quite noisy, and at higher x , the random $(1/P)x^2 \sim p$ is approached (here $x = 1/16$ in the above figure).

For coding regions, the third base of a triplet is a reflection mainly of the mean base frequency in the full sequence, and so one might expect a significant $k = 3$ peak for any normal distribution (the observation has in fact been used to detect genes). The trigonometric functions in the Fourier transform are not especially appropriate for functions that are always small integers, e.g., 0 or 1 for occupancy by a base. More suitable are the Walsh

functions, essentially discrete sine and cosine, see, e.g., Tufts and Chellings, 1985. The first five Walsh functions are shown, defined on the interval $[0, 1]$ more generally by



$$W_0(x) = 1, \quad x \in [0, 1], \quad W_0(x) = \begin{cases} 1 & x \in [0, \frac{1}{2}) \\ -1 & x \in [\frac{1}{2}, 1] \end{cases}$$

$$W_1(x) = \frac{W_0(2x)}{(-1)^x W_0(2x-1)} \begin{cases} 1 & x \in [0, \frac{1}{2}) \\ -1 & x \in [\frac{1}{2}, 1] \end{cases}$$

$$W_{2k+1}(x) = \frac{W_0(2x)}{(-1)^{k+1} W_0(2x-1)} \begin{cases} 1 & x \in [0, \frac{1}{2}) \\ -1 & x \in [\frac{1}{2}, 1] \end{cases}$$

For discrete x as well, say $x = j/N$, where $N = 2^k$, we set

$$W_k(j/N) = w(k, j) \quad \text{for } 0 \leq j, k < N = 2^k,$$

and then we define the Walsh transform for the sequence s_j belonging to the base w (e.g., $w = 0$, 1 as the base of j is not, or is, w) by

$$w_j = \frac{1}{N} \sum_{n=0}^{N-1} w(k_n, j) s_{j_n}$$

If this is really to be the analog of the Fourier transform, we will want the inverse to take the same form:

$$s_j = \sum_{n=0}^{N-1} w(k_n, j) w_j.$$

To prove this, we need orthogonality. For this purpose, the representation

$$w(k, j) = (-1)^{\sum_{n=0}^{N-1} (k_n - j_n) + k} = w(j, k),$$

where

$$J = \sum_{k=0}^{N-1} kZ, \quad I = \sum_{k=0}^{N-1} k_i Z, \quad (k, k_i = 0, 1),$$

is used, which is readily seen to satisfy the preceding recursion relation. Then

$$\begin{aligned} \sum_j w(k, j)w(k', j) &= \sum_{j=0}^{N-1} (-1)^{\sum_{i=0}^{k-1} k_i k'_{i+1} + \sum_{i=0}^{k'-1} k'_{i+1} k_i + \sum_{i=k}^{k'-1} k_i k'_{i+1}} \\ &= \prod_{i=0}^{k-1} \sum_{j=0}^{N-1} (-1)^{k_i k'_{i+1} + k_{i+1} k'_{i+2} + \dots + k'_{N-1} k_N} \\ &= N! \prod_{i=0}^{k-1} k_{i+1}(k_{i+1} + k_{i+2} + \dots + k'_{N-1} + k'_N) \\ &= N! \prod_{i=0}^{k-1} k_{i+1} = W(k, k'), \end{aligned}$$

which is all we need.

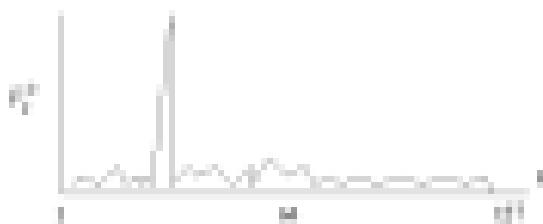
At this point, we can proceed exactly as in the Fourier case, thus defining

$$a_k(x) = \frac{1}{N} \sum_{j=0}^{N-1} w(k, j) b_j(x)$$

for each base and noting that $\sum_j b_j(x) = 1$ implies that $\sum_{k=0}^{N-1} a_k(x) = b_0(x)$, and also that $a_k(x)$ is the frequency p_k . Then comes the power spectrum

$$P_k^2 = N \sum_{j=0}^{N-1} |a_k(x)|^2.$$

This is a bit harder to interpret, but in the same example as before, period 3 certainly



comes up as a very sharp peak with respect to it. In practice, a combination of Fourier and Walsh produces a signature that is not necessarily unique, but scarcely unique.

An analytic technique designed to be particularly sensitive to structures at a given spatial scale in a hierarchical set of scales is that of wavelet

analysis (see, e.g., Meyer, 1992, for a high-level treatment). For a prototype, we might imagine replacing the transform kernel $e^{i\omega x}$ in a Fourier transform with $e^{i\omega x}e^{-i\omega(x-X)}$, creating a "local" Fourier transform that depends on both frequency ω and location X . Most (continuous) wavelet transforms that are used take the form

$$f(k, X) = k^{1/2} \int \phi(ikx - X) f(x) dx$$

The spatial ϕ , and the (k, X) are arranged, typically with $k > 2^4$, such that the $\phi_{k,x} = k^{1/2} \phi(ikx - X)$ form an orthonormal set of functions, leading to the usual sorts of expansion theorems. The $f(k, X)$ are useful descriptive parameters even if the $\phi_{k,x}$ are not orthonormal, and of course extend at once to discrete variables. For a typical application to DNA, see Trifunovic et al. (1995).

The task of discovering the words and the phrases of DNA, not to mention the proteins, from internal—i.e., sequence—indications alone prove difficult in one respect: There are a number of unequaly spaced exact and inverse repeats of long (periodic) and very long (aperiodic) identifiable subsequences. For example, there are Aata, 300 bp, 5% of the human genome on L1, every long, 4% of the human genome. These are transposable elements, RNAs mediated with reverse transcription, and are sufficient, e.g., to identify human DNA in a nonhuman cell. There are also shorter frequent motifs or "tules," picked up out of noise, e.g., by the r -mean technique of Karlin et al. (1989), in which the statistic of distances between identical or almost identical subsequences, spaced r subsequences apart, is examined for significance. However, most of the progress made in finding the longer words of DNA, RNA, and protein, has come from comparison of sequences, and it is to this that we will now turn our attention.

3.2.2. Aberrant Criteria

Special methods (see, e.g., Li, 1991) are particularly good at picking up repeats and near repeats; the human genome, for example, is estimated to consisting of approximately half repeated subsequences. Three repeats, as has been noted, occur in many places—aside from the obvious codon repeats in genes—from small tandem repeats, through 3-DNA (short interspersed equal segments, $10^2 - 10^3$ bp) and on to 12-Mb-long interspersed repeat sequences, $\sim 10^7$ bp. There are effective algorithms (Apostolico, 2000) for detecting repeats without the impossibly laborious checking of sequences after

sequences. Also, repeats are but one example of unambiguously represented subsequences; underrepresented ones occur as well. The existence of a core of strictly avoided words in fact sets up a complex structure of avoidance, which has recently been investigated (Liu et al., 2000). For a survey of the statistics of such events under geometrical random placement, see Kerlin et al. (2000). Of course, repeats are not the only special constraints found in the DNA "language"—reverses palindromes are frequent, and "splices" that are only found together are common as well (Trifunovic, 1998).

Thus the DNA language is far from random, both in its coding and in its noncoding subsequences, and a first estimate of where the information resides might proceed by finding the information content, the negative entropy, of selected regions. The entropy parallel,

$$H_1 = - \sum_{x=0}^A p_x \ln(p_x),$$

where p_x is the relative frequency of base x , is a very primitive indicator. Clearly $0 \leq H_1 \leq 2$, with $H_1 = 2$ for the purely random $p_x = 1/4$, and even GC-rich regions can reduce this substantially; but for sizable genomes, $H_1 = 1.9 - 2.0$ is unusual. More informatively, we can look at motifs of all lengths, occurring at frequency p_n^{ext} , and construct

$$H_n = \sum_m -p_m^{ext} \ln(p_m/p_n^{ext}).$$

The $n = 2$ codons in coding regions typically produce $H_2 \approx 1.5$, closer to the purely random H_1 but at large $n < n_c$ below realization in excess entropy:

$$H_n = H_{n+1} + R_n.$$

In fact, it is even at once that

$$\begin{aligned} R_n &= \sum_m -p_m^{ext} \ln_2(p_m^{ext}/p_n^{ext}) \\ &= \langle \log_2 P(\mu) | \omega^{ext} \rangle. \end{aligned}$$

averaged over a word and the following base ν , so that R_n is a measure of the unpredictability of the next base that is added. More sophisticated definitions of excess entropy (Lauwens and Yerushalmy, 1999) make below a substantial reduction, to $R_n \approx 0.6$ for a variety of genomes.

A relatively refined measure of excess entropy is obtained from

$$S = \lim_{n \rightarrow \infty} H_n.$$

where “ $\delta\ell$ ” is a bit smaller than the genomic length L , at b_2 would be last minute at b_n for an n -memory Markov chain. The full dependence of b_n on α is a better indicator (Herd et al., 1994) of a repetitive structure. We can see this quickly without going into great detail, but by forming an n -word that can overlap the left edge of an k -by- k repeat. If ℓ bases, where $1 \leq k \leq n - 1$. Suppose the repetition rate is ρ per site, so that in length L there are ρL repeats. Then, because k bases are fixed, the expected number of occurrences of the word, i.e., of matches with the pattern, is clearly $\rho L / 4^{k-1}$ for those that overlap the left edge of an k -by- k bases and $L / 4^k$ for those that do not and hence are free. $\ell^2 - L$ is the number of sites for n words that do not overlap an k . Thus the frequency of such a word is $\rho (\ell^2 - L) / 4^{k-1}$. Now if we append a base to the right end of the word, only one base δ is allowed when the word does overlap, but the expected number of $(n + 1)$ -word matches in the rest of the genome is now $L / 4^{k+1}$. Hence the entropy contribution to b_2 is

$$\begin{aligned} -\log_2 \Pr(\delta | \text{wt}) &= \log_2 \left(4^{k-1} + \frac{1}{4^k} \right) / \left(\frac{\rho}{4^{k-1}} + \frac{1}{4^{k+1}} \right) \\ &= \delta - \log_2 \left(1 + 4^{k+1} \rho / (1 + 4^k \rho) \right), \end{aligned}$$

which changes from 2 at small δ to 0 at large δ , a rapid-change occurring for $\delta \approx \log_2 (1/\rho)$. So if $n < k$, then dependence will be smooth, but pass through b_n a rapid shift occurring, giving us an estimate of the frequency ρ .

Sequence Comparison

Having examined aspects of the general structure of the language of DNA, we continue to focus in on the words, phrases, etc., of the language. A word of course is a subsequence that occurs in the same or another sequence, either exactly, or distorted, or in squeeze form; a phrase consists of key words together perhaps with filler. We will save for Section 4.3 the question of how potential words or phrases are located in the first place in our new more general context, and concentrate now on the degree of confidence with which we can assert that these objects have indeed been found. It must be emphasized that we are only attending to linear reading, the primary structure of the molecule, so that proteins, for example, correspond simply to coding linear subsequences; the feature in which the distinctive three-dimensional structure arises, clearly crucial for proteins and hardly irrelevant for DNA, is not being addressed.

4.1. Basic Matching

The principal problem situation is by analogy to this: Two linear chains of length ℓ (nucleotides, amino acids, . . .), when aligned, are found to have a common contiguous subsequence of r units. What is the probability that this was a random event and not an indicator of a functional or co-evolutionary relationship between the chains? At the most primitive level of resolution, random means independent selection of the units at the overall frequencies, p_n , for the n th type of unit, $n = 1, \dots, n$ ($n = 4, 20, \dots$). In this case, assuming the same as that of Subsection 3.2.1, the probability of a match at a given location will be

$$p = \sum_{n=1}^N p_n^2,$$

and the probability of not matching will be $q = 1 - p$. Note that the term "match" depends very much on the way equivalent units are defined; the

20 amino acids, for example, can for most purposes be simplified to four or five equivalence classes.

4.2.1. Mixed-Correlation Model

Our previous question will be: what is the probability P , under the assumption of randomness, that a match of length $\geq r$ will occur? Then, if this event occurs, randomness can be rejected in favor of significance at a confidence level of $1 - P$. The key to our present analysis (Pearson and Pearson, 1994) is that the result will be interesting precisely when P is small, so that any of the probabilities of events of which it is composed will be very small. This means that although the different matches that can occur are not mutually exclusive, i.e., if $A \cap B = A \cdot B$, we have $P(A \cdot B) \ll P$; conversely, $P(A + B)$ will be several orders of magnitude, and we can add at all:

$$P(A \cup B) = P(A) + P(B).$$

i.e., just add the component probabilities.

Back to $\geq r$ matches ("successes") out of l : the possibilities are that

1. at least r matches in sequence (no term this as r match) start at the left end, probability p^r .
2. an r match starts at one of $a = 2, 3, \dots, l+1-r$; hence a mismatch at location $a+1$ is followed by r matches, probability $(1-p)^a p^r$.

The single-events of type 1, and $l-a$ of type 2 give us, under the assumption of additivity,

$$\begin{aligned} P_{\geq r} &= p^r + (l-a)p^r \\ &= (1+p)^{-a}p^r \quad \text{for } r \geq 1. \end{aligned}$$

Example: At 0.0001 confidence level, a tail probability of $P = 0.0001$, and with $p = 0.04$, a match of $r \geq \log(0.0001)/\log(1-0.04) = 16.790$, i.e., $r \geq 17$. But $l = 10^6$, would be significant.

Now, more realistically, we consider chains of lengths $l_2 \leq l_1$ and ask for a match of $\geq r$ units in a row in one chain with $\geq r$ in a row elsewhere along the other, i.e., not necessarily in register. We do this by sliding one chain over the other, and we treat the lengths in succession as a pair of chains in register, to which we can apply the above result. We have the following possibilities:

1. The l_2 chain fits inside the l_1 chain at $l_1 + 1 - l_2$ different starting places, each contributing $(1+p)^{-a}p^r = r!p^r(1-p)^a$.

2. The A_1 chain starts to the left of B_1 , with $i = l_1$ (and $j \geq r$) also non-overlapping, a contribution of $|i + j - r|q(|j|)$.
3. 2 is repeated on the right. Hence

$$\begin{aligned} P_{1r}^{(2)} &= l_1 + 1 - \text{Ad}(1 + (l_1 - r)q(|j'|) + 2 \sum_{i=1}^{l_1-1} (1 + i_1 - r)q(|j'|)) \\ &= |1 + (l_1 - r) + (l_1 - r) + (l_1 - r)(l_1 - r)q(|j'|), \quad 1 \leq r < l_1. \end{aligned}$$

Example: Two DNA fragments, of lengths $l_1 = 154$ and $l_2 = 103$, with nominal $p = 1/4$, would have, if exactly matched, $P_{11} = 0.008$. The above estimate yields $P_{11} = 0.048$, which is very close.

Of course, some substitutions can be tolerated without changing functionality of biological subsequences. Suppose for example that the matching score is $|j'|^r$, but only $m = r$ sites have to match, including as well the first and the last that define r . Now the elementary composite match probability p , instead of being p^r , is determined by one match, $m - 1$ out of $r - 1$ matches, then another match, a probability of $p \sum_{m=1}^{r-1} q^{m-1} q^{r-m} p$. We conclude that the only effect is the replacement:

$$\text{for } m \text{-out-of-} r, \text{ replace: } p^r \rightarrow \binom{r-1}{m-1} p^m q^{r-m}.$$

More reasonable might be that at least m -out-of- r match, leading to

$$\text{for } \geq m \text{-out-of-} r, \text{ replace: } p^r \rightarrow \sum_{m=1}^r \binom{r-1}{m-1} p^m q^{r-m}$$

The assumption that all successful r matches are "space filling," and therefore mutually exclusive, is clearly an assumption that fails only if the elementary events have sufficiently small probability that the mutual probabilities are fully negligible. If we are not sure, some estimate is needed. The simplest goes like this. Consider the basic matching problem. If A_{ij} is the event that a match, at $\geq r$ contiguous sites, starts at j , for $j = 1, \dots, l + 1 - r$, then

$$P_{1r}^I := P(A_{11} \cup A_{12} \cup \dots \cup A_{1l+1-r}).$$

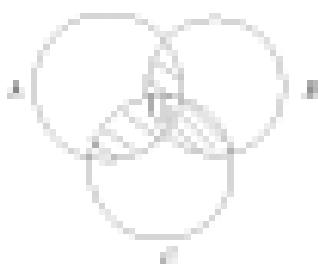
However, now we can use the first of the Bonferroni inequalities (see, e.g., Peleg, 1996, p. 186). These are most easily obtained in terms of Venn diagrams on the space of elementary events. Representing the composite events A, B, C, \dots by circles and simply counting, we have $P(A \cup B) =$



$P(A) + P(B) = P(A \cup B)$, so that

$$P(A) + P(B) \leq P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

At the next stage,



$$\begin{aligned} P(A \cup B \cup C) &= P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) \\ &\quad - P(B \cap C) + P(A \cap B \cap C), \end{aligned}$$

so that

$$\begin{aligned} P(A_0) + P(A_1) + P(A_2) &\geq P(A_0 \cup A_1 \cup A_2) \geq P(A_0) + P(A_1) + P(A_2) \\ &\quad - P(A_0 \cap A_1) - P(A_0 \cap A_2) - P(A_1 \cap A_2), \end{aligned}$$

and, more generally (Przytak, 1940),

$$\begin{aligned} \sum_{j=0}^{\min(r,t)} P(A_j) &\geq P\left(\bigcup_{j=0}^{\min(r,t)} A_j\right) \\ &\geq \sum_{j=0}^{\min(r,t)} P(A_j) - \sum_{A_j \text{ not } r \text{-match}} P(A_j \cap A_k). \end{aligned}$$

Applying this to the basic matching problem, we have $P(A_0) = p^r$, and $P(A_{j,k}) = qp^j$ for $j \geq 0$. If $j < k \leq j+r$, then $A_{j,k}$ consists of those $A_{j,k}$ such that k cannot be the start of an r match, and $P(A_{j,k} \cap A_k) = 0$. If $k > j+r$, A_j and A_k are independent, so $P(A_{j,k} \cap A_k) = P(A_{j,k})P(A_k)$. Thus the non-matching count that goes to the correction comes from $j = 0, r+1 \leq k \leq l+r-1$, a total of $(l-r)qp^rqp^{r-k}$, and $0 \leq j \leq l-2r, r_j < k \leq l+1-r$, a total of $\sum_{j=0}^{l-2r} (l-j+1-2r-j)qp^rqp^{r-k} = [(l-2r+1)(l-2r)qp^r]^2$. We conclude that (see, e.g., Uspensky, 1937):

$$\begin{aligned} (l+r-1)qp^rqp^r &\geq P_{\text{err}}^l \geq [(l+r-1)qp^r]qp^r \\ &\quad - (l-2r)\left[1 + \frac{2}{3}(q-p-1)\right]qp^{2r}. \end{aligned}$$

Example. For a "fair" coin, $p = \frac{1}{2}$, tossed 2000 times, the actual result to four places is $P_{111} = 0.1793$. From our inequality, $0.1249 \leq P_{111}^{(200)} \leq 0.1772$, indicating as well the rapid convergence of the alternating series of which we have used only the first two terms.

The general result is that, for the first-order error, $\delta \hat{P} \sim P^2$. However, even more importantly, the extremely simple approximation that we get with the multinomial index model is precisely what we need to ensure nonoverexposure at a given level of confidence.

4.2.2. Independence Model

A different kind of approximation assumes that successive matches are independent events, as indeed most of them are. Then $1 - P = P(A_1 \cdot A_2 \cdots) \approx \prod_i P(A_i) = \prod_i (1 - P(A_i))$, for small $P(A_i)$.

$$1 - P \approx \exp\left(-\sum_i P(A_i)\right)$$

In Poisson distribution form, here, too, the corrections can be arranged to supply bounds on the error made. Let us introduce the variable $X_i := (0, 1)$, with $A_i := 1 - X_i$, and use the notation

$$\begin{aligned} A_i \text{ occurs} &\iff X_i = 1, \\ A_i \text{ does not occur} &\iff X_i = 0. \end{aligned}$$

We can then, quite generally, proceed as in Subsection 3.2.2, but we shall do so with a little more detail. We want to find $1 - P = \Pr[\prod_i X_i = A_i | \prod_i (1 - X_i)]$, but instead we look at

$$\delta(x) = \ln\left(\frac{\prod_i (1 - x X_i)}{\prod_i (1 - x X_{i'})}\right).$$

Now

$$\begin{aligned} P(x) &= -\left(\sum_i X_i \prod_{i' \neq i} (1 - x X_{i'})\right) \Bigg/ \left(\prod_i (1 - x X_i)\right), \\ P(x') &= \left(\sum_{i'} X_{i'} X_i \prod_{i'' \neq i, i'} (1 - x X_{i''})\right) \Bigg/ \left(\prod_i (1 - x X_i)\right) \\ &\quad - \left[\left(\sum_i X_i \prod_{i' \neq i} (1 - x X_{i'})\right) \Bigg/ \left(\prod_i (1 - x X_i)\right)\right]^2. \end{aligned}$$

Clearly

$$\begin{aligned} J(0) &= 0, \quad J'(0) = -\sum_i (X_i), \\ J''(0) &= \left\{ \sum_{i,j} X_i X_j \right\} - \left\{ \sum_i X_i \right\}^2 \\ &= \sum_{i,j} (X_i - \bar{X})(X_j - \bar{X}) = K_1(1 - \sum_i K_i)^2, \end{aligned}$$

so that if

$$\lambda := \sum_i (X_i),$$

a Taylor series expansion gives

$$\begin{aligned} J &= 1 - e^{J(0)} = 1 - e^{-\lambda} e^{\lambda \sum_i K_i}, \\ &\approx 1 - e^{-\lambda} + \frac{1}{2} e^{-\lambda} J''(0), \dots \end{aligned}$$

Still being very general, suppose there is a distance $d(i, j)$ defined with respect to indices such that X_i and X_j are independent for all $|i - j| > r$. Then $J''(0)$ simplifies to

$$\begin{aligned} J''(0) &= \sum_{\substack{0 \leq i, j \leq n \\ |i - j| > r}} (X_i X_j) - (X_1)(X_r) = \sum_i (X_i)^2 \\ &= k_2 - k_1 \end{aligned}$$

where

$$k_1 = \sum_{\substack{0 \leq i, j \leq n \\ |i - j| \leq r}} (X_i)_+ (X_j)_+, \quad k_2 = \sum_{\substack{0 \leq i, j \leq n \\ |i - j| > r}} (X_i X_j)_+$$

The Chen–Stein theorem (Chen, 1975) guarantees this is a somewhat more conservative but rigorous bound:

$$k^2 - (1 - e^{-\lambda}) \lambda \leq \frac{1 - e^{-\lambda}}{\lambda} (k_1 + k_2).$$

In particular, in our typical applications, X_1 and X_r will be mutually exclusive for all $|i - j| \geq r$ but $i \neq j$, and we are well have $k_2 = 0$.

Computation for the Chen–Stein estimate (Pensia et al., 1999) requires us precisely the same information as we needed previously in Subsection 4.1.1. We must make sure the the required independence and mutual exclusivity are satisfied, and we will formalize our previous assumptions a bit for this purpose.

Consider first the longer-than- r match problem (no mismatching). Let the independent variables $C_{ij} = 0$ if there is failure or success of a match at site i , then X_j for a success run $\geq r$ starting at site j is simply

$$X_j = (1 - C_{j,j+1}) \prod_{i=j+2}^{j+r-1} C_{i,i+1} \quad \text{if } C_j = 0,$$

Clearly X_j is independent of X_l if $j > l + r$ (they have no C_{ij} in common), whereas $E(X_j) = 0$ if $j < l + r$ (X_j contains $C_{j,j+1}$ but X_l contains $1 - C_{j,j+1}$). Furthermore, $b_1 = (\sum X_j) = [1 + (l - r)p]p^r$, as in Subsection 4.1.2, and $b_2 = \sum_{1 \leq j < l \leq r} (X_j)(X_l) = [(l - r)p]^2 + \frac{1}{2}(l - 2r)(l - 2r - 1)p^2p^{2r}$, for which a very good estimate is (show this)

$$b_2 \approx \frac{2r + 1}{4} b_1 + 2rp^r.$$

Example: Fair coin, $l = 2040$, $p = \frac{1}{2}$, $r = 14$:

$$|P_{(1,1)}^{\text{full}} - 0.04890| \leq \delta \approx 10^{-6}.$$

Continuing to two sequences, (X_i) of length b_1 and (Y_j) of length b_2 , we define $C_{ij} = 1$ if $X_i = Y_j$ and $p = \langle C_{ij} \rangle$. Now the indices are pairs, and, for an r run starting at (i_1, j_1) ,

$$X_{ij} = (1 - C_{i+1,j+1})(C_{i,j} C_{i+1,j+1} \cdots C_{i+r-1,j+r-1}).$$

Here $\delta(i_1, j_1) = \min(|i_1 - i|, |j_1 - j|)$, and again $b_1 = 0$. In the same fashion as before, we find

$$\begin{aligned} b_1 &= [(l_1 + b_1 - 2r + 1) + (l_1 - r)(l_1 - r)p]p^r, \\ b_2 &\approx \frac{(2r + 1)b_1^2}{(l_1 - r + 1)(l_1 - r + 1)} + 2rp^r. \end{aligned}$$

4.2.5. Direct Asymptotic Evaluation

On the assumption of independent base placement in DNA (mutation by low-order Mutation changes very little (Kaslin et al., 1989)), many matching problems can be solved exactly, in principle, and in rapidly convergent asymptotic series. In practice, Consider once more the prototypical run of n r successes in a chain of length l . We have seen in Subsection 3.2.1 that for the complementary probability $P_{(r,r)}^{(l)}$ of having no run of r or more, thus decomposing the sequence into successive runs of less than r , separated by

before, some results.

$$\begin{aligned} G(x) &= \sum_{k=0}^{\infty} P_{k,r} x^k = \sum_{k=0}^{\infty} \sum_{i=0}^{k-1} q x^{k-i} p x^{k-i} p x \cdots (px)^i \\ &= [(1 - (px)^r)(1 - x + px^r)]. \end{aligned}$$

The simplest expansion is in $y = px$ at fixed r :

$$\begin{aligned} P_{k,r} &= \text{coeff}_y^k \ln(1 - px^r/(1 - x + px^{r+1})) \\ &= \text{coeff}_y^k \ln \frac{1}{1-x} \frac{1 - px^r}{1 + px^r q x/(1-x)} \\ &= \text{coeff}_y^k \ln \frac{1}{1-x} (1 - px^r) \left[1 - px^r \frac{q x}{1-x} + p x^{2r} \frac{q x^2 r^2}{(1-x)^2} + \dots \right] \\ &= \text{coeff}_y^k \ln \frac{1}{1-x} \left[\left(\frac{x^r}{1-x} + q \frac{x^{r+1}}{(1-x)^2} \right) y \right. \\ &\quad \left. + \left[\frac{x^{2r+1}}{(1-x)^2} + \frac{x^{2r+2}}{(1-x)^3} \right] y^2 + \dots \right]. \end{aligned}$$

picking:

$$\begin{aligned} P_{k,r}^j &= 1 - P_{k,r}^j \\ &= [(1 + y) - r](y - \left[q(1 - 2x) + y^2 \binom{r - 2x}{2} \right]) y^j, \dots \end{aligned}$$

just as before.

However, we can just as easily have more structured matching events (Peng et al., 1997). Suppose we demand at least k matches sequences of length $\geq r$ to declare a match. This is almost as easy. Now we “tag” any subsequence of r or more matched by a variable y in order to recognise it. Then we replace the strict failure condition $(\sum_{i=0}^{k-1} p^i x^i) = (1 - p^k x^k)(1 - px)$ with $\sum_{i=0}^{k-1} p^i x^i + k \sum_{i=0}^{k-1} p^i x^i = (1 - p^k x^k + p^k x^k)(1 - px)$, giving us instead

$$G(x,y) = \frac{1 - p^k x^k + k p^k x^k}{1 - x + q p^k x^{k+1}(1 - y)}.$$

Failing configurations are those in which only the powers $1, 2, 3^2, \dots, j^{k+1}$ are present. Because

$$(\text{coeff}(q^k + \text{coeff}(q^1 + \dots + \text{coeff}(x^{k-1}) q))_j(y,1) = \text{coeff}(x^{k-1} \frac{G(x,1)}{1-y}),$$

we conclude that now

$$P_{\text{DP}}^{(1)} = \text{mult}(t^2) q^{r+2} \ln \frac{1}{1-t} \frac{(1-pz) - p(z^2)}{(1-z + pqz^{r+1})},$$

from which

$$P_{\text{DP}}^{(2)} = p^2 q^2 \text{mult}(t^2) \ln \frac{q^{2(r+1)} - p^{2(r+1)}(1-t)}{(1-z + pqz^{r+1})^2},$$

and so, to leading order, which is p^2 , we have

$$P_{\text{DP}}^{(2)} = t^2 q^2 \left[\binom{1+1-r}{r} - p \binom{1-r}{r} \right] + \dots.$$

In particular, if $r > n_1$, then $\binom{1+1-r}{r} \sim (1-r)/n_1$, $\binom{1-r}{r} \sim (1-r)/n_1$, $q^{2(r+1)} \sim 1 - r^2/n_1$, and therefore

$$\begin{aligned} P_{\text{DP}}^{(2)} &\sim p^2 q^2 (1 - r^2/n_1) \\ &\sim p^2 q^2 t^2 e^{-r^2/2n_1}. \end{aligned}$$

4.2.4. Extreme-Value Techniques

When an analysis of the significance of an observed sequence is carried out, it is certainly preferable to focus first on a single criterion and then examine the confidence with which we can assert significance. Such a single criterion, as was mentioned in Subsection 4.2.1, might be an extreme value like the length of the longest match found, to be assessed by its random reference-expectation value, variance, and so forth. Now, at fixed L ,

$$\Pr(\tau_{\max} = R) = P_{\text{DP}} = P_{\text{DP},L},$$

Hence

$$\begin{aligned} E(\tau_{\max}) &= \sum_R R P_{\text{DP},L} = P_{\text{DP},L} \bar{\tau}_R \\ &= \sum_{k=1}^n k P_{\text{DP},L} = \sum_{k=1}^n k (R - 1) P_{\text{DP},L} \\ &= \sum_{k=1}^n k P_{\text{DP},L}, \end{aligned}$$

and, in the same fashion,

$$E(\tau_{\max}^2) = \sum_{k=1}^n k^2 (R - 1) P_{\text{DP},L}.$$

The asymptotic analysis of Definition 4.1.1 is not valid in this paper, as typical probabilities such large than 1 are small. The Chernoff bounds are not valid for such values of α , and indeed Definition 4.1.1 necessarily be considered asymptotic theory, the asymptotic law theorem, especially regarding the case when asymptotic posterior estimates are stated and cannot be assumed to converge with probability one, e.g., Wasserman (1993), Chapter 11.

In asymptotic theory consider the standard coin problem. As we recall that Bernoulli's law, under H_0 both the two-sided hypothesis distribution has large λ . Following a similar thought it follows that the prior law probability $p\mu^{\lambda}$, under probability of type H_0 is $\sum_{k=0}^{\infty} p\mu^{\lambda} = 1 - p^{\lambda}$. However, then

$$\lim_{\lambda \rightarrow \infty} \frac{P_{H_0}}{P_{H_1}} = 0,$$

meaning that the H_0 belief is zero. Hence $P_{H_1} = P_{H_0 \cup H_1} = 1 - P_{H_0} = 1 - p^{\lambda} \approx 1 - p^{\lambda}$. P_{H_1} is given by $P_{H_1} = P_{H_1 \cap H_0^c} = P_{H_1} \cap H_0^c = 1 - (1 - p^{\lambda})^2 \approx 2p^{\lambda}$, the former classical Bayesian probability. Then, assuming normality hypothesis, we obtain

$$E_{H_0 \cup H_1} = \int_{\mathbb{R}} (1 - p^{\lambda})^2 d\mu.$$

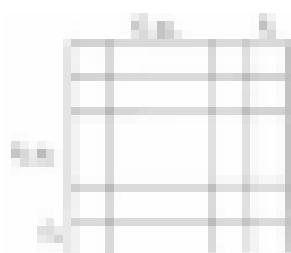
If we define a large λ , then $E_{H_0 \cup H_1} = \int_{\mathbb{R}} (1 - p^{\lambda})^2 d\mu$, or $E_{H_0 \cup H_1} \approx 2p^{\lambda}(1 - p^{\lambda})^2 \approx 2p^{\lambda}(1 - p^{\lambda})$. The large λ 's are near the integral $\int_{\mathbb{R}} (1 - p^{\lambda})^2 d\mu$, which is $\int_{\mathbb{R}} (1 - p^{\lambda})^2 d\mu = \int_{\mathbb{R}} (1 - p^{\lambda})^2 d\mu = \lambda p^{\lambda} e^{-\lambda} = \lambda p^{\lambda} e^{-\lambda} \approx \lambda p^{\lambda}$.

$$E_{H_0 \cup H_1} \approx \lambda p^{\lambda} e^{-\lambda} \approx \lambda p^{\lambda}.$$

The evaluation of the maximum value function, but easier to show that

$$E_{H_0 \cup H_1} \approx \frac{\lambda^2}{2} p^{\lambda} e^{-\lambda} \approx \dots$$

Now let us consider the hypothesis H_0 and H_1 the same from the previous section and compute the resulting sufficient statistic we will perform here. It compares all regions, and τ stands for our sufficient statistic of



$(1 - \alpha) r$ out of r successive comparisons. Consider how α there are already, on the average, $b_1 p_1 \cdot b_2 p_2$ possible matches of bins i between the two datasets, and $\sum_i b_i p_i \cdot d_{\text{true}} = 1/p_1 p_2$ matches of any bins; hence there are $(1 - b_1 p_1) \cdots (1 - b_r p_r)$ failures in the complete set of comparisons. Again, our match begins after a failure. We then need $\Pr[\text{Match} \leq R] = \Pr[\text{Match} \leq r]$ if $R \leq R_0$. If the bins matches were contiguous, this would be the standard Poisson with $r = R$, and no restrictions on exceeding sizes), and although this is only an approximation for $n \neq 0$, we will use it anyway. Because there are $(1 - \alpha)R$ successes and αR failures required, the latter quantity will be given by $(\frac{R}{R_0})^R e^{(1-\alpha)R} q^{R_0}$, or now $\Pr[\text{Match} \leq R] = \Pr[\text{Match} \leq R_0] = (1 - (1 - \alpha)^R q^{R_0})^{-1} = \exp[-\Pr[\text{Match} \leq R_0] q^{R_0 - R}]$. We can use the central-limit theorem to approximate the binomial by a Gaussian:

$$\binom{R}{R_0} e^{(1-\alpha)R} q^{R_0} \sim \frac{1}{\sqrt{2\pi R p_1}} e^{-(1-\alpha)R q_{\text{true}} - \alpha^2}.$$

Setting $a = b_1 p_1 (1 - \sqrt{2\pi R p_1})e^{-(1-\alpha)R q_{\text{true}} - \alpha^2}$, we have $\alpha/\sqrt{n} = -(1/2R) + \beta p_1 = \alpha^2/(2\pi R p_1) \alpha R$, or dropping the $(1/\alpha R)$ in the dominant large R regions, but taking it multiplying on the binomial, we have $\Pr[\text{Match} \leq R] = \int_0^R \Pr[\text{Match} \leq R] dR = (\sqrt{2\pi p_1}/\beta p_1 - \alpha^2/2) \int_0^{\alpha R} (1 - e^{-x^2}) dx/\alpha$, so that

$$\Pr[\text{Match}] = \frac{\Pr[\text{Match} \leq R_0]}{(\alpha - \alpha^2)} \ln(\alpha/\beta) + \dots$$

with $\text{Var}(\text{Match})$ again (βp_1) independent.

Two points are to be made. First, the coefficient of $\ln(\alpha/\beta)$ diverges as $\alpha \rightarrow \beta p_1$, so that weakening the match criterion leads to the breakdown of the law term. Second, however, the central-limit theorem only applies to a large neighbourhood of $q = \alpha$, not including the tail, which dominates for small α ; thus it is incorrect for $\alpha = 0$. This is trivially avoided: instead of the normal approximation by the binomial, we just insert Stirling's approximation $n! \sim \sqrt{2\pi n} (n/e)^n$ into the binomial coefficient. This gives us, after a little algebra,

$$\binom{R}{R_0} e^{(1-\alpha)R} q^{R_0} \sim \left[\left(\frac{\beta}{1-\alpha} \right)^{1-\alpha} \left(\frac{\alpha}{\beta} \right)^{\alpha} \right]^R \int \sqrt{2\pi R \alpha(1-\alpha)}$$

and changing nothing but the coefficient of $\ln(\alpha/\beta)$, we now find

$$\Pr[\text{Match}] = \ln(\alpha/\beta) \alpha / \left[(1-\alpha) \ln \frac{1-\alpha}{\beta} + \alpha \ln \frac{\alpha}{\beta} \right] + \dots$$

This is indeed correct as $\alpha \rightarrow 0$, and coincides with the normal approximation as $\alpha \rightarrow \beta p_1$.

Assignment 8

- Suppose that base correlations were fixed in time; what effect would this have on the Fourier power spectrum and on the Walsh power spectrum?
- Apply the mutual-exclusion model to joint $m \times r$ matches for simultaneous comparison of three sequences.
- For comparison of two sequences, how does amalgamating units explicitly affect the statistics of r_{max} ? How does grouping units, i.e., into nonoverlapping doublets or triplets, affect it?

4.2. Matching with Imperfections

4.2.1. Score Distribution

The general procedure in comparing two sequences, protein, DNA, ..., is to algorithmically produce a best alignment and then assess the alignment for significance. If found, we deduce functional similarity, an old example being between platelet growth factor and the κ -sis oncogene product, suggesting a growth factor in the latter. Of course, because of ambiguous replication errors, mutations, evolution, etc., the matches need not be perfect to reflect similarity. We now study how to include this possibility.

When the characterization of an imperfect match is no longer as simple as the length of a sequence of perfect matches (the r_1 function of 4.1), it is convenient to define a single quantity, here called S , to epitomize the quality of the match. The (local) score of a two-sequence comparison, *i.e.*, S , will be defined as the maximum score of aligned subsequences conformations, *i.e.*, of $s(I, J)$, where $I \subseteq A$ and $J \subseteq B$; at this stage, I is required to be a contiguous subsequence, as is J . (Please note that, in the basic score $s(I, J)$, the crucial statistical datum for the two-sequence comparison is the cumulative tail probability

$$\begin{aligned} P(S) &= \Pr[s(I, J) \geq S] \\ &= 1 - \Pr[s(I, J) < S] \\ &= 1 - \Pr[s(I, J) < S, \quad \forall (I', J') \subset (I, J)] \\ &= 1 - \left(\prod_{(I', J') \subset (I, J)} [1 - \Pr[s(I', J')]] \right), \end{aligned}$$

where $\Pr[S] = 1$ is the indicator of the event that $s(I, J) \geq S$; otherwise $\Pr[S] = 0$. As a rule, the failure of some intelligently chosen subset denoted by (I, J) , of the full set of (I', J') , is sufficient to imply that of the full set. Now

If the corresponding $\lambda_{i,j}(S)$ are small in probability and nearly independent, the effect of the Chen–Stein theorem, and of any number of approximation methods (Goldstein, 1993), is to reduce the preceding probability to the Poisson form:

$$P(S) = 1 - \exp \left\{ - \sum_{(I,J) \in \Omega_{\leq S}} \Pr[\alpha(I, J) \geq S] \right\}.$$

For this equation to be valid and useful, it is important to select the pairs (I, J) to avoid redundancy. This of course does not imply the absence of further collision between subsets of pairs; the assumption that it does, effectively restricts us to “groupy” match criteria.

In the primitive case of exact matching of units, in which the score is the maximum matching length, we can use the same strategy as in the extreme-value technique of Subsection 4.1.4. Suppose that I and J have lengths l_1 and l_2 . Then we choose a pair (i_1, j_1) of starting positions of I and J , which are of lengths precisely S , restricted so that $(i_1 + 1, j_1 + 1)$ is not a matching pair. Of the $(l_1 + 1 - S)(l_2 + 1 - S)$ possible pairs, just $q(l_1 + 1 - S)(l_2 + 1 - S)$ on average, but sharply distributed, will have the desired property. A complete match of I and J satisfies $\Pr[\alpha(I, J) = S] = p^S$, independently of I and J ; the corresponding events both exhaust all possibilities of $\alpha(I, J) \geq S$ and are nearly independent of each other. Neglecting S compared with l_1 and l_2 , we see that the desired form:

$$P(S) = 1 - \exp(-k(S) q p^S)$$

follows at once.

We can relax the strict definition by allowing up to d mismatched binary subsequences comparisons, but we retain the score α as the total length of the comparison made; the test is then parameterized by k . With this criterion, we can still proceed by choosing (i_1, j_1) to always follow a failure to match, if a pair is at, and then impose the condition of $\leq d$ mismatches on the next S units. The probability is $\Pr[\alpha(I, J) = S] = \sum_{m=0}^d \frac{S!}{m!(S-m)!} q^{S-m} p^m$, which again includes all $\alpha \geq S$ in larger-subsequence pairs. Although the weak dependence of the I, J matches is further increased, we can still approximate $P(S) = 1 - \exp(-k(S) \sum_{m=0}^d \frac{S!}{m!(S-m)!} q^{S-m} p^{m+1})$. A function of altered mismatched patterns in the same fashion.

The above tests are all characterized by a representation in which the only quantities required are

$$\Pr[\alpha(I, J) = S] = C(S) q p^S$$

for a selected subsequence $\{I, J\}$, and there is no explicit I, J dependence. Here $C(S)$ is a slowly varying function (at least slower than exponential) for large S ; ρ is to be computed as

$$\rho = \lim_{T \rightarrow \infty} (\Pr[\ln T, I]) = T_0^{1/C},$$

and $C(S)$ is the remaining ratio. If $\Pr[\{I, J\} = S] = \sum_m a_m(S)$ is the sum of a series of at most B^m (positive) terms for fixed m , then $\max_m a_m(S) \approx \Pr[\{I, J\} = S]$, $a_m(S)$, and we have simply

$$\rho = \lim_{S \rightarrow \infty} \left[\exp a_m(S) \right]^{1/C}.$$

For an exact match, of course $\rho = p$, $C = 1$; for $\leq k$ mismatches, $\rho = p$ but $C(S) = (S_0/p)^2/4k$. For a fraction $a < p/q$ of mismatches, $\rho = (p/1 - a)^{1/C} (p/q)^2$ and $C = \infty$. In all of these cases, we have seen (Subsections 3.1–4.2) that

$$\begin{aligned} D(\text{Q}_\text{match}) &= \ln(1/\ln(1/\rho)), \\ \text{Var}(\text{Q}_\text{match}) &= \rho(p/q)^2/(p\ln(1/\rho)^2), \end{aligned}$$

where \mathbb{E} is the number of I, J pair requests. Of course, this characteristic $\ln(1/\rho)$ dependence is valid only if $\rho < 1$; the $\rho = 1$ situation must always be handled separately.

A more intuitively satisfying way of defining a similarity score is by rewarding matches and punishing mismatches. The prototype is the linear expression

$$x(I, J) = r(I, J) - \lambda \ln T, I, J,$$

where r is the number of matches, λ of mismatches, and the x_match is parameterized by b , $b = \infty$ would reward only runs of exact matches, whereas $b = 0$ would accept matches no matter how far they are spread apart. We can analyze this very similarly, while realizing that now $x = S$ for a pair I, J no longer implies $x \geq S$ for supersequences. Thus the common length of I and J , say n , is another parameter. After choosing the starting pair (I_0, J_0) , we must find the probability that $x_\text{match} \geq S$ over all (I', J') starting at (I_0, J_0) . If their common length is n , the probability that $x \geq S$ is simply $\max_m \sum_{\substack{I' \in \mathcal{I}, J' \in \mathcal{J}, \\ |I'|=n+1, |J'|=n+1}} C(I', J') q^4$, but there are two regions to examine. If $\lambda = p/q$, then the maximum of the summand, at $T = n/p$, $R = n/q$, occurs outside of the cross-hatched region $R - \lambda b \geq S$. Thus the maximum term in the sum occurs at $R + b = n$, $R - \lambda b = S$ (i.e. within a bracketed rectangle) and the probability that $x \geq S$

For some value of n it follows

$$\sum_{r+k \geq n+1} \binom{r+k}{r} p^r q^k \leq \Pr_{\text{max}} \leq 2 \leq \max_{r+k \leq n} \binom{r+k}{r} p^r q^k.$$

Clearly then

$$\rho = \lim_{n \rightarrow \infty} \Pr_{\text{max}} \leq D^{(2)} = \lim_{n \rightarrow \infty} \left[\max_{r+k \leq n} \binom{r+k}{r} p^r q^k \right]^{1/n}.$$



For large k , it must be small, so observing that $\rho = p$ in the special case $r = 0, k = 0$, we see that $\rho \leq p$. But k not far from p/q , we are close to the maximum of the binomial probability (at the maximum if $r, k = \lfloor np \rfloor$), which can likewise be approximated by the normal distribution

$$\binom{r+k}{k} p^r q^k \approx [(2\pi(p+1/p))^{1/2} \exp(-\frac{1}{2})] [r + p(r+1/p)^2/(r+1/p)]^k.$$

Setting $\delta = SK$, we have $r = S(\delta + 1/K)$, and we find that this becomes

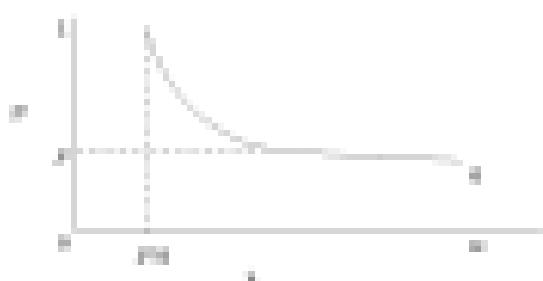
$$(2\pi(pqS)(1+1/K)^2)^{-1/2} \times \exp(-(S^2(pq/\delta p + (\delta p - p)K^2)/S) + (\delta + 1/K)).$$

whose asymptotic behaviour is $\exp(-S^2(2pq/\delta p + (\delta p - p)K^2)/S + (\delta + 1/K)^2)$. On maximising over K , we find

$$\rho \approx \exp \left[- \frac{2}{(1+2p/q)} \left(\frac{S}{p} (\delta - 1) \right) \right].$$

On the other hand, if k is large compared with p/q , only $r \approx 2$ and $k \approx 0$ contribute to ρ , which tends to the approximation above.

If $k < p/q$, the maximum in the previous figure is in the allowed region, so that $\Pr_{\text{max}} \geq 2^{-1/n}$ for large enough n ; this is because the penalty is small enough that the match sequence will simply become as large as



possible), and we expect that, for sequence length $\rightarrow \infty$, $P(\Delta_{\text{match}}) = 2p - 2q = 2p - 2p^2/(p+q) = \lambda$. There is thus a "phase transition" (Waterman et al., 1980) from a local to a global dependence when $\lambda = p/q$.

A nonparametric score incorporating mismatch/match assessment has been suggested and used (Dotsch, 1992), based on the local likely configurations sequence of matches and mismatches. In this version of this very general strategy, we set

$$J = -\ln \left(\frac{r+k}{r} \right) p^r q^k,$$

using the same notation as before, and seek the maximum J over r and k , starting at $(1, 1)$. We estimate

$$P(\Delta_{\text{match}}(i, j); r, k) \approx J \text{ as } \sum_{\text{mismatches}} \binom{r+k}{r} p^r q^k,$$

where J_k is the set of (r, k) defined by $\binom{r+k}{r} p^r q^k \leq e^{-k}$. But to compute the parameter p_0 , only the maximum term in the sum, e^{-k_0} , is required. Hence $p = 1/\lambda$, $q = 1$, allowing us to conclude, as in Section 3.1, that

$$J_0(i) \sim \ln \lambda / \lambda, \quad \text{Var}(J) \sim 1.$$

To check the adequacy of the predicted distribution, $P(J) = \exp(-J_0(i)/\lambda e^{-k_0})$, for random sequences, a large simulation was carried out, and the distribution of the function $F = P(J)$ plotted. Because $F = \sum_i P(\Delta_{\text{match}}(i))$, then we should have $P(\Delta_{\text{match}}) = P(\Delta_{\text{match}})F = P(\Delta_{\text{match}})P(\Delta_{\text{match}})/F$, or $P(\Delta_{\text{match}})/F = 1/F$. This was verified to high accuracy.

Mismatches are not the only imperfections we should expect. The insertion of bases – e.g., of three-base codons – or deletions may have only a marginal effect on the resulting protein coding or regulatory sequences. These are termed blurs. For two-sequence comparison, we can allow for either one by inserting blurs in either sequence when testing an oligonucleotide position such that just one blur is at the same site, so there was no reason to

include them in start with). In other words, although the linear parametric form is not necessarily optimal (Joumard et al., 1992), and in particular, the interior blocks of a continuous sequence of blocks will certainly "cost" less than the first and last blocks, we now set

$$\pi(\ell, \ell) = \pi(\ell, \ell) - \lambda \delta(\ell, \ell) - \mu h(\ell, \ell),$$

where λ is the number of Puspaioli's blocks in the alignment. Now, $\mu = \mu_0$ reduces us to the maximal case and $\lambda + \mu = \alpha$ is the permissable case, whereas $\lambda + \mu = 0$, taking neither about matches nor insertions, should produce very high scores, of the order of the sequence length, with a line of phase transitions in between the extremes. The $\lambda = \mu = 0$ situation, with the anticipated noninformative $E(S) \approx L$, is nonetheless one that we can say something about.

4.2.2. Penalty-Free Limit

Suppose $A = a_1 \dots a_n$ is a sequence and $B = b_1 \dots b_m$ is a sequence with $m \leq n$. If both matches and insertions are to be ignored, we should now define S as a subsequence of A if for some $0 \leq i_1 < \dots < i_m \leq n$, we have $a_{i_j} = b_j$. For sequences A and B , taken for convenience to having a common length n , we then define

$$\pi(A, B) = \max \text{ length of subsequence common to } A \text{ and } B,$$

identical with the $\lambda = \mu = 0$ score. We now look for

$$\lim_{n \rightarrow \infty} E[\pi(A, B)]/n,$$

assuming a 3-letter-alphabet with equal probability of unit selection (Cleivald and Bandelt, 1993; Delam, 1997). Oddly enough, the inspection concludes that the limit is bounded from below in easy to show. We simply observe that the maximum occupation of each sequence by some unit must exceed the average, which is $n/3$, and tag each sequence by its dominant unit, say a . Clearly any two a -sequences have a common subsequence of length $\pi(a)$. If K_a is the number of a -tagged sequences, then there are K_a^2 pairs with $\pi(a, a) > \pi(a)$. Hence $E[\pi(A, B)]/n \geq \sum_{a=1}^3 K_a^2 (n^2 - \pi(a))$, where $N = \sum K_a$. However, K^2 is convex, so that $\sum K_a^2 \geq N^2 / 3$. We conclude that

$$E[\pi(A, B)]/n \geq \frac{N}{3} - 1$$

of course, it is possible to do better.

It is harder to show that the limit is bounded from above by something less than 1. To do so, we define $P(m, k, l)$ as the number of length m sequences containing a fixed sequence S of length n . We claim that

$$P(m, k, l) = \sum_{j=0}^m \binom{m}{j} (k - 1)^{m-j}.$$

This is certainly true for $n = 1$, and so we shall carry out induction on n . Note first that the assertion is also true for $m = n$ and $k = m + 1$ (the number containing a given unit is the total number minus the total number lacking that unit, i.e., $(k - 1)^m - 1^m = (k + m - 1)^m - k^m = l^m$, which agrees with the $n = 1$ case). Then if $1 < m < n$ and $S \subset A$ is a subsequence, we define $A' = a_1, \dots, a_{n-1}$, $E = E_1, \dots, E_{n-1}$. The set $\{A\}$ is divided into (1) those sites for which $a_i = a_m$ have $S \subset A'_i$, in which case $P(m - 1, m, l)$ possibilities for A'_i , and (2) sites for which $a_i \neq a_m$ have $S \subset A'_i$, and a_i is free to have $k - 1$ different values, a total of $(k - 1)P(m - 1, m, l)$ possibilities. Hence

$$P(m, m, l) = P(m - 1, m - 1, l) + (k - 1)P(m - 1, m, l),$$

which is indeed satisfied by the assumed formula.

Actually, the expression for $P(m, k, l)$ is steeper than we need. Note that if $j \geq m/k$, then $\binom{m}{j} (k - 1)^{m-j} \leq \binom{m}{j} (k - 1)^{m-k}$; it follows that $\binom{m}{j} (k - 1)^{m-j} \leq \binom{m}{j} (k - 1)^{m-k}$ for $j \geq m/k$, from which

$$P(m, m, l) \leq m \binom{m}{m} (k - 1)^{m-m} = m! (k - 1)^m, \quad \text{when } m \geq n/k.$$

The next step, in defining $\phi = m/n = 1/k_1$, is to show that for large n , the proportion $A_1^{(n)}(k_1)$ of pairs (A, B) of lengths n with $n(A_1, B_1) \geq k_1$ is bounded by

$$A_1^{(n)}(k_1) \leq M_1(n)^{k_1}$$

where

$$M_1(n) = \frac{n^{n/2-1}(k_1 - 1)^{k_1}}{\binom{n}{k_1} - (k_1)^{k_1-1}}.$$

We do this by "overcounting," i.e., ignoring the fact that the same pair may have more than one common subsequence and just using subsequences to build up pairs. In other words, if $p(n, m, k)$ is the number of pairs (A, B) with $n(A, B) \geq m$ and $G(n, m, k)$ is the number of triples (A, B, S) with A and B of lengths n , S of length m , then certainly $p(n, m, k) \leq G(n, m, k)$. It

follows that $\partial P/\partial t = \pi_1(\lambda)$,

$$\begin{aligned}\frac{\partial \ln \frac{P(x, m, d)}{P(x)}}{\partial t} &\leq \frac{G(x, m, d)}{d^2} \\ &= \sum_k \frac{P(x, k, d)^2}{d^2} \leq d^{m-2} \left[\pi \binom{n}{dm} (k - 1)^{m-2} \right]^2 \\ &= \left[d^{m/2-1} (k - 1)^{m-2} \left[\pi \binom{n}{dm} \right]^{1/m} \right]^{2m}\end{aligned}$$

which, by virtue of Stirling's approximation, can be shown to imply the preceding equation.

Now we use this result. First, by examining $\partial^2 \ln \frac{P(x, m, d)}{P(x)}$, we find that the equation $P_0(d) = 1$ has a unique solution $d = T_1$ in the interval $[1/M, 1]$ and that $M_0(d) < 1$ for $d > T_1$. Thus, for any d , $T_1 < d < 1$, we divide the d^{2m} pairs of sequences into two categories: (1) those for which $k \in \pi(A, d)$; (2) $\in \bar{\pi}_1$, and (3) those for which $k \in \pi(A, d)$, $k \notin \pi$. With the above definition of $d_0^{(2)}(d)$ then $\partial[\pi(A, d)] = \text{Re}[1 - d_0^{(2)}(d)] + i\text{Im}[d_0^{(2)}(d)] \leq \pi[d + M_0(d)]^2$; but $M_0(d) < 1$, and we can let d approach T_1 . Hence, $\lim_{d \rightarrow T_1} \partial[\pi(A, d)]/d \leq d \rightarrow T_1$, the desired linear upper bound.

4.2.3. Effect of Insertion Probability

Let us return to the more informative and controllable score for sequences with gaps, $n = r - k; k = p/d$. The needed score distribution is not known exactly, even for the Markov models that we have been considering, but has been approximated in numerous occasions (see, e.g., Zhang and Marr, 1993; Mori and Truta, 1999; Banavar et al., 2001). Evolutionary models to assess the confidence with which biological similarity can be inferred from such scores have appeared as well (Hsu and Lusby, 1996), and this is obviously a more realistic direction to take. Here, however, we will attend to a quick and dirty extension of the argument of Subsection 4.2.1 to the situation of gapped alignment of two thoroughly random sequences. The subsequences being compared now constitute a sequence of matches, mismatches, and unpaired vacancies with respective probabilities p , q , and r —now $p + q + r = 1$ —and rather than maximize over $m = r - k + d$, we will simply sum:

$$P(k|m) \geq M_0(d) \sum_{r+k+d \leq n} \frac{r+k+d}{r!k!M!} p^r q^k r^d.$$

We conveniently the inequality restriction by inserting a unit step (Heaviside) function $H(x) := \int_x^\infty e^{-t}/dt$ and, where ν is a real free variable (defined below)

the origin in the complex plane. Hence

$$\begin{aligned} \Pr[\text{align}] \geq S(\mu) &= \int_{\mathbb{C}} e^{2\pi i \lambda \operatorname{Im} z} \frac{(p+q+z)}{z^{\alpha+\beta+1}} |p^z q^{\bar{z}} z^{\bar{z}}|^2 dz/(2\pi) \\ &= \int_{\mathbb{C}} d^2z \int (1 - \mu z^{\alpha\bar{z}} - \bar{q} z^{\beta\bar{z}} - p z^{\bar{\alpha}\bar{z}}) dz/dz, \end{aligned}$$

where we have used the identity $\sum_{n_1, \dots, n_\ell} (a_1 + \dots + a_\ell)(b_1, \dots, b_\ell) = a_1^{\alpha_1} \dots a_\ell^{\alpha_\ell} = 1 = a_1 \dots a_\ell z^{\alpha z}$.

It follows by a standard argument that if $p = q^2$ is the root of

$$1 - p(z) - q(z^2) - \bar{q}(z^2) = 0$$

of maximum amplitude, then $\lim_{\mu \rightarrow 0} \Pr[\text{align}] \leq S(p) = \mu$. For μ close to 0, we expand each $p^z = \exp(-\log p)$ in several orders in $\log p$, and use $p + q + \bar{q} = 1$ to obtain

$$\log p = \frac{2}{p + \lambda/q + \bar{q}/p} (p - \lambda q - \bar{q}p).$$

If λq , $\bar{q}p$, and p are of the same order of magnitude and small, we can write this as

$$-\log p = 2 \left(\lambda \frac{q}{p} + \bar{q} \frac{p}{q} - 1 \right),$$

which indeed reduces properly to the previously analyzed $p = O(\mu q)$. When $p \neq 0$, we see that the general structure of the distribution is unchanged, except for the replacement $\lambda \rightarrow \lambda + p/(q/\bar{q}p)$ in the determination of p . For λ and/or μ not small, the original equation is required, and p asymptotic to μ . Also when, at small λ, μ , we have

$$2q + pq = \mu,$$

p reaches 1, and the linear regime takes over. To be sure, here, as in the case of the gap-free score, the full distribution function depends as well on the slowly varying $O(\delta)$ and the number of relevant pairs (I, J) .

4.2.4. Score derivatives

The maximum-score criterion for comparison of sequences A and B is useful provided that the score $s(I^*, J^*)$, $I^* \subset A$, $J^* \subset B$, can be quickly computed and that alignments can be generated rapidly. At the very least, we should be able to verify that a score is locally optimal, so that $s(I^*, J^*) \geq s(I^*, J^*)$ for (I^*, J^*) in a suitably defined neighbourhood of (I^*, J^*) . All of this depends very much on the precise nature of the score that is used. We have seen

that, roughly speaking, scores may have a periodic linear alteration with respect to sequence length ℓ that goes between $m\ell$ and $m'\ell$. Different considerations apply to the two extremes.

Let us first focus on the c regime, where the constant c has been estimated (Buckheit and Haubrille, 1981) for more agreeable letters, k , in the alphabet – and in the case of mismatch and match penalty-free criteria – as $2/3k^2$. In fact, the matching criterion can be weakened even further if the two subsequences are allowed to be permuted arbitrarily for the purpose of the test. Hence the maximal subsequence for a sequence A containing w_1^A, \dots, w_k^A letters of types $1, \dots, k$, and similarly for B , is represented as $w_1 = \min(w_1^A, w_1^B), \dots, w_k = \min(w_k^A, w_k^B)$. We can certainly regard $\sum_i w_i$ as the length of this match, but how do we penalize the remaining errors, $(w_1^A - w_1^B), \dots, (w_k^A - w_k^B)$ in number, that do not match with anything? Such a penalty should be nonnegative and vanish only for the identical letter content of the two subsequences. Perhaps the simplest is just a weighted rms form

$$\begin{aligned} s(A, B)^2 &= \sum_{j=1}^k w_j (w_j^A - w_j^B)^2, \\ (w_j) &= \frac{1}{2} (w_j^A + w_j^B), \end{aligned}$$

where $w_j = 1/k(j)$ would produce a χ^2 form.

If the letters are bases, the aggregate information contained by the above $s(A, B)$ – just that of base frequencies – is not a very incisive measure of sequence similarity. However, in the work of Tonney et al. (1990), the subis are taken as the complete set of $k = 4^n$ words of n nucleotides, or even those between any α and any β in length. The predictive power of this score is claimed to be very strong.

Most of the scores traditionally used have been in the last category (Bilgram, 1985; 1987) that these scores tend to rely on the consensus form:

$$\begin{aligned} s(\bar{A}, \bar{B}) &= \max_{1 \leq i \leq k} s(\bar{i}, \bar{j}), \\ s(\bar{i}, \bar{j}) &= \rho s(\bar{i}, \bar{j}) - d(\bar{i}, \bar{j}), \\ d(\bar{i}, \bar{j}) &= \min_{\lambda} d(\bar{i}, \bar{j}, \lambda), \end{aligned}$$

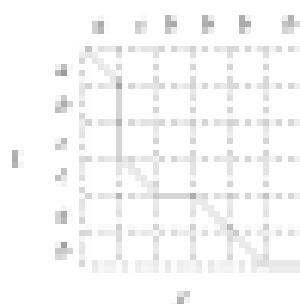
where $s(\bar{i}, \bar{j})$ measures the length of the (\bar{i}, \bar{j}) pair (e.g., the mean length of i and j), $d(\bar{i}, \bar{j}, \lambda)$ is a dissimilarity function depending on which alignment λ is being considered, and ρ is a relative weight. These can be much more strict in equidurating a match, but do have the option of gap insertions in either

sequence, like those previously referred to. How should we organize the possible alignments? A convenient representation, similar in concept to that of Subsection 3.2.1, for the allowed gap-inserted subsequences, goes like this. We consider, for example, two subsequences, $\bar{I} = \text{aibccb}$ and $\bar{J} = \text{acbbcc}$, aligned as

$$\begin{array}{ccccccc} \bar{I} & = & \text{a} & \text{b} & \text{c} & \text{b} & \text{c} \\ & & \downarrow & & \downarrow & & \downarrow \\ \bar{J} & = & \text{a} & - & - & \text{c} & \text{b} & \text{b} & \text{c} \end{array}$$

formally we have a vertical line for a gap in sequence \bar{J} , horizontal line over in sequence \bar{I} , and diagonal otherwise; each diagonal line can be a match or a mismatch. Below sets

$$d(\bar{I}, \bar{J}, \Delta) = \# \text{mismatches} + 2 \times \# \text{gaps}.$$



(the first corresponding to mismatches, the second to evolutionary deletions), which is a special case of the Smith-Waterman form of Subsection 3.2.1; in the above alignment, $d(\bar{I}, \bar{J}, \Delta) = 1 + (2 + 2) = 5$.

With this representation, there is a 1:1 correspondence between $(\bar{I}, \bar{J}, \Delta)$ alignments and paths (of south, southeast, and east links) from upper left to lower right. There are a number of algorithms for obtaining the dissimilarity $d(\bar{I}, \bar{J}) = \min_{\Delta} d(\bar{I}, \bar{J}, \Delta)$. The simplest is to follow a neighborhood of Δ in any intersecting path, because $d(\bar{I}, \bar{J}, \Delta)$ (addition of subpaths, comparison) is easily made, but may result in only a local minimum. There are, however, dynamical programming methods that solve recursion relations for the scores of truncated subsequences. The Needleman-Wunsh (1970) prototype takes the form:

$$\begin{aligned} d'(\bar{I}', \bar{J}') &= \min \{ d'(\bar{I}' - 1, \bar{J}) + \Delta(a_1, -), d'(\bar{I}' - 1, \bar{J}' - 1) + \Delta(a_1, b_1'), \\ &\quad d'(\bar{I}', \bar{J}' - 1) + \Delta(-, b_1') \}, \end{aligned}$$

where \bar{I}' and \bar{J}' are truncated to their first i' and first j' elements, respectively

($\beta = 0$ for negative argument), and A is the penalty for the pair or link mismatch. There is as well a more recent technique (Zhang and Marr, 1993) that uses statistical mechanical methods to compute directly a weighted average of $s(\mathcal{A}, \mathcal{B}, A)$ over all A .

We now have a global assessment of the match between \mathcal{A} and \mathcal{B} , the remaining problem is to determine the optimal pair $(\mathcal{A}', \mathcal{B}')$ of $(\mathcal{A}, \mathcal{B})$. Rather than tally all subsequence pairs $\mathcal{A}', \mathcal{B}'$, we can at least accept $(\mathcal{A}', \mathcal{B}')$ only if the optimal alignment function satisfies $s(\mathcal{A}', \mathcal{B}', A_{\text{opt}}) = p(\mathcal{A}', \mathcal{B}') - d(\mathcal{A}', \mathcal{B}') > 0$; then $\rho > d/\lambda$, the mismatch density. Following this, accept only if either (1) the gap separation $d_{\mathcal{A}'} := A_{\mathcal{A}'} + A'_{\mathcal{B}'} - d(\mathcal{A}', \mathcal{B}') \leq g$ as well, or (2) if $A_{\mathcal{A}'}$ is overlapped by $A_{\mathcal{B}'}$, that $p_{\mathcal{A}'} - d_{\mathcal{A}'} \geq p_{\mathcal{B}'} - d_{\mathcal{B}'}$. In fact, a large database can be compared by a modified dynamical programming routine. If $s(\mathcal{A}', \mathcal{B}')$ is the maximum for sequences ending at $\mathcal{A}', \mathcal{B}'$, it takes the form

$$\begin{aligned} s(\mathcal{A}', \mathcal{B}') &= \max\{p_{\mathcal{A}'}(i-1, j-1) + d_{\mathcal{A}'}(i, j_1), p_{\mathcal{B}'}(i-1, j) - d_{\mathcal{B}'}(i_1, j)\}, \\ &\quad + s(\mathcal{A}', j-1) - d(\mathcal{A}', \mathcal{B}'). \end{aligned}$$

However, agreement on the parameters to be used in A is not universal.

4.3. Multisequence Comparison

We continue with the question of how to locate and/or characterize sequences with functional similarity measured by structural similarity, but now at the multi-sequence level.

4.3.1. Locating a Common Pattern

Suppose that we have a collection of sequences associated with a common function, such as containing a binding domain for a protein-activating transcription. How do we characterize this domain, which may very well consist of diverse subunits—conserved subpieces—but not necessarily with invariant spacings? Presumably there will be an optimal subset of each sequence with respect to the others for lining up the subpieces, but the number of possible off-together rearrangements is astronomical. If we have some idea of the length ℓ of the subdomains involved, there are shortest inserted as follows (Stormo and Hartwell, 1989). We take each sequence, say of common length L , and discompose it into $L/A + 1 - \ell$ nonoverlapping constituent words of length ℓ . Now we start with sequence #1, split into $L/A + 1 - \ell$ words. We take each word of sequence #2, similarly discomposed, and test it for maximum match among words of #1, applied it to its optimal partner in #1, and, in case of a

is, keep both pairs. Then we test the words of \mathcal{B}^1 against the set of pairs and append to the closest pair, and so forth, until the sequences are exhausted.

However, what do we mean by the best match of a word to an n -tuple of words? A standard assessment is in terms of "information content." If p_i is the relative frequency of base i in the genome and p_{ij} that of base i at location j in the whole set of inserted bacterial fragments, then

$$I(j)(p) = \sum_{i=1}^n p_{ij} \ln(p_{ij}/p_i)$$

is the information content of the set of locations j . For a word $j = 1, \dots, R$, we correspondingly set $J = \sum_{i=1}^R J_{ij}$. In the preceding procedure, then, the best matching new word is that which gives the highest information content when combined with the previous set. After all sequences are combined, we end up with a set of $(L+1) - k$ weight matrices (see also Becher, 1990), each being a $d \times d$ array of J_{ij} , the most informative of which, as in the figure, signifies a dinucleotide motif. This scheme applies of course to proteins as well as to DNAs, but the controlling amino acid interactions depend more on physical category than specific identity (Miyazawa and Janish, 1985), so that the "threading" techniques (Lathrop and Smith, 1990) need to allocate sequence to structure in function of the hydrophilicity or hydrophobicity classes of amino acids. Note that, if $J_{ij} = p_i$, then a Taylor expansion yields



$$I_j = \sum_{i=1}^n \frac{1}{Z} \frac{(p_{ij} - p_i)^2}{p_i},$$

related to p^2 . Note too that

$$\begin{aligned} d(f - g) &= I\left(\frac{f}{p}\right) + I\left(\frac{g}{p}\right) \\ &= \sum_i (p_i - p_i) \ln\left(\frac{p_i}{p_i}\right) \end{aligned}$$

is a metric on probability distribution functions, always ≥ 0 and only $= 0$ for identical distributions.

We can generalize in another way. Suppose that $\{A_i\}$ is a grid-like set of features asserted to characterize (by presence or absence), e.g., the existence of a type of binding domain, and that these occur at relative frequency $f(A_i)$ in the sample. In "random," modeled by a parameter set $\{\beta_k\}$, the model frequency would be $(p_{\text{ref}}/\lambda)^k$. Is $f(A_i)$ a significant departure from randomness? In the above, A_i would be taken as a configuration of all bases at all positions in the length k window. If we imagine the data as the result of $N \rightarrow \infty$ trials with independent selection of the $\{A_i\}$, then the probability of $f(A_i) = N_i^k/\lambda^k$ would be $P_k = (\lambda^k N! / k! N_i^k) \prod_i p_{A_i}^{N_i} / \lambda^k$. The normalized (negative) log likelihood:

$$\ell = \lim_{N \rightarrow \infty} -\frac{1}{N} \ln P_k,$$

which, by means of Stirling's approximation, works out to

$$\ell = \sum_i f(A_i) \ln(p_{\text{ref}}/f(A_i)/p_{A_i}/\lambda),$$

would then be a legitimate measure of the deviation from randomness. Note that if A_i is decomposed into mutually independent A_j (say the base at site j), so that $p_{A_i}(A) = \prod_j p_{A_j}(A_j)$, and only $f(A) = \prod_j f(A_j)$ is computed, then we have instead

$$\ell = \sum_i \sum_j f(A_j) \ln(f(A_j)/p_{A_j}/\lambda),$$

which was used in the above.

4.3.2. Assessing Significance

If there is expected to be a consensus motif, appearing in various encrypted forms, but short enough – a minimal functional unit – that insertions and deletions need not be considered, this should be picked up by similarity in – analog of register – multiple occupancy. Because we are looking at a sequence of locations, it is really the sign feature of an unlikely "run" that is being assessed, and there is no reason why the pair-matching technology should not be extended to this case. Again, let us take the random reference as an independent choice of letters at each site, with probability p_{ref} for type a . Suppose that we have obtained an optimally aligned sample of n sequences. The probability of a motif that mismatch need not be negligible, in the sense that it is a match at a given location in $y[k]$ of the n sequences containing the same letter. If we observe a contiguous r -site match in the entire just described extension by L repeated r matches in object y , then a first estimate of significance would

In a comparison with $A[t_{\text{last}}]$, after which we can then have by comparison with the distribution of T_{last} ,

The question then is that of $E(T_{\text{last}})$ over all (out-of-register) comparisons of sequences of length $k = l_1$. Proceeding exactly as we did in the two-sequence comparison, we find that the expected number of letters with at least k out of k matching letters would be the same as permutations (i.e., when j are the same)

$$\begin{aligned} & \sum_{j=0}^{k-1} \sum_{\mu_1, \dots, \mu_k} \sum_{\nu_1, \dots, \nu_k} P(\mu_1 \mu_2 \dots \mu_k | \nu_1 \nu_2 \dots \nu_k) \dots P(\mu_j \mu_{j+1} \dots \mu_k | \nu_1 \nu_2 \dots \nu_k) \\ &= \sum_{j=0}^{k-1} \sum_{\mu_1, \dots, \mu_k} \binom{k}{j} P(\mu_1 \mu_2 \dots \mu_k | \nu_1 \nu_2 \dots \nu_k) \\ &= p l_1 \dots l_k. \end{aligned}$$

Again, a run starts after a failure at any of μ_1, \dots, μ_k multisequence points, which tells us that

$$E(T_{\text{last}}) = \log_2(p l_1 \dots l_k) + \dots$$

where

$$p = \sum_{j=0}^{k-1} \sum_{\mu_1, \dots, \mu_k} \binom{k}{j} P(\mu_1 \mu_2 \dots \mu_k | \nu_1 \nu_2 \dots \nu_k).$$

The technical problem of computing p for large k is not trivial and is often solved by large-deviation theory. However, we can be quite direct. The quantity we need is

$$f(x) = \sum_{j=0}^k \binom{k}{j} x^j (1-x)^{k-j},$$

where $0 < x < 1$ is the region of interest, so that a normal approximation is unsuitable. A Poisson approximation can be used, but still better, let us put $f(x)$ in integral form, so that it is completely controllable. To do so, we simply observe that, from $\binom{k}{j} (m - j) = \binom{k}{j+1} (k + 1)$, then

$$\begin{aligned} f(x) &= \sum_{j=0}^k \binom{k}{j} (x^{j+1}) (1-x)^{k-j} = \sum_{j=0}^k \binom{k}{j} m - f(x^j) (1-x)^{k-j} \\ &= \sum_{j=0}^{k-1} \binom{k}{j} (x^{j+1}) (1-x)^{k-j} - \sum_{j=0}^{k-1} \binom{k}{j} (x^{j+1}) (1-x)^{k-j} \\ &= (x/k)(k-1)m - k(x^{k+1})(1-x)^{k-1}. \end{aligned}$$

Hence

$$\Gamma(x) = \frac{\pi^{\frac{1}{2}}}{\theta - 1(x - 1)} \int_0^\infty y^{x-1} (1 - y)^{\theta-1} dy$$

in terms of the incomplete beta function.

Because the integrand is maximum at $y_0 = 1 - 1/\theta = 1$, but $x < y_0$, the integrand rises rapidly as y approaches x , so we can approximate

$$\begin{aligned} & \int_0^x y^{x-1} (1 - y)^{\theta-1} dy \\ &= \int_0^x (x - y)^{x-1} (1 - x + y)^{\theta-1} dy \\ &= x^{x-1} (1 - x)^{\theta-1} \int_0^x \left(1 - \frac{y}{x}\right)^{x-1} \left(1 + \frac{y}{1-x}\right)^{\theta-1} dy \\ &= x^{x-1} (1 - x)^{\theta-1} \int_0^x \exp\left[-\left((1-y)\frac{x}{x} + (x-y)\frac{1}{1-x}\right)\right] dy \\ &= x^{x-1} (1 - x)^{\theta-1} e^{-1/(x(1-x))}. \end{aligned}$$

We conclude that p_x has the very compatible form

$$p_x = \left[\sum_{k=0}^{\infty} \frac{p_k^{x-1} q_k^{\theta-1-k}}{x - 1 - p_k(x - 1)} \right] B\left(\frac{x}{x}, \frac{\theta}{x}\right)$$

by which, e.g., the Poisson approximation, can also be applied.

Assignment 7

In a model of random evolution, an organism is defined by a string of T composite sites, each one of which can assume one of two forms, labeled by 0 or 1. Transitions from 0 to 1 or 1 to 0 occur at a rate of ν per site per generation, from an initial proportion. After T generations, two numbers of the population are compared, according to the score

$$s = r - k d,$$

where r is the number of matches in homologous site structures of the two numbers, k is the number of mismatched, and d_{\max} is the maximum of r over all s .

1. What is the probability that two corresponding sites are the same? Different?
2. Plot the distribution of s at fixed n , the distribution of d_{\max} .

3. Show that the dependence of $P(D_{\text{max}})$ on t changes qualitatively when $b = \text{const} \cdot t^{1/2}$.

4.3.3. Category Analysis

In attempting to understand the language of DNA, we know that there are broad categories of function that constitute the basic structure. These are best in the form of general instructions: splice out the intron that substitutes, and transmethylate another histone, . . . , and translated proportion. This is, in fact, this will be involved in phosphorylation, . . . In numerous cases, we have a collection of subsequences known to have a common characteristic, and the objective is to find out in what fashion this can be used either as a common characteristic of the base sequence involved, or then the property could be readily identified from sequence data alone.

An important example is that of splicing signals. Introns are removed from precursor-RNA by RNA splicing. First the pre-RNA is cleaved at a 5'-splice site of the intron, separating the sequence into \cdots exon + intron + exon \cdots , and then cleaved at a 3'-splice site of the intron, with the intron "ariat" removed by a spliceosome of small ribosomal proteins. A great deal of study in higher eukaryotes has led to the conclusion (Mount, 1982) that *consensus subsequences*:

$$\begin{aligned} S' - \text{exon} &= \binom{C}{A} \text{ACGTC} \binom{A}{G} \text{AGT} - \text{intron} - R', \\ S - \text{intron} &= \binom{U}{C} \text{UAG} N \binom{C}{U} \text{AGUC} - \text{exon} - S' \end{aligned}$$

are involved; here $\binom{C}{A}$ indicates that either C or A is required, N stands for any base, and the cut takes place at R . The problem is to decide what combinations of the suitable bases indeed give a splice command. A traditional technique is to choose the characteristic to being tested like, say identification as a splice signal, construct the corresponding weight matrix ($d_{i,j}$) of relative frequencies of bases i at site j for a large sample of cases that have been found and then repeat

$$W = \sum d_{i,j}(x)$$

as the score for the sequence being tested, which must exceed some threshold to be evaluated as a "splice." Here, again, $d_{i,j}(x) = 1$ if base i is at site j , else 0. There are a number of intelligent modifications that have been made to accord with presumed composite structure. For example (Shapiro and Burzynski, 1982) in the 5'-splice problem, there are 11 sites in the left of 5'

and δ to the right that appear to be significant. Let b_1 be the highest sum of the highest 8 weights out of 10 in the "training" sample, b_2 the lowest sum of the lowest 8 weights, b_3 the highest sum of 4 weights, b_4 the lowest. Then if t_1 is the highest weight of 8 out of 10 in the test sequence, t_2 the sum of the 4 weights, the score

$$W = \frac{b_1 - b_2}{b_1 + b_2} + \frac{2(t_1 - b_3)}{2(t_1 - b_4)}$$

has empirically proven quite effective.

The above is more than a little bit ad hoc. What we are really trying to do is to separate the sample "points", i.e., subspace configurations, into two clusters or categories, so that a test "point" can be assigned as to which cluster it is closest. More generally, the nonweighted cluster problem is to separate n distinguishable interclusterers, which can be translated into the following dual tasks:

$$J(w, c) = \sum_{i=1}^n \sum_{a=1}^c w_{ia} d_{ia}^2$$

subject to

$$\sum_{a=1}^c w_{ia} = 1,$$

where $w_{ia} = 1$ if point i is assigned to cluster a and d_{ia} is the Euclidean distance between point i and the centroid of cluster a . This can be solved by "steepest-ascent" minimization techniques.

An alternative approach is by means of classical linear discriminant analysis. Pielot and Tong (1992) looked at the problem of distinguishing between coding (exons) and noncoding segments (introns) of DNA. Among various suggested criteria, they found that dioder usage – the frequency of the 4000 different hexamers in a segment – was as effective as any more sophisticated criterions. The formulation produces a coefficient vector such that for a window characterization vector f , the window is "coding" if $c \cdot f > r$ for some threshold r . Here, we determine c by maximizing the ratio of the between-population variation of $c \cdot f$ to the within-population variation of $c \cdot f$. Specifically, suppose that $f_{w,j}$ is the dioder $-j$ frequency for the w^{th} window element of class c (coding or noncoding), $\bar{f}_{w,j}$ is the mean over w , and \bar{f}_j is the mean over w and c . Then the total covariance matrix V is given by

$$T_{jk} = \sum_{w=1}^W (f_{w,j} - \bar{f}_j)(f_{w,k} - \bar{f}_k)$$

the within-covariance \mathbf{W} by

$$\mathbf{W}_{ij} := \sum_{k \in \Omega} M_{ijk} = f_{ij}(\mathbf{R}_{\text{test}} - \bar{\mathbf{R}}_{\text{obs}})$$

and the between-covariance by $\mathbf{B} = \mathbf{T} - \mathbf{W}$. Maximizing the ratio of variances $\sigma_w/\sqrt{\lambda_{\text{min}}(\mathbf{W})}$ means that σ belongs to the maximal eigenvalue λ_1 of \mathbf{B} with weight $\mathbf{W} = \lambda \mathbf{B} = \lambda \mathbf{W}_{\text{obs}}$, or

$$\mathbf{B}^{-1} \mathbf{B} \sigma = \lambda \sigma,$$

and the threshold r is generally determined so that the fraction of errors in coding windows equals the fraction in noncoding windows.

Computationally, the evaluation of \mathbf{B}^{-1} can be a problem, especially because no knowledge of information makes \mathbf{B} nearly singular. We can avoid this by replacing \mathbf{B} with its diagonal part \mathbf{W}_{B} , resulting in what is equivalent to the Pearson discriminant. It is this version that was used, resulting in a 90% accuracy in the above study.

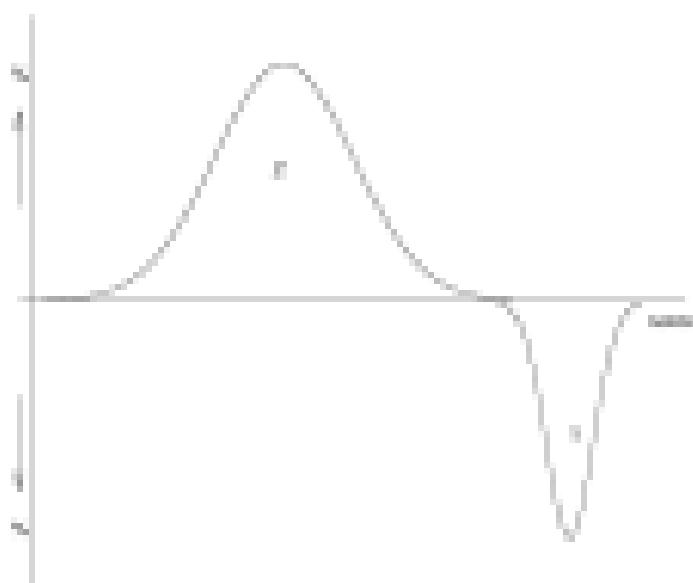
A similar technique has been used (Hayashi, 1992; Iida, 1997, 1998) for splice signal identification. Here the 3-splice signature was taken as the total nucleotide stretch corresponding to the constraint $\binom{r}{2} \geq N_1^2 N_2^2 / 4(1 + \epsilon)$. The training sample was described by $\delta_j^{(n+1)}(x) = \binom{r}{2}$ if base x is at least present at site j in the n^{th} base from a reference N_1 splice ($n = 1$) and N_2 complete ($n = 2$) sequence. A linear rule

$$g^{(n+1)} = \sum_{j \in \Omega} C_{j,n} \delta_j^{(n+1)}(x)$$

was set up for the term (y, n) , and the problem was expressed as that of finding the 34 weights, not simply as a frequency-weight matrix, but so that the $n = 1$ scores are clustered as far from the $n = 2$ scores as possible. Again, the traditional criterion was used. We define

$$\begin{aligned} g^{(1)} &:= \frac{1}{K_1} \sum_{n=1}^{N_1} g^{(n+1)}, \quad \beta := \frac{1}{N_1} \sum_{n=1}^{N_1} g^{(n+1)}, \\ \sigma_{\text{tot}}^2 &:= \frac{1}{N_1} \sum_{n=1}^{N_1} (g^{(n+1)} - \beta)^2, \\ \sigma_{\text{err}}^2 &:= \frac{1}{N_1} \sum_{n=1}^{N_1} N_1 (D^{(n+1)} - \beta)^2 \\ &= (K_1 K_2 / N^2) (\beta^{(0)} - \beta)^2, \end{aligned}$$

and maximize $\sigma_{\text{err}}^2 / \sigma_{\text{tot}}^2$, i.e., we maximize the distance between centers of mass relative to the total standard deviation. The $\langle C_{j,n} \rangle$ evaluated in this way of course tell us which sites are important and which contribute mainly



to the noise. The corresponding scores, $S = \sum_{j,k} C_j(\mu)(\delta_j(k))$, then applied to the training data, show a very good separation between type 1 and type 2, yielding greater than 90% selection accuracy, and have been used to predict new splice sites as well.

4.3.6. Adaptive Techniques

(Bayesian Analysis)

There is a more direct utilization of category [see, e.g., Pfeifer et al., 1991; Balaji and Breslow, 1994]. Suppose that the datum to be classified consists of a sequence $Z = \{z_i\}$ of bases, or of codons, or of dicodons, or of amino acids, etc., each unit of which may be specified in various redundant ways by a bit sequence. For example, 3 bits for an amino acid or spacer (end of protein)/will suffice, but we can also use 21 bits, all of which vanish but one—a typical neural net specification, which was also used in the 3-bit example. If the categories are labeled by μ , Bayes's theorem tells us that the probability that a sequence Z being tested belongs to μ can be written as

$$P(\mu | Z) = P(Z | \mu)P(\mu) / P(Z)$$

where

$$P(Z) = \sum_{\mu} P(Z | \mu)P(\mu).$$

Here, $P(\mu)$ and $P(M \mid \mu)$ are in principle estimated by analysis of the prior training data, but this is not feasible if S is too large. A given Σ may never have occurred in the training data, making $P(\Sigma \mid \mu)$ impossible to estimate. The simplest way of avoiding this problem is by modelling $P(\Sigma \mid \mu)$. We can, for example, regard the units x of Σ as independent (but then these units should be composite, to afford more information), i.e.,

$$P(\Sigma \mid \mu) = \prod_x P(x \mid \mu),$$

where, of course,

$$P(x \mid \mu) = E[\delta(x_0) \mid \mu]$$

(the event that x_0 is all 0) in our previous notation.

If the units are not only independent but also identically distributed under change-of-location, with exchangeability criteria, then the above equation becomes:

$$P(\Sigma \mid \mu) = \prod_x P(x \mid \mu)^{N_x},$$

where N_x is the number of times that x occurs in $S = \{x\}$, so that

$$P(\mu \mid S) = \frac{P_{\text{prior}} \exp \sum_x N_x \ln P(x \mid \mu)}{\sum_{\mu'} P_{\text{prior}} \exp \sum_x N_x \ln P(x \mid \mu')}.$$

In particular, for dichotomic μ , i.e., true or false, we can write

$$\begin{aligned} P(T \mid S) &= \frac{P(T) \exp \sum_x N_x \ln p_x}{P(T) \exp \sum_x N_x \ln p_x + P(\bar{T}) \exp \sum_x N_x \ln q_x} \\ &= \frac{1}{1 + \exp \left(\sum_x T_x K_x + \delta \right)}, \end{aligned}$$

where $p_x = P(x \mid T)$, $q_x = 1 - p_x$, $T_x = \ln p_x / p_x$, and $\delta = \ln(P(T)/P(\bar{T}))$, all of which are estimable from the training set. This sigmoid form, δ for large argument, 1 for small – or the reverse – is a standard transformation from quantitative data to approximate yes or no.

A more sophisticated model uses a Markov chain:

$$\begin{aligned} P(\Sigma \mid T) &= P(x_1) P(x_2 \mid x_1) \cdots P(x_{n-1} \mid x_1 \cdots x_{n-2}) \\ &= P_{\text{prior}}(x_1) P_{\text{prior}}(x_2, x_1) \cdots P_{n-1,n}(x_n, x_{n-1}) P_{\text{prior}}(x_1) \cdots P_{n-1,n}(x_{n-1}) \\ &= P_{\text{prior}}(x_1) \cdots P_{\text{prior}}(x_{n-1}, x_{n-2}) \cdots P_{n-1,n}(x_n, x_{n-1}), \end{aligned}$$

where $p(x, t) = P(x \mid t) / P(x) P(t)$ is the weight or correlation coefficient. For

Identically distributed singletons and pairs, we have, by obvious notation,

$$P(X|T) = \prod_i p_i^{K_i} \prod_{ij} p_{ij} K_{ij},$$

leading, by equally obvious notation, to

$$P(X|B) = 1/(1 + \exp\left(\sum_i T_i K_i + \sum_{ij} T_{ij} K_{ij} + b\right)).$$

Neural Networks

There is certainly more significant information available in the independently but not identically distributed site model. Suppose that μ is not necessarily deterministic, and in conformity with computational realizations, let us instead imagine that we have adopted a binary coding representation, with the n probabilities at each composite site represented by a 1 in a substring of a bit, all of the others being 0. Then if $p_\mu(x) = \text{prob}(x_1 = 1 | \mu)$, we have

$$P(X|B) = \prod_j P_j(X_j)^{w_{j1}},$$

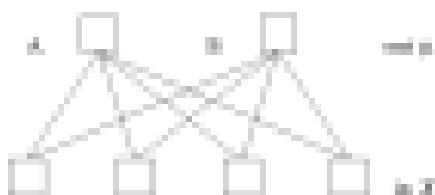
so that if $w_{j1} = \ln p_\mu(x)$ and $b_{j1} = \ln P(x)$, we can write

$$P(x|B) = e^{b_{j1}} \int \sum_i e^{w_{j1} x_i},$$

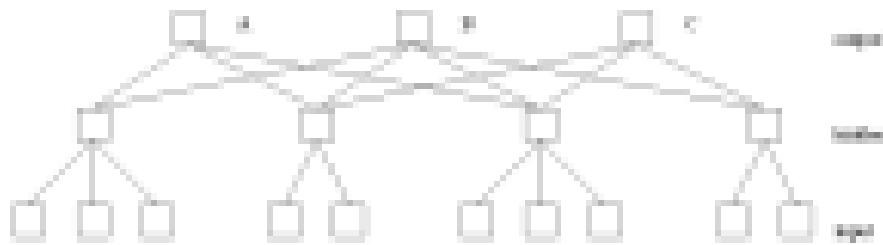
where

$$w_{j1} = \sum_i w_{j1i} x_i + b_{j1}.$$

The problem at issue is to use training information to determine the parameters w_{j1} , b_{j1} which can thereafter be used for testing mystery data. Once this is done, the μ to be allocated to the data would, e.g., be that which maximizes $P(x|B)$. Ideally, the maximum would be close to 1, and the others close to 0. In the real, we will have developed in this way a vector of functions $[L_i(X)]$ whose values will be the type μ probabilities associated with the input data X . The specific form, $P(x|B)$, given above is reasonable but unless overly simplifying assumptions be able to tolerate the unknown parameters in the function from limited training data. A class of representations functions, more complex (but usually general) and with similar components, is that termed "neural networks." Here, the unit is the multiple input-one output function $P(x|B)$, which by itself represents all networks of the two-layer "perception" architecture. If the units are stacked by "hidden layers"



between actual input and output, then it is termed a neural network. In either case, the probability at any level can be quantitated if any $P(\mu \mid \Sigma)$ is replaced with $\mu^k = \frac{1}{Z} e^{-\mu k} / \sum_{\mu'} e^{-\mu' k}$, which at $\mu \rightarrow \infty$ is instant 1 for the maximum μ , otherwise 0. In strictly binary representation, which is the norm, this is equivalent to $P(1 \mid \Sigma)$, $P(1 \mid \Sigma) = [1, 1, \dots, \exp(-\mu)]/Z = \mu_1/(1 + \mu)$, and would certainly be used for the output layer if the training data were used one piece at a time, as its value is then deterministic, not probabilistic.



There are then many internal parameters, which we can designate as the set $\{\mu_i\}$, that enter into the input-output function denoted by $P_{\mu}(\mu \mid \Sigma)$. In the process of "training" the network, we want to home in on optimal values of the parameters μ_i to best represent the input-output relationship. In practice, some feeling as to how the inputs should be clustered in a qualitative fashion means that many connections will not be used, i.e., the corresponding μ_i are not equal to 0 and do not appear in $P_{\mu}(\mu \mid \Sigma)$. Thus if training sets $\Sigma^{(j)}$ are entered, whose output characteristics, say $p_j^{(k)}$, are known, the μ_i are to be adjusted so that the errors between the $P_{\mu}(\mu^{(j)} \mid \Sigma^{(j)})$ and the $p_j^{(k)}$ are minimized. This depends on how we define the error. One definition is simply to sum the mean square error over the whole training set

$$E = \sum_{j=1}^J \|P_{\mu}(\mu^{(j)} \mid \Sigma^{(j)}) - p_j^{(k)}\|^2$$

and minimize to obtain the optimal μ . Minimization of a many-variable function is a slow art, the most primitive version of which ("conjugate gradient")

In one iteration in which a pattern \mathbf{y} is corrected by setting

$$\Delta \mathbf{w}_n = -\mu \Delta E / \Delta \mathbf{w}_n$$

If this is really a small change, then $\Delta E = \sum_i M_i(\mathbf{w}_n) \Delta w_{ni} = -\mu \sum_i M_i(\mathbf{w}_n)^2 < 0$, as claimed. If the uncorrected \mathbf{y}_n is the result of the previous training set and the correction is that which is due to inclusion of further sets, we speak of backpropagating the error to modify the \mathbf{y}_n . Other techniques, such as stochastic annealing, allow a fraction-of-the-time in E to well, to prevent any local-minimum traps.

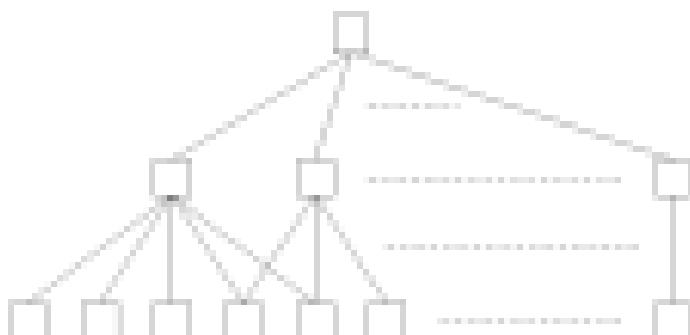
A simple modification often used is to balance the total weights of the categories in the training set, e.g.,

$$E = \sum_{\mathbf{y} \in S} \{ p_{\mathbf{y}}^{(0)} - P_{\mathbf{y}}(p_{\mathbf{y}}^{(0)}, \mathbf{w}) \}^2 / N_S$$

where N_S is the number of times that \mathbf{y} occurs in the training set. A rather different form maximizes the mutual information previously introduced. For this purpose, we first need the relative frequency $p_{\mathbf{y},i}$ of inserting a training pattern characterized by \mathbf{y} , and concluding that the pattern is \mathbf{y}_i for a given (\mathbf{w}_n) . The mutual information is then

$$M(\mathbf{y}) = \sum_{\mathbf{y} \in S} p_{\mathbf{y},i} \ln(p_{\mathbf{y},i}/P_{\mathbf{y}}),$$

Neural networks are not simple, they deal with a very small subset of parameterized output functions of the available input. However, if we have empirical, even anecdotal, information as to the important paths from input data to output characteristics, this can be incorporated into the neural network structure (consequently, the optimum set of $\{\mathbf{w}_n\}$ determined in the process may give a hint as to the biologically valid paths involved – this is certainly the



use for the organization of analogous linguistic data in the associated correspondence terms. For example, the prediction of an α -helix, a β sheet, or a coil secondary structure of proteins from the primary amino acid sequence, by use of the three-output network structure given above (Sokolov et al., 1992), does no better than 67%–77% accuracy whereas a corresponding highly structured network for predicting *A. castellanum* transcriptional promotion (Arenzana et al., 1991) has close to 100% accuracy in prediction and is very effective in prediction. Here, the 4 bases of CGT were given a 4-bit representation, and the significant information was that of 6 bases from the –10 region, 6 from the –26 region, which however could be separated from the –10 region by 15–21 bases. Thus, in the neural network, the 7 possible overlapping 6-base regions, together with the 6-base region at –10, were represented by 22 input units and connected to 8 hidden units, 1–24 to the first, 25–29 to the second, etc. The 8 were then connected to 3 output unit. In the perceptron version, all 12 were connected to the single output. Both variants had the high accuracy quoted above, showing how important assessment of relevant information becomes.

Applying Networks

In the feed-forward network design, the training of the network requires feedback as well, which we do "by hand" by setting the network parameters to minimize the difference between input and desired output. More complex processing of the input could be carried out if feedback existed in the network itself, creating what is termed a Hopfield network; see, e.g., Bodde and Dukens, 1990. The power of a Hopfield network is more evident in a somewhat different, but also somewhat related, application, in which the network, completely untrained, needs to make, starts out with a "memory" of many possible outputs. Given any input, it then proceeds in an unsupervised way to find the most closely related output.

Specifically, suppose we want to associate a set of patterns

$$\phi^d = \{a_i^d\}; \quad i = 1, \dots, M \}, \quad d = 1, \dots, D$$

where

$$a_i = \pm 1,$$

Then an Hebbian disease-dynamics might update the setting according to

$$w_j^d = \text{sign} \left(\sum_{i=1}^M x_{ij} a_i^d \right).$$

How should the multi-modal j competition matrix be determined? An obvious choice would be:

$$P_{ij} := \frac{1}{M} \sum_{k=1}^M \alpha_k^j \alpha_k^i$$

because if the patterns were orthogonal,

$$\sum_j \alpha_j^i(\pi) = N \delta_{ii}$$

then we would have $\sum_i P_{ij} \alpha_j^i = (1/M) \sum_{k=1}^M \alpha_k^i \alpha_k^j(\pi) = (N/M) \delta_{ij}$, and the updating would stop as soon as any pattern ρ_j^i is reached. In practice, a sizable fraction of M nonorthogonal patterns still serve as attractors of the dynamics.

5

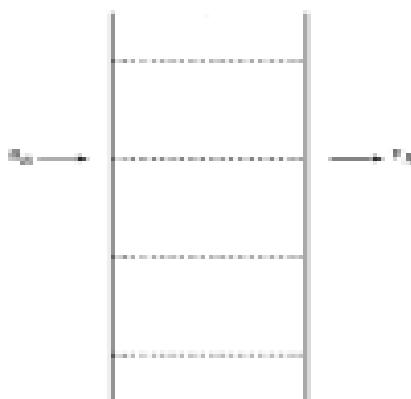
Spatial Structure and Dynamics of DNA.

3.1. Thermal Behavior

We have paid little attention to the way in which DNA transmits its information. There are many physicochemical steps involved, with a common theme of energy propagation and localization. Much activity has been devoted to the large-scale behavior of DNA, modeled, e.g., as a belt with elastic properties and giving rise to the "supercoiled" configurations responsible for placing linearly distant sections in close proximity spatially. The double strand under such motions is still "extended" (at a resolution of some 30 bp). On the other hand, the most exhalation of information transmission is that of transcription in RNA involving a "transcription bubble" or strand separation of only some 20 bp, activated by RNA polymerase contact. Of course, strand separation should be minor where it is required, and if we had DNA - as a metaphor for uniform energy transfer - the "walking" is stepwise and non-uniform, presumably a function of the local properties of the jiggly molecule. In fact, we might anticipate that such heterogenous dynamics would appear even in homogeneous DNA, as this would imply greater sensitivity to external stimuli.

Quite different aspects are involved in the dynamics of DNA, RNA, and proteins. The first minimally realistic model to be adopted along these lines was that of local denaturation (H-bond breaking) of DNA, a requirement for strand separation to allow for transcription (Watson, 1966; in Watson, said in effect that this work was the greatest discovery since the zipper). The model (Trotter et al., 1989; Peypoch and Bishop, 1989) is literally that of H-bond connections that lose their integrity when stretched too far by chain vibrations, taken here as thermally excited. As model energy we will choose the ultramodified

$$H = \sum_{n=1}^{N-1} \left[\frac{m}{2} (\dot{\theta}_n^2 + \dot{\psi}_n^2) + \frac{k}{2} [(r_{n+1} - r_{n-1})^2 + (r_n - r_{n-2})^2] + V(r_n - r_{n-1}) \right]$$



where

$$V(\alpha - \beta) = D[e^{-\alpha^2 - \beta^2/2} - 1]^2.$$

Only the out-of-phase $y_n = \frac{1}{\sqrt{2}}(u_n - v_n)$ motions stretch bonds, so we can eliminate the $x_n = (1/\sqrt{2})(u_n + v_n)$, retaining only the y energy:

$$H_T = \sum_n \left[\frac{m}{2} \dot{y}_n^2 + \frac{K}{2} (y_n - y_{n-1})^2 + D(e^{-\alpha y_n} - 1)^2 \right].$$

In this reduced model, in which one pays homage to the strong directionality of the hydrogen bond, the twist of the two-strand ladder, and indeed any motion transverse to y_n , becomes irrelevant – but the basic phenomenology is not disturbed.

A first step in the analysis is traditionally that of assuming that the system is in thermal equilibrium, and we accomplish this by asserting that, at reciprocal temperature β , the unnormalized probability of a configuration $(y_1, \dots, y_N, \dot{y}_1, \dots, \dot{y}_N)$ is given by the Boltzmann factor

$$\rho(y, \dot{y}) = e^{-\beta H_T(y, \dot{y})}.$$

The basic construct in statistical mechanics is the normalization factor, or *partition function*, subsequent to whose computation all of the system's properties are readily found. Here then, the partition function Z is

$$Z = \int \int e^{-\beta \frac{m}{2} \sum_i \dot{y}_i^2} dy_1 \cdots dy_N \\ \int \cdots \int e^{-\beta \frac{K}{2} (y_n - y_{n-1})^2} e^{-\beta D(e^{-\alpha y_n} - 1)^2} dy_1 \cdots dy_N.$$

The velocity integrations separate out, so if we are interested in positional

amplitudes alone, it suffices to consider

$$\mathcal{Z}^t = \int \dots \int e^{-\beta \left(\mu_1 - \mu_2 + k_B T_{\text{ext}} (x_1 - x_2) \right)} dx_1 \dots dx_N.$$

Of course, we have not really defined the problem until we supply information as to what happens at the end of the chain. If the chain is long enough, this should not affect any significant properties, so we will choose periodic boundary conditions. With coordinates x_1, x_2, \dots, x_N , but not mind $x_{N+1} = x_1$ is a mathematical convenience. Hence we can write

$$\mathcal{Z}^t = \int \dots \int T(x_1, x_2) T(x_2, x_3) \dots T(x_{N-1}, x_N) T(x_N, x_1) dx_1 \dots dx_N$$

where

$$T(x, x') = \exp \left[-\beta \left(\frac{k}{2} (x - x')^2 + D x^{(N-1)} - 1 \right)^2 \right].$$

Indeed the transition operator can be regarded as a matrix with continuous indices, so that

$$\mathcal{Z}^t = \text{Tr } T^N.$$

Now what are we looking for? Presumably, for a qualitative change in the distribution of visual superpositions as temperature is varied. However, the system can be regarded as a set of one-dimensional particles x_1, \dots, x_N under an external potential $V(x)$ and with next-neighbour interaction $\frac{1}{2} k (x_n - x_{n-1})^2$, and a theorem of van Hove asserts in that, even in the limit $N \rightarrow \infty$ needed to show a sharp phase transition, $(1/N) \ln \mathcal{Z}^t$ and associated physical quantities remain analytic in all parameters. Indeed, if we seek a state in which the x_i are localized — so that we can delocalize them by raising the temperature, this would require a trapped phase in equilibrium with a phase of zero vapor pressure, which certainly will not happen in one dimension. One way of enabling a qualitative change is by a tuning operation, imagining K as arbitrarily large, so that elasticity appears only on a large length scale. Without actually scaling everything so that the limit can be taken, let us see how this works out.

Leaving the potential $D(x^{(N-1)} - 1)^2 = \delta(x)$ unspecified for the moment, we need to carry out the operation

$$\begin{aligned} T(x) &= \int T(x, x') U(x') dx' \\ &= e^{-\beta V(x)} \int e^{-\beta k (x - x')^2} \rho(x') dx'. \end{aligned}$$

Because the Fourier transform $\int e^{ikx} \langle e^{-\beta H(t)} e^{-itf} f(x') \rangle dy' = \int e^{ikx} e^{-\beta H(t)} f(y') dy' = (2\pi/\beta K)^{1/2} e^{-\beta H(t)} \int e^{iky'} f(y') dy' = (2\pi/\beta K)^{1/2} \int e^{iky'} g(y') dy' = (2\pi/\beta K)^{1/2} f(e^{iky})$, we can write

$$Tf(y) = e^{-\beta H(t)} e^{iky} e^{-\beta H(t)} f(y) / (2\pi/\beta K)^{1/2}$$

so that

$$\begin{aligned} \mathbb{E}' &= (2\pi/\beta K)^{1/2} T \left[e^{-\beta H(t)} e^{iky} e^{-\beta H(t)} f^2 \right] \\ &= (2\pi/\beta K)^{1/2} T \left[e^{-\beta H(t)} e^{iky} e^{-\beta H(t)} e^{-\beta H(t)} f^2 \right]. \end{aligned}$$

According to a slightly modified Baker-Campbell-Hausdorff expansion (Hausser and Schwartz, 1968) we have

$$e^{A+B} e^B e^{AB} = \exp \left(A + B + \frac{1}{2!} [B, [A, B]] - \frac{1}{3!} [A, [B, [A, B]]] + \dots \right)$$

where

$$[A, B] = AB - BA,$$

so

$$\begin{aligned} &e^{-\beta H(t)} e^{iky} e^{-\beta H(t)} e^{-\beta H(t)} f^2 \\ &= \exp \left(-\beta H(t) + \frac{1}{2!} \beta^2 H(t)^2 \right. \\ &\quad \left. - \frac{1}{3!} \left[\beta^2 H(t)^2 (y^2 + 1.0008K) \left[\frac{d^2}{dy^2} D^2(y) \right. \right. \right. \\ &\quad \left. \left. \left. + 2 \frac{d^2}{dy^2} D^2(y) \left(\frac{d}{dy} + D^2(y) \frac{d^2}{dy^2} \right) \right] \right] + \dots \right) \end{aligned}$$

Neglecting the $1/K^2$ term, we therefore have made the replacement

$$T = (2\pi/\beta K)^{1/2} \exp \left[- \left\{ \frac{-1}{2\pi K} \frac{d^2}{dy^2} + D^2(y) + D^2(y)^2/(24K) \right\} \right].$$

Now it is easily seen that the expectation of y_1 (equivalent to that of any y_j) is given by

$$\mathbb{E}(y) = \text{Tr}_y T^2 / \text{Tr} T^2.$$

However, we know that, for large K , if φ_0 is the largest eigenvalue of the real symmetric operator T , then $(T/\varphi_0)^2 \langle \psi, \psi' \rangle = \varphi_0 \delta(\psi, \psi')$, where φ_0 is the

non-localized eigenfunctions belonging to λ_0 , and so

$$\delta(\gamma) = \int p d\gamma \psi^2 d\gamma.$$

Finally $\psi(\gamma)$ depends of course on the precise potential $\tilde{U}(\gamma)$ that is used, but the two general possibilities are easily seen. If λ_0 is part of the discrete spectrum, then $\delta(\gamma)$ (localized and relative $\tilde{U}(\gamma)$) results; if there is no discrete spectrum, $\psi(\gamma)$ can be normalized only if the domain of γ is restricted to being finite, and then $\delta(\gamma)$ will diverge as the size of the domain increases. There will thus be a phase transition to unbounded transverse motion if the discrete spectrum disappears at some value of β . Clearly,

$$\lambda_0 = e^{-\beta E} (\ln(\beta E))^{1/2}$$

where

$$-\frac{1}{2\pi i \epsilon} \partial_x^2 U(x) + \beta \tilde{U}(x) \psi(x) = \lambda_0 \psi(x),$$

where λ_0 is the lowest eigenvalue of the accompanying operator and \tilde{U} is the modified potential, as above.

Now an even one-dimensional potential (the crucial point is that $\tilde{U}(x) = U(x-x_0)$ with a trough below the lowest asymptotic value) will always have a localized discrete eigenstate. However, U is not of this form, and so for sufficiently high temperature — low β — there will be no bound state, leading to a γ -probability distribution stretching out uniformly to infinity. A quick estimate of the temperature at which rigid delocalisation in this sense occurs is given by the standard JWKB method of solving the Schrödinger equation satisfied by $\psi(\gamma)$ and gives precisely the same result as a "semiclassical" tunneling estimate, which is as follows. In the usual Schrödinger equations of quantum physics, say for unit mass, the parameter $\beta E'$ is equated with $1/\hbar^2$, where \hbar denotes $(1/2\pi i)$ times Planck's constant, and then, not distinguishing between \tilde{U} and U' , the classical mechanics being analysed has the energy $E' = [\hbar^2 p^2 + \beta U'(x)]$. Each discrete eigenstate of energy E' covers a volume $2\pi R = 2\pi/\sqrt{\beta E'}$ in (x, p) space, a shell whose thickness is precisely $[\hbar^2 + \beta U'(x)] = \beta$. Hence the number of discrete eigenstates up to the energy value E' is given by

$$\begin{aligned} N(E') - \frac{1}{2} &= \iint_{\{x^2 + p^2 \leq 2\pi R\}} dy dp \psi(2\pi/\sqrt{\beta E'}) \\ &= \frac{1}{\pi} \sqrt{\beta E'} \int p dx = \mu(\beta) \nu(U'^{-1}) dx, \end{aligned}$$

integrated over the range in which the square root is real, and so the number of states up to the continuum, which starts at $E = \hbar\gamma/\omega_0$, is given by 1. Then

$$\frac{1}{2}\pi\sqrt{\hbar E} = \iint D\eta D\psi(\psi - D\eta)^{1/2} d\eta,$$

or

$$\frac{1}{2} = (\sqrt{E}/\pi) \int W(\psi) = D\eta^2 e^{D\eta^2/2},$$

In particular, for the (Dihedral)potential of the Peacock-Bishop model,

$$\begin{aligned} \frac{1}{2} &= \sqrt{\hbar E}/\pi \int_{-\pi/2\pi/2} D\eta^2 [C e^{-D\eta^2} - e^{-D\eta^2}]^{1/2} d\eta \\ &= (4/\pi\omega_0)\sqrt{\hbar D} \end{aligned}$$

A.2. Dynamics

The fact that the original Ising model "respects" no separation in thermal equilibrium (Trotter) is a consequence of the energy fluctuations available in a thermal ensemble defined as being supplied by a heat bath, e.g., the aqueous environment. A more relevant question might be this: Suppose that the temperature is low enough that unpaired pairs of strands are at least metastable, and localized energy is supplied, e.g., by RNA polymerase; what then will be the time-development of the resulting strand separation pattern? To start with, let us anchor the pair of strands in its potential minimum ($\eta_0 = \beta$) and look at the small-amplitude motion in the vicinity of this state. For this purpose, note

$$H_0 = \sum_i \frac{m_i^2}{2} \dot{\eta}_i^2 + \sum_i D_i \eta_i + \frac{1}{2} \sum_i D_i^2 \eta_i^2 \phi_i - E^2 + \frac{K}{2} \sum_i (\eta_i - \eta_{i-1})^2,$$

so the equations of motion for $(\eta_i = \eta_0 + \delta_i)$ take the form $\ddot{\delta}_i = -m_i^2(\phi_{i+1} - 2\phi_i + \phi_{i-1}) + m_i^2\phi_i = 0$ (where $m_i^2 = K/m_i$ and $\phi_i^2 = D_i^2/(E^2\omega_i)$). Rather than solve this easily solvable linear equation, let us observe that, on summing over i ,

$$\frac{d^2}{dt^2} \sum_i \delta_i + m_0^2 \sum_i \delta_i = 0,$$

so if, e.g., we add pure kinetic energy at time 0, with $\sum_i \delta_i(0) = 0$, then

$$\sum_i \delta_i(t) = \Delta m_0 \cos \omega t.$$

On the other hand, multiplying by e^t and summing, we obtain

$$\frac{d^2}{dt^2} \sum_{n=1}^{\infty} n^2 q_n + m_0^2 \sum_{n=1}^{\infty} n^2 q_n = 2m_0^2 \sum_{n=1}^{\infty} q_n,$$

with the solution

$$\sum_{n=1}^{\infty} n^2 q_n(t) := -A \frac{m_0^2}{m_0^2 - \omega_0^2} \cos(\omega_0 t).$$

In other words, not only does an initial localized excitation oscillate as expected, but it spreads spatially as well.

The above linear analysis need not be valid beyond a short time. At longer time, the spreading of any initial distribution is very much affected by the excited nonlinear terms. There are many ways of seeing this, but perhaps the simplest is that of equivalent linearization (Krylov and Bogoliubov, 1947) an intelligent modification of the familiar variation-of-constants approach. Suppose that the potential minima successive at $t = 0$, as it does in the Pöschl-Teller model. The basic solution of the linearly truncated dynamics is of course exponential, in both the variables x and t : $q_n = e^{i\omega_n t}$, where substitution into the linearized equation shows that

$$\omega_n^2 t^2 = \omega_0^2 + 2m_0^2(1 - \cos t).$$

We now create a time-dependent envelope:

$$g_n(t) = \frac{1}{2} [E_n(t) e^{i\omega_n t} + E_n^*(t) e^{-i\omega_n t}],$$

and ask when this can satisfy the full equation.

A neat procedure is to work with the Lagrangian rather than the Hamiltonian.

$$L_T = \sum \frac{m}{2} \dot{x}_n^2 + \frac{K}{2} \sum (\dot{x}_n(t) - \dot{x}_{n+1}(t))^2 + U(x_n),$$

and, because we are interested in the slowly varying envelope $E_n(t)$, we average the Lagrangian over a few cycles of the basic oscillation. That by writing

$$\phi = \omega_0^2 - \omega_n^2(t)$$

and then using $(e^{iKt})^k = 0$ for any integer $K \neq 0$, we have

$$\begin{aligned} \left[\sum \dot{x}_n^2 \right] &= \frac{1}{4} \sum [E_n^2 - i\omega_n E_n \phi^2 + (E_n^*)^2 + i\omega_n E_n^* \phi^2] \\ &= \frac{1}{2} \sum (E_n^2 - i\omega_n E_n) (P_n^2 + i\omega_n P_n^2) \end{aligned}$$

and similarly

$$\begin{aligned} \left[\sum (x_n - x_{n-1})^2 \right] &= \frac{1}{4} \sum ((R_n - R_{n-1}) e^{i\theta_n})^2 + (R_n^* - R_{n-1}^*) e^{i\theta_{n-1}} (R_n^* - R_{n-1}^*) \\ &= \frac{1}{2} \sum (R_n - R_{n-1}) e^{i\theta_n} (R_n^* - R_{n-1}^*). \end{aligned}$$

Therefore for $D(x) := (x^{1/2} - 1)^2$, we see

$$x_n = R_n \left(\cos \left(\theta + \frac{1}{2\pi} \ln R_n / R_0 \right) \right)$$

to obtain (J_0 and J_1 are the usual modified Bessel functions)

$$\begin{aligned} \left[\sum D(x_n) \right] &= D \sum (e^{-i\theta_n} - 2e^{-i\theta_n} + 1) \\ &= D \sum [J_0(2\pi |R_n|) - 2J_0(\pi |R_n|) + 1]. \end{aligned}$$

Hence

$$\begin{aligned} \langle R_n \rangle &= \frac{1}{D} \sum (R_n - \ln R_n) (R_n^* + \ln R_n^*) \\ &= \frac{D}{4} \sum (R_n - R_{n-1}) e^{i\theta_n} (R_n^* - R_{n-1}^*) \\ &= D \sum [J_0(2\pi |R_n|) - 2J_0(\pi |R_n|) + 1], \end{aligned}$$

yielding the equation of motion $i\partial_t(\partial_x^2 \phi)(J_0(\pi |R_n|)/\pi |R_n^*|) = \phi(J_0(\pi |R_n|)/\pi |R_n^*|)$, or

$$\begin{aligned} \ddot{R}_n - 2\omega R_n - \omega^2 R_n &= i\phi(R_{n+1}) e^{i\theta_n} - 2R_n + R_{n-1} e^{i\theta_n^*} \\ &\quad - \frac{4\pi D}{\pi} R_n \frac{1}{|R_n|} [J_0(2\pi |R_n|) - J_0(\pi |R_n|)]. \end{aligned}$$

Because we know that $R_0 = \text{const.}$ is a solution at small R_n , we subtract $-i\omega^2 R_n = 2\omega [j_0(\pi |R_0|) - 1] R_n - (4\pi D/\pi) \frac{1}{|R_0|} R_n$, giving us

$$\begin{aligned} \ddot{R}_n - 2\omega R_n &= i\phi[\cos \theta(R_{n+1} - 2R_n + R_{n-1}) + (\omega^2/\sin \theta)(R_{n+1} - R_{n-1})] \\ &\quad - \frac{4\pi D}{\pi} R_n \frac{1}{|R_n|} \left[J_0(\ln |R_n|) - \left[J_0(\ln |R_0|) - \frac{1}{2} \omega |R_0| \right] \right]. \end{aligned}$$

If the envelope changes slowly in space on the scale of bare-bare separation, we can replace this with

$$\begin{aligned} \ddot{R}_n - \omega_0^2 \cos \theta R_n^* &= 2\omega \dot{R}_n + 2\omega_0^2 \sin \theta R_n^* \\ &\quad - \frac{4\pi D}{\pi} R_n \frac{1}{|R_n|} \left[J_0(\ln |R_n|) - J_0(\ln |R_0|) - \frac{1}{2} \omega |R_0| \right]. \end{aligned}$$

in which form some qualitative aspects are apparent. For example, imagine a frequency-shifted traveling envelope (Renaissenet, 1986)

$$F_n(t) = e^{i\omega t} F(t - nt),$$

so that

$$\begin{aligned} (\omega_0^2 \cos \theta - v^2) F'' + 2i [\omega_0^2 \sin \theta - (\omega - \delta)v] F' \\ = [2\omega_0^2 + f(\alpha|F|) - \delta^2] F, \end{aligned}$$

where

$$\omega_0^2 = 2\pi^2 D/m, \quad f(x) = \frac{1}{x} \left[I_0(2x) - I_1(x) - \frac{1}{2} x \right].$$

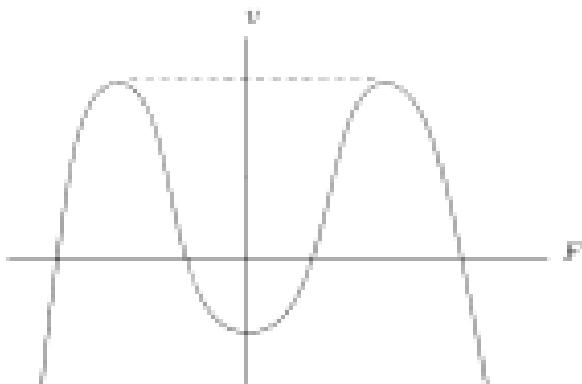
This will have a real solution if

$$v = \omega_0^2 \sin \theta / (\omega - \delta)$$

(note that because $\omega^2 = \omega_0^2 + 2\omega_0^2(1 - \cos \theta)$, then $\omega(d\omega/d\theta) = \omega_0^2 \sin \theta$, so $v_g = \omega_0^2 \sin \theta / \omega$ is the group velocity of a wave packet), and then will take the form

$$MF'' = -V'(F),$$

with the double-well potential $V(F)$, as shown. This, we can see, means that at finite amplitude we have a solution (dashed line in figure) with a moving "kink" if $M > 0$, but otherwise (and for $M < 0$) only a moving wave train.



5.3. Effect of Heterogeneity

Of course, the essence of DNA lies in its heterogeneity, and we typically might expect that varying the chain-chain interactions, e.g., having a random set

$(V_0(\beta_0))$ rather than a constant one, $V(\beta_0)$, would impart a detailed structure to the dynamics, as well as to the thermally equilibrated state. Numerical work has shown that this is not the case. However, a hint as to what might be going on is supplied by the observation that if the stiffness K in the basic model is given a β_0 dependence, e.g.,

$$K \rightarrow K_1 + (K_2 - K_1)e^{-\beta_0 t},$$

thereby accounting the effect of any fluctuations in the $\{\alpha_i\}$, the hysteresis "melting" transition appears to sharpen to first order (Peyraud, 1990). This suggests that it is the variation in stiffness that is responsible for spatially preferred domains of separation. Simulations with the random noise applied to $K(\beta)$ instead (Cale and Hwa, 1990) have also verified this supposition. No corresponding analytic work exists.

Assignment 8

1. Examine the statistics of the traditional score $W = \sum f_{ij} \delta_j(\alpha)$ of Subsection 3.3.3.
2. Discuss the relevance of the Arnold-Kolmogorov theorem on the recursive representation of many-variable functions to design of neural networks.
3. How is the Peyraud-Bishop analysis of Hagedorn breaking modified by heterogeneity in the field strength?

Bibliography

(A very incomplete list, restricted to those I happened to find useful)

- Abrams, A., Serein, K. M., and Lapidus, A., "Application of neural networks and information theory to the identification of E. coli transcription promoters," Los Alamos Report LA-UR-91-721 (Los Alamos National Laboratory, Los Alamos, NM, 1991).
- Alberts, B., Bray, D., Lewis, J., Raff, M., Roberts, K., and Watson, J. D., *Molecular Biology of the Cell* (Garland, New York, 1983).
- Allegood, P., Baloni, M., Chignell, P., and Rice, R. L., "Non-Poisson statistics of mammalian tissues: the DNA sequence of prokaryotes," *Phys. Rev. A* 58, 3680-3693 (1998).
- Almeida, T. and Pernici, A., "Long distance-range correlations in genomic sequences," *J. Stat. Phys.* 47, 135-155 (1987).
- Appelacio, A., Brock, M. L., Lewand, S., and Xu, X., "Efficient detection of genomic words," *J. Comput. Biol.* 7, 71-94 (2000).
- Arenz, R., Landau, E. N., Turner, S., and Wasserman, M. H., "Genomic mapping by averaging random clones: a mathematical analysis," *Genomics* 11, 486-497 (1991).
- Arenz, R., Chikudate, C., and Davies, L., "Primer representation and the Chen-Davis Method," *Stat. Sci.* 8, 403-434 (1993).
- Bach, F. and Bengio, Y., *Machine Learning* (MIT Press, Cambridge, MA, 1998).
- Bentley, D. J. and Bentley, D. C., "The design of pooling experiments for obtaining a clone map," *J. Planta Genet.* 14, 565-585 (1997).
- Bentley, R. and Bentley, T., *Assembly Computing* (Oxford, London, 1998).
- Bent, G. L., "Evolution of single sequence repeats," *Comput. Chem.* (1994).
- Berman, S. A., Cohen, M. A., and Glaser, G. H., "Empirical and structural models for insertion and deletion in the divergent evolution of proteins," *J. Mol. Biol.* 229, 1061-1082 (1992).
- Bermudo-Garcia, P., Arrieta, L., Chapela, P., Oliva, L. L., Rojas-Perez, R., and Stenberg, H. H., "Hunting for links between coding and noncoding DNA regions by entropy representation method," *Phys. Rev. Lett.* 85, 1363-1366 (2000).
- Bernardi, G., "The nucleic organization of the human genome," *Annu. Rev. Genet.* 29, 437-462 (1998).

- Borodkin, G. I., Glazier, J. D., and Krichevsky, M. H., "Global fractal dimension of human DNA sequences treated as pseudorandom walks," *Phys. Rev. A*, **45**, 2922-2942 (1992).
- Bortnick-Goldstein, P., Konon-Kostka, R., and Weiss, L. L., "Compositional representation and long-range correlations in DNA sequences," *Phys. Rev. E*, **53**, 3491-3499 (1996).
- Bishop, D. E., et al., "Number of polymorphic DNA classes required to map human genome," in *WGS* (1993).
- Bonfield, D., White, R. L., Krichevsky, M., and Davis, R. W., "Characterization of a genetic linkage group by measuring restriction fragment length polymorphisms," *Am. J. Med. Genet.* **32**, 314-321 (1992).
- Bordoli, R., "Weight matrix description of four eukaryotic RNA polymerases II promoter elements," *J. Mol. Biol.* **212**, 365-379 (1990).
- Casella, V. and Sankoff, D., "Length correlation interpretation of DNA sequence repeats," *J. Appl. Probab.* **12**, 366-375 (1975).
- Casella, V. and Sankoff, D., "An upper-bound technique for lengths of common subsequences," in *Sankoff and Seshai* (1993), pp. 273-278.
- Chen, L. H. T., "Poisson approximation for dependent trials," *Ann. Probab.* **3**, 534-551 (1975).
- Churchill, G. A., "Stochastic models for heterogeneous DNA sequences," *Scand. Statist. Probab. Model.* **11**, 79-94 (1984).
- Churchill, G. A., "A coding tree containing systems for the optimal placement representation of the entire DNA genome," *Cell* **9**, 11-16 (1976).
- Cooper, N. G., "The human genome project," *Los Alamos Sci.* **19**, 103-105 (1992).
- Cox, D. R., Burdette, M., Price, B. R., Kim, S., and Myers, R. M., "Haplotype linkage mapping," *Science* **266**, 143-146 (1994).
- Durrett, R., *Mathematical Methods in Population Genetics* (Princeton Univ. Press, Princeton, NJ, 1991), Chap. 10.
- Elliott, D. and Fries, T., "Decimation of heterogeneous DNA," *Phys. Rev. Lett.* **79**, 2373-2376 (1997).
- Fernandes, R., Rangwala, G., Maxted, M. V., Parsons, A. D., and Freaney, D. C., "Weight matrix for optimally predicting DNA sequence borders," in *Proceedings of the Annual ACM SIGART Symposium (SIGART) Society for Industrial and Applied Mathematics*, Philadelphia, 2000, pp. 444-459.
- Fisher, J., "Probabilistic behavior of longest common subsequence length," in *Sankoff and Seshai* (1993), pp. 309-320.
- Fordell, E. L., Sankoff, D., and Orlitzky, M., "How big is the Universe of DNA?," *Science* **258**, 1277-1282 (1992).
- Frost, D. H. and Lippert, M., "Heating trees and similarity detection in sequence alignment with gaps," *J. Comput. Biol.* **7**, 115-141 (2000).
- Grimm, R., Otto, S., Krug, A., and Stoye, J., *Biological Sequence Analysis* (Cambridge Univ. Press, New York, 1998).
- Gruber, B. A., "Assembly of open form surface transposable genetic elements: implications for the evolution of protein-protein interactions," *J. Theor. Biol.* **196**, No. 1, 11-27 (1998).
- Hillis, D. A., "Theoretical Models for Heterogeneity of Base Composition in DNA," *J. Theor. Biol.* **40**, 313-333 (1973).

- Pelizzetti, D., Laprade, A., and Kostadin, K. M., "Generalization of evolutionary coding schemes," Los Alamos Report LA-UR-94-134 (Los Alamos National Laboratory, Los Alamos, NM, 1994).
- Penttila, D., Sung, T., Kirk, R. M., Suttorp, S. J., and Thygesen, E., "An algorithmic approach to multiple sequence digest mapping," *J. Comput. Biol.* **8**, 187-207 (1995).
- Petkov, W., *Probability Theory and its Applications* (Wiley, New York, 1976), Vol. I.
- Pielou, E. C., Turner, D. C., and Reid, D. R., "Some computational aspects of processes," *Canadian J. Zool.* **33**, 2093-2102 (1955).
- Pielou, E. C. and Sung, C. H., "Assessment of protein coding measures," *Proc. British Roy. Soc.* **26**, 605-610 (1952).
- Pielou, E. C., in *Probabilistic Approaches in the Sciences: Combinatorics of Distributions*, Hermann et Cie, Paris 1949.
- Pratt, R., "The code theory of genes," Cold Spring Harbor Symp. Quant. Biol. **LII**, 389-393 (1957).
- Quinton, P., *Mathematical Methods of Analysis of DNA Sequences*, Vol. II of IJEMACS Series (American Mathematical Society, Providence, RI, 1993).
- Quinton, P., "Pattern approximation and DNA sequence matching," *Commun. Stat. Theory Meth.* **17**, 1167-1189 (1988).
- Quinn, S. J. and Hause, H., "New method for mapping genes in human chromosomes," *Nature Genetics* **29**, 589-594 (1993).
- Quinton, P. and Ohyama, R., "Long repetition patterns in cDNA sequences," *Z. Naturforschung-B* **45**, 341-350 (1990).
- Ramsey, J., *The Theory of Ramsey* (Columbia University Press, 1956).
- Ramsey, J., *Algorithms on DNA: Strings and Sequences* (Cambridge Univ. Press, New York, 1993).
- Han, B., Xie, H., Yu, Z., and Chen, H., "A combinatorial problem related to ranked and under-represented strings in bacterial complete genomes," *Int. J. Comp. Bio.* (2000).
- Hausser, M. and Schmid, J. T., "Lie groups, Lie Algebras" (Dordrecht and Boston, 1989).
- Hayashi, C., "On the prediction of processes from qualitative data," *Ann. Inst. Stat. Math.* **A30**, 43-56 (1978).
- Hausser, H., Böring, W., and Schmid, J. T., "Principles of biosequences: the role of repeats," *Phys. Rev. E* **50**, 5061-5071 (1994).
- Hausser, H. and Schmid, J., "Correlations in DNA sequences: the role of protein coding segments," *Phys. Rev. E* **54**, 826-830 (1996).
- Hausser, H. and Lusby, M., "Similarity detection and localization," *Phys. Rev. Lett.* **76**, 2293-2294 (1996).
- Soh, T., "Probabilistic sequences of DNA and their structural analysis," *Bull. Chem. Soc. Am.* **51**, 2077-2082 (1955).
- Soh, T., "Compositional distribution analysis of T-oligonucleotide signals of self-5% processes in higher eukaryotic genes," *J. Theor. Biol.* **100**, 109-118 (1983).
- Kudin, S. and Kostadin, K. M., "Probabilities and correlations in DNA sequences," *Science* **268**, 472-480 (1995).
- Kudin, S., Gao, F., and Kostadin, K. M., in *White noise* (1995), Chap. 6.
- Kudin, S. and Blasius, C., "Some statistical problems in the assessment of inherent properties of DNA sequence data," *J. Am. Stat. Assoc.* **86**, 27-35 (1991).

- Karp, R. M. and Shamir, R., "Algorithms for optimal mapping," *J. Comput. Biol.* **7**, 383-394 (2000).
- Kasai, T. and On, P., "Maximum length of common words among random letter sequences," *Ann. Probab.* **14**, 333-354 (1986).
- Katz, S. W., "Observe-estimate-estimate: Iteration of Master-Slave," *Technometrics* **23**, 243-249 (1981).
- Kim, S. and Segal, A. M., "EMMAP: a refined pattern-matching approach to align protein sequences," *J. Comput. Biol.* **6**, 155-166 (1999).
- Kirkhoff, P. and Bogolyubov, N., *Brownian Motion in Non-Linear Mechanics* (Princeton Univ. Press, Princeton, NJ, 1937).
- Lander, E. S. and Waterman, M. L., "Genome mapping by fingerprinting and analysis: mathematical analysis," *bioRxiv* 2, 231-239 (1998).
- Lander, E. S., "Genotype with no heterozygotes," in: *Waterman's Problems*, Chap. 2.
- Lange, K. and Roederer, J., "How many polypeptides grow with it takes to open the human genome?" *Am. J. Hum. Genet.* **34**, 842-843 (1984).
- Lubotzky, R. H. and Baum, T. O., "Mixed optimality genetic decoding with general alignment and partial pair score functions," *J. Mol. Biol.* **285**, 931-945 (1999).
- Li, W., "Shared information between correlation functions," *J. Stat. Phys.* **85**, 1031-1037 (1996).
- Li, W., "The study of correlation structure of DNA sequences: a critical review," *Comput. Chem.* **21**, 259-272 (1997).
- Li, W., "Statistical properties of spin pairing theory in complete genome sequence," *Comput. Chem.* (1998).
- Lodewijks, B. and Thielke, D. M., "Significantly lower entropy estimate for natural DNA sequences," *J. Comput. Biol.* **4**, 113-141 (1997).
- Li, W. and Gaoxin, K., "Long range correlation and partial 1/ ν spectrum in a noncoding DNA sequence," *Complexity Lett.* **12**, 303-307 (1992).
- Li, W., Shao, T.-G., and Kaoxin, K., "Understanding long-range correlations in DNA sequences," *Physica* **20** (1994).
- Lindström, H., *J. Theor. Biol.* **18**, 289-308 (1969).
- Li, X., Liu, Z., Chen, H., and Li, Y., "Characterizing self-similarity on bacterial DNA sequences," *Phys. Rev. E* **53**, 1535-1544 (1996).
- Madden, J., *Nature* London **199**, (1963).
- Moresco, G. and Stocchi, R., "Construction of physical maps from oligonucleotide fingerprint data," *J. Comput. Biol.* **6**, 271-281 (1999).
- Moyen, Y., *Biostatistics* (Cambridge Univ. Press, New York, 1992).
- Mitra, R. M., *Statistical Mathematical Quantities in Biology*, Vol. II of *DNA Sequence Analysis Series* (American Mathematical Society, Providence, RI, 1995).
- Miyazawa, S. and Karplus, K. L., "Formation of helices in interacting protein surfaces: How protein crystal structures?," *bioRxiv* **14**, 534-550 (1995).
- McLennan, A. K. and Narayan Rao, A. V. S. S., "Theoretical moment analysis shows a characteristic length scale in DNA sequences," *Phys. Rev. Lett.* **84**, 1451-1455 (2000).
- Milnor, J. M., *Amer. Journ. Math.* **64**, 1-16 (1942).
- Moss, R. F., Kikuchi, T. B. L., and Kaoxin, K. Y., "An accurate approximation to the distribution of the length of the longest matching word between two random sequences," *Bull. Math. Biol.* **52**, 773-794 (1990).

- Muth, R. and Weber, R., "Hypergeometric statistics of gapped alignments," *J. Comput. Biol.*, **6**, 311-332 (1999).
- Nathans, D. and Smith, B. D., "Nucleotide substitution rates," *Annu. Rev. Biochem.*, **44**, 273-298 (1975).
- Neidhardt, R. H. and Weinstock, C. D., "A general method applicable to the search for similarity in the amino acid sequence of two proteins," *J. Mol. Biol.*, **48**, 443-453 (1970).
- Ong, G., New York University, New York, NY 10012 (personal communication, 1992).
- Ping, C. K., Baldwin, R. V., Goldberger, A. L., Berlin, B., Schatzkin, B., Stevens, M., and Stanley, H. E., "Long-range correlations in nucleotide sequences," *Nature (London)*, **356**, 126-129 (1992).
- Ping, C. K., Baldwin, R. V., Berlin, B., Stevens, M., Stanley, H. E., and Goldberger, A. L., "Scaling properties of DNA nucleotide," *Phys. Rev. A*, **49**, 1483-1493 (1994).
- Perez, O. R. and Perez, J. R., "String matching for the terrorist," *Am. Math. Monthly*, **102**, 544-563 (1995).
- Perez, J. R., *Computational Methods* (Springer-Verlag, New York, 1991).
- Perez, O. R., Perez, J. R., Bruck, and Tresser, D. C., "The impact of pooling design performance," *J. Appl. Prob.*, **36**, 591-606 (1999).
- Perez, O. R. and Perez, J. R., "Median length distribution in genome sequencing," *J. Math. Stat.*, **28**, 269-288 (1999).
- Perez, J. R., Perez, O. R., and Perez, A. R., "Universality of self-similarity distributions," *AIHPA/ASPA Conf.* **11**, 33-36 (1993).
- Peyrard, M., "Nonlinear energy localization in Noncopolar" in *Biological Physics*, Prigogine, I., Nicolis, G., and Babloyantz, A., eds. (American Institute of Physics, New York, 1998), p. 147-152.
- Peyrard, M. and Bishop, A. R., "Statistical mechanics of a nonlinear model for DNA denaturation," *Phys. Rev. Lett.*, **61**, 2711-2714 (1988).
- Perc, E., Sos, P., Sherris, D., and Waterman, M. S., "Chemical mapping by multi-dimensional random walks: a mathematical analysis," *Chemistry*, **26**, 99-120 (1997).
- Perez, A. et al., "Biopolar approach to maximum position," Cold Spring Harbor Symp. Quant. Biol. **51**, 131-136 (1986).
- Pelton, A. E., "A model for high order Markov chains," *J. R. Stat. Soc. Ser. B*, **42**, 529 (1974).
- Ramsey, V. R., Makarev V. Jr., Royston, M. Jr., and Tsuneyama, T. G., "DNA sequencing through the Bayesian approach," *J. Comput. Biol.*, **7**, 219-231 (2000).
- Ramond, M., "Low-complexity feasible and optimal solution via quasi-one-dimensional physical models," *Phys. Rev. B*, **38**, 2385-2392 (1988).
- Ramond, M., "On the nature of a codon-ant," *Ann. Math. Stat.*, **13**, 93-114 (1942).
- Ramond-Rothlin, B., Dumas-Girardin, F., and Oliver, J. L., "Sequence computational complexity of DNA through an entropy representation method," *Phys. Rev. Lett.*, **83**, 1344-1347 (1999).
- Rand, R. K., Delcher, A. L., Salzberg, S., Schatz, G. J., Higgins, R., Horn, G. T., Heller, K. H., Blattler, B. A., "Primer-directed enzymatic amplification of DNA with a thermostable DNA polymerase," *Science*, **268**, 497-501 (1995).
- Ranjeva, P. et al., "DNA sequencing with chain terminating inhibitors," *Proc. Natl. Acad. Sci. USA*, **94**, 3463-3467 (1997).

- Randell, D., and Kunkel, J. A., *Flow Cytometry, Living Cells and Biomaterials*, John Wiley, New York, NY, 1982.
- Randell, D., and Matzke, S., "Common subsequences and maximal subsequences," in Randell and Kristal (1982).
- Schmidt, S., Howard, P., and Turner, R., "The effect of nucleotide sequence distribution on the progress of an STS mapping project," *J. Comput. Biol.* 5, 41-51 (2002).
- Schubert, H. E., "Determination of DNA fragment size from gel-electrophoresis mobility," in Wink (1982), pp. 1-14.
- Schubert, H. E., and Randell, D. R., "Improved estimation of DNA fragment lengths from agarose gels," *Anal. Biochem.* 118, 113-123 (1982).
- Selkoe, P., and Whittemore, M. S., "Multiple solutions of DNA restriction mapping problems," *Ann. Appl. Stat.* 12, 412-427 (1998).
- Selkoe, D. C., Li, X., Hernandez, L. L., Hernandez, S. P., Hsu, R. J., and Wang, T. H., "Unlinked restriction maps of meiotic stages chromosomes constructed by optical mapping," *Biochem. Biophys.* 110, 111-114 (1993).
- Selkoe, D. C., and Abelson, K., "A simulated annealing algorithm for the clustering problem," *Pattern Recogn.* 24, 1003-1010 (1991).
- Selkoe, P. H., "The mathematical modeling of pattern similarities in genetic structures," in *Mathematics and Computers in Biological Applications*, Rosenfeld and G. Lin, eds., 23-27 (Marcel Dekker, New York, 1982), Chap. 23.
- Selkoe, P. H., "Pattern recognition in DNA," in *Recent Mathematics Survey* (1980), Chap. 19.
- Shapiro, M. D., and Sonnedecker, E., "RNA splicing junctions of different classes of eukaryotes," *Mol. Androcl. Res.* 25, 7055-7074 (1987).
- Sternin, O., et al., "Predicting human genome maps with radiolabel hybrids," *Proc. Roy. Soc. B* 277, 277-288 (1997).
- Stevens, Lopatin, and Zha, "Predicting protein-protein interactions using neural net and statistical methods," *J. Mol. Biol.* 1992.
- Tarolli, A., Stevens, D. R., Salas, M., Lopatin, I. M., and Everett, M. F., "Testing the specificity of genes: the evidence from protein interaction," *Science* 268, 100-103 (1995).
- Stevens, O. D., and Howell, G. W., "Identifying protein-building sites from unaligned DNA sequences," *Proc. Natl. Acad. Sci. USA* 86, 1163-1167 (1989).
- Tomas, R., and Chittenden, R. W., "Some statistical aspects of the primary structure of nucleic acids sequences," in Whetstone (1980a), Chap. 8.
- Trotter, M., Chaitin, L. L., and Pollock, R. W., "Statistical model of the DNA molecule," *Phys. Rev. A* 40, 1828-1842 (1989).
- Tung, R., "Determination of the order of a Motzkin chain by dynamic programming iteration," *J. Appl. Probab.* 18, 488-497 (1981).
- Tunney, D. C., Davis, C., Stevens, O., and Stevens, M. M., "Computation of π^* : a measure of sequence dissimilarity," in *Computers and DNA*, G. Bell and T. Margolin (Eds.), John Wiley, New York, NY, 1992, pp. 139-159.
- Tunney, D. C., "Mapping using unique sequences," *J. Mol. Biol.* 217, 189-202 (1991).
- Tunney, D. C., "Phenotype responses in a language morphological class of "words,"" in *Classification and Related Methods of Data Analysis*, H. B. Neudeck (Ed.), New York, 1988, pp. 31-44.

- Turrisi, A. A., Klassen, P., Elizondo, J. H., and Turrisi, R. A., "Wigner analysis of DNA sequences," *Phys. Rev. E* 53, 3828-3834 (1996).
- Uspensky, J. V., *Introduction to Mathematical Probability* (McGraw-Hill, New York, 1937).
- Yano, R. F., "Illustration of long-range fractal correlations and 1/f noise in DNA base sequences," *Phys. Rev. Lett.* 64, 3625-3630 (1990).
- Wolenski, M. S., "Properties of repetition sites," *Proc. Acad. Sci. U.S.S.R.* 11, 1051-1056 (1963).
- Wolenski, M. S., "Probability distribution for DNA sequence comparisons," *Zhurnal Biostat. Lit.* 17, 28-36, American Math Soc., Providence (1986).
- Wolenski, M. S., ed., *Mathematical Methods for DNA Sequences (CD)*, Boca Raton, FL, 1999a.
- Wolenski, M. S., *Commonalities Between DNA Sequences*, in Wolenski (1999a), Chap. 4.
- Wolenski, M. S., *Commonalities Method for Finding Simple Repeated Motifs*, in Wolenski (1999a), Chap. 6.
- Wolenski, M. S., *Introduction to Computational Biology* (Chapman & Hall, London, 1999).
- Wolenski, M. S., Gordon, L., and Jernstol, R., "Phase transitions in sequence matches and motifs and structure," *Proc. Natl. Acad. Sci. USA* 84, 1136-1141 (1987).
- Yano, R. F., ed., *Statistical Analysis of DNA Sequence Data* (Marcel Dekker, New York, 1992).
- Yano, R. F., *Statistical Data Analysis* (Duxbury, Belmont, Mass., 1990).
- Zhang, J. J. Q. and Mao, T. G., Cold Spring Harbor, Long Island, NY (personal communication, 1992).
- Zhang, M. Q. and Mao, T. G., "Genome mapping by random ordering: a discrete theoretical analysis," *J. Stat. Phys.* 73, 615-633 (1994a).
- Zhang, M. Q. and Mao, T. G., "Statistical analysis of the human path genome DNA sequence and gene recognition," *J. Stat. Anal. Res.* 22, 1139-1159 (1994b).
- Zhang, M. Q. and Mao, T. G., "Molecular sequence alignment and its relation path analysis," *J. Statist. Anal.* 124, 119-129 (1995).
- Zhang, T., Zhang, W., Lin, J., and Chen, S. Z., "Theory of DNA matching based on the Frequency-Matching Model," *Phys. Rev. E* 56, 700-715 (1997).

Index

- stable, 14
stability, 19–21
- back propagation, 113
Balas-Campbell-Boland expansion, 126
bias, 2–9
 discretization, 11
Bayes' theorem, 111
Berkman inequality, 111
Bellman equation, 117
- Chebyshev bounds, 11
clustering, 10
clustering validity, 11
class
 theory, 11, 12
 and dimensionality, 11
 theory, 11, 12
 weighted regions, 11
- citation, 2
 cluster maps, 119
coherence, 101
coevolutionary product method, 114
coding, 13, 16, 20
convolution
 condition, 112
 cortex, 20
 diagram, 11
 functional measures, 111
 genotypes, 111
 phenotype, 111
 rate, 11
 valence, 11, 12
- convolutional networks, 24, 27–29
convergence, 6, 13, 19
convolution expansion, 11
- design, experimental, 29
diffusion, 11
- DNA, 1, 2
 complementarity, 117
 cytosine, 123–4
 functional regions, 117
 interconvertibility, 117, 124
 strand symmetry, 117
 unspecified, 117
 visual behavior, 118–123
- dissimilarity, 11, 112
dynamical programming, 116
- energy landscape information
 linear clusters, 11
 PCA, 116, 117
- ergodic distribution, 113
expansion, 113
extreme, stochastic, 117
PCA, 1, 10, 102–11
hypercube solution model, 26–27
- hypercube stability, 11–12, 26–27
Hölder-Poincaré, 112
Ising lattice, for repeated patterns, 11
lognormal size distribution, 1
- multichromosomal, 1
prior, 1, 4, 14, 27
partition function, 11, 16–17, 21
- Julia sets, 113
logistic growth, 113
Möbius ladder, 117
- negative labels, 11
negative-exclusive clusters, 11
noise, 10, 101, 102
information landscape entropy
 constant, 118–19
 natural, 119
 prior spectrum, 119–21

- intra, 4, 100–1
short, 10, 10
vertical, 10, 10–11
- PiYG method, 111
- process, 10
- Lipopeptide form, 124–5
medium, 10, 10
lipopeptide protein, 12
- loop
 bulge, 12
 cyclic, 12
 rotation, 12
- Master model:
 beta, 12
 hidden, 12–13
 higher-order, 12
 walking, 12–13
- Motif motif, 12
aligned, 12, 12
reciprocal exchange, 12–13
reverse protein, 12
gap, 12
strand, 12
superfamily, 12
state, 12, 12
independent model, 12
exchange, 12
coupled evolution model, 12
P-loop, 12
reciprocal form, 12
- Motif protein, 12
- mixed motifs, 111
- motif, 12
ortholog protein, 12
- partial digestion, 12
position specific, 12
protein model, 10, 10
- Protein digest motif, 12
- proteins, 111
- protein exchange, 12–13
walking, 12–13
- Protein signature, 12–13
protein distribution, 12
- probability distribution
 chi square, 12
- cluster size, 12
extreme value, 12
Gaussian, 12, 12
poisson, 12
Poisson, 12–13, 12
process, 12–13, 12
prior, 12–13
posterior, 12, 12
- random evidence, 12
realigning, 12
 open, 12
- realign
 open, 12
 target, 12
single sequence, 12
oligonucleotide, 12
feature matching, 12
- RNA, 12
 splice signal, 12
- secondary structure, 12, 12
- secondary structure based polymerization, 12
- RNA, nucleic acid, 12–13
- secondary site sequence compatibility:
 competition, 120–1
 distribution, 121
 nonrigid, 120–1
 rigid regions, 121, 121
 superstate, 121
 placement, 121
 phase transition, 121
- secondary structure (see also matching)
 distribution, 121
 competing, 121
 shifted, 121
- secondary superstate, 12
- superposing methods
 walking, 12
 supercoil, 12
 hybridization, 12
 optimal, 12
 changes, 12
- supercoil form, 12
- statistical methods of sequence correlation
 cluster distribution, 12
 signature matrices, 12–13
 prior model, 12–13
 Hidden Markov, 12–13
 walking Markov, 12–13

- binding approximation, 29, 41, 44, 105
biomass analysis
 categories, 102
 changes, 99
 linear functions, 109
 systems, 10
- biting set, 134
- biomass, 2
 analysis, 102
- biotopes
 Fischer, 14, 76
- biophase, 24
- black, 16
- blacklist, 77
- bounding regions, 120
- boutiques, 2
- box-flow theorem, 139
- bottom, 14
- boxe diagram, 10
- boundary
 multiple points, 103
- border
 correlation matrix, 61–2
 correlation between, 64