

heringa@few.vu.nl

Content

- · Problem example: Approaches to repeats findina
 - Background
 - FFT
 - Repro
 - Transitivity (TRUST)
- Parameterization of methods – GA
- Analyzing data Statistics
 - PCA
- · Comparing methods
- Objective function
- Standard of truth

Repeats

- · Evolution reuses developed material developed
- Multiple stoichiometric and spatially close combined structure-function relationships
- · In proteins, repeats vary from a single amino acid (e.g. poly-GIn) to complete domain sequences o combinations thereof.
- Many types of (near)identical repeats exist in genomes (Human genome > 50%) (next slide): Micro- and mini-satellites



- VNTRs Interspersed repeats
 - (http://globin.cse.psu.edu/courses/spring2000/repeats.html)
 - · LINE and SINE repeats · LTR retroposons, also called retrovirus-like elements
 - · DNA transposons



Protein repeats and disease

A number of neurodegenerative diseases have been found to be strongly associated with proteins containing a poly-glutamine stretch. A conformational change in the expanded polyglutamine stretch is believed to form the molecular basis for disease onset.

Five neurodegenerative diseases (Huntington disease (HD). spinocerebellar ataxia type 1 (SCA1), dentatorubral-pallidoluysian atrophy (DRPLA), Macado-Joseph disease (MJD), and spinobulbar muscular atrophy (SBMA)) have been found to be strongly associated with a protein containing a polyglutamine stretch which is greatly expanded in affected individuals (for a review, see D.E.Housman, Nature Genet. 8, 10 -11, 1995). For the five diseases, the mean length of the glutamine repeat in unaffected individuals is approximately twenty, and the cutoff for pathology is about forty (the cutoff may be higher for MJD (Housman, 1995). Furthermore, long polyglutamine stretches have been found in many transcription factors.



Genome Repeats

- Types of genome repeats Microsatellites, 2-3bp (e.g. (CA),)
- Minisatellites, 10-100bp, occurring at more than 1000 locations in the human denome)
- Variable number tandem repeats (VNTR) range from 14 to 100 nucleotides long that is organized into clusters of tandem repeats, usually repeated in the range of between 4 and 40 times per occurrence. Clusters of such repeats are scattered on many chromosomes.

VNTRs have been very effective in *forensic crime* investigations. When VNTRs are cut out, on either side of the sequence, by restriction enzymes and the results are visualized with a gel electropresis, a pattern of bands unique to each individual is produced. The nurber of times that a sequence is repeated varies between different individuals and between maternal and paternal loci of an individual. The likelihood th von cividuals having the same band pattern is extremely improbable.

- Interspersed repeats
- erspersed repeats SINE (short interspersed nuclear element), LINE (long interspersed nuclear element) (next slide). These are also called non-LTR or poly-A retro(trans)posons LTR retroposons Elements of several hundred bp in length, called the long terminal repeat, that appears at each end. Some autonomous elements are cousins of retroviruses (e.g., HIV) but are unable to survive outside of the cell, and are called endogenous retroviruses. None are known to be currently active in humans, though some are still mobile in mice. The so-called MaLR (mammalian LTR) elements, which arose before the mammalian radiation, seem to be non-autonomous repeats that move via proteins from endogenous retroviruses.
- DNA transposons. Full-length autonomous elements encode a protein, called transposase, by which an element can be removed from one position and inserted at another. Transposons typically have short inverted repeats at each end.

LINE and SINE repeats

(elaboration of preceding slide)

- A LINE (long interspersed nuclear element) encodes a reverse transcriptase (RT) and perhaps other proteins.
 - Mammalian genomes contain an old LINE family, called LINE2, which apparently stopped transposing before the mammalian radiation, and a younger family, called L1 or LINE1, many of which were inserted after the mammalian radiation (and are still being inserted).
- A SINE (short interspersed nuclear element) generally moves using RT from a LINE.
 - Examples include the MIR elements, which co-evolved with the LINE2 elements. Since the mammalian radiation, each lineage has evolved its own SINE family. Primates have Alu elements and mice have B1, B2, etc.
- The process of insertion of a LINE or SINE into the genome causes a short sequence (7-21 bp for Alus) to be repeated, with one copy (in the same orientation) at each end of the inserted sequence. Alus have accumulated preferentially in GC-rich regions, L1s in GC-poor regions.

How to delineate repeats

- 1. Supervised: if you have a repeat motif, use profilebased methods or the like
- 2. Nonsupervised
 - You want to find a single repeat type
 - You want to find tandem repeats
 - You want to find interspersed repeats (intervening sequence stretches
 - You want to find multiple repeat types

Fast Fourier Transformation

A **fast Fourier transform** (**FFT**) is an efficient algorithm to compute the discrete Fourier transform (DFT) and its inverse. FFTs are of great importance to a wide variety of applications, from digital signal processing to solving partial differential equations to algorithms for quickly multiplying large integers.

FFT is an intuitive algorithm for detecting repeats because it analysis **periodicity** in data



FFT

The Fourier transform converts a time domain representation of a signal into a frequency domain representation. The Fast Fourier Transform (FFT) is an optimized implementation of a DFT that takes less computation to perform. The Fourier Transform is defined by the following equation:

$$X(f) = F{x(t)} = \int_{-\infty}^{\infty} x(t) e^{-j2\pi ft} dt$$
,
(1)

However, a digitizer samples a waveform and transforms it into discrete values. Because of this transformation, the Fourier transform will not work on this data. Instead, the Discrete Fourier Transform (DFT) is used, which produces as its result, the frequency domain components in discrete values, or "bins." So, the Discrete Fourier Transform (DFT) maps discrete-time sequences into discretefrequency representations. DFT is given by the following equation:











Limitations of FFT-based approaches with respect to repeats finding:

- 1. Repeats can be interspersed
- 2. Multiple repeat types
- 3. Incomplete repeats



Szklarczyk, R. and Heringa, J. (2004) Tracking repeats using significance and transitivity. Bioinformatics 20 Suppl. 1, i311-i317.































poble
TRUST 43% 82% 24
TRUST -force 75% 85% 34
RADAR 63% 64% 86

















					Example Setup
ear-opt	imal p	ar	ап	netrisa	ation by Genetic Algorithm
enes (O	4] add	op [.]	t ra	andon	n values, for example from 0 to 9.
ey are	transi	foi	rm	ed to	yield the parameter space.
Gene	es			Par	ameters
wei	aht			=	0.05 * parameter[0];
gap	-open			=	(0.5 * parameter[1]) + 10;
gap	exte	nc	ł	=	0.5 * parameter[2];
add	-cons	ta	int	: =	0.4 * parameter[3];
Nr	ge	nc	ome	a	score (= fitness)
	[5	5	8	5]	153.6867
0:	[3	8	6	4]	153.6446
0: 1:		1	6	3]	153.3365
0: 1: 2:	[4	ю	~		
0: 1: 2: 3:	[4 [7	2	8	0]	153.2564

Exan Near	nple Setup -optimal param	etrisation by Genetic Algorithm
Gene They Gene	es [0-4] adopt ra are transforme es Parameters	ndom values, for example from 0 to 9. d to yield the parameter space.
weigh gap-c gap-e add-c	nt = 0.05 * para open = (0.5 * pa extend = 0.5 * p constant = 0.4 *	meter[0]; rameter[1]) + 10; arameter[2]; parameter[3];
Nr	genome	score (= fitness)
0:	[5 5 8 5]	153.6867
1:	[3 8 6 4]	153.6446
2:	[4 6 6 3]	153.3365
3:	[7 2 8 0]	153.2564
4:	[4 6 7 4]	153.0322



Typical Genetic Algorithm Scheme Iteration scheme of Genetic Algorithm

- 1. generate 200 random genomes ->
- 2. run alignments with paramaters of each gene ||
- 3. score result (fitness) | |

Г

- 4. sort (according to fitness) | <-
- 5. select top genomes and create new generation [5 5 8 5] [3 8 6 4]

X [3 8 8 5]

Genetic Algorithm (III)																						
erati	ion	re	es	ul	ts																	
	Iteration 1								15		Iteration 8											
0:	[5	5	8	5	1]	153	. 68	0:	[5	5	8	3	1]	154.90	0:	[5	5	8	3	1]	154.90	
1:	[3	8	6	4	0]	153	. 64	1:	[5	4	8	2	2]	154.68	1:	[5	4	8	3	2]	154.90	
2:	[4	6	6	3	0]	153	. 33	2:	[5	4	8	2	2]	154.68	2:	[5	4	8	3	2]	154.90	
3:	[7	2	8	0	2]	153	. 25	3:	[5	4	8	2	0]	154.64	3 :	[5	4	8	2	2]	154.90	
4:	[4	6	7	4	2]	153	. 03	4:	[5	4	8	3	0]	154.64	4:	[5	4	8	3	2]	154.90	
5:	[3	8	7	1	8]	152	. 48	5:	[4	5	8	3	3]	154.44	5:	[5	4	8	3	2]	154.90	
6:	[6	3	8	2	5]	152	. 37	6:	[4	5	8	3	0]	154.39	6:	[5	4	8	3	2]	154.90	
7:	[6	5	3	3	0]	152	. 37	7:	[4	5	8	3	0]	154.39	7:	[5	4	8	3	2]	154.90	
8:	[6	4	5	3	3]	152	. 30	8:	[4	5	8	3	0]	154.39	8:	[5	4	8	2	2]	154.90	
9:	[7	7	8	0	1]	152	. 01	9:	[4	4	8	3	2]	154.35	9:	[5	4	8	3	2]	154.90	
10:	[5	5	4	5	2]	151	.91	10:	[5	5	8	3	0]	154.32	10:	[5	4	8	3	2]	154.90	

Multivariate statistics – Principal Component Analysis (PCA)

Principal component analysis (PCA) involves a mathematical procedure that transforms a number of (possibly) correlated variables into a (smaller) number of uncorrelated variables called *principal components*. The first principal component accounts for as much of the variability in the data as possible, and each succeeding component accounts for as much of the remaining variability as possible.

Traditionally, principal component analysis is performed on a square symmetric matrix of type SSCP (pure sums of squares and cross products), Covariance (scaled sums of squares and cross products), or Correlation (sums of squares and cross products from standardized data).

The analysis results for objects of type SSCP and Covariance do not differ, since these objects only differ in a global scaling factor. A Correlation object has to be used if the variances of individual variates differ much, or if the units of measurement of the individual variates differ.

The result of a principal component analysis on such objects will be a new object of type $\ensuremath{\mathsf{PCA}}$

Multivariate statistics – Principal Component Analysis (PCA)

Objectives of principal component analysis

To discover or to reduce the dimensionality of the data set.

To identify new meaningful underlying variables.

Multivariate statistics - Principal Component Analysis (PCA) How to start

We assume that the multi-dimensional data have been collected in a table. If the variances of the individual columns differ much or the measurement units of the columns differ then you should first standardize the data.

Performing a principal component analysis on a standardized data matrix has the same effect as performing the analysis on the correlation matrix (the covariance matrix from standardized data is equal to the correlation matrix of these data).

Calculate Eigenvectors and Eigenvalues

We can now make a plot of the eigenvalues to get an indication of the importance of each eigenvalue. The exact contribution of each eigenvalue (or a range of eigenvalues) to the "explained variance" can also be queried: You might also check for the equality of a number of eigenvalues.



Multivariate statistics – Principal Component Analysis (PCA)

Determining the number of components

There are two methods to help you to choose the number of components. Both methods are based on relations between the eigenvalues.

Plot the eigenvalues: If the points on the graph tend to level out (show an "elbow"), these eigenvalues are usually close enough to zero that they can be ignored.

Limit variance accounted for and get associated number of components

Multivariate statistics – Principal Component Analysis (PCA)

Getting the principal components

Principal components are obtained by projecting the multivariate datavectors on the space spanned by the eigenvectors. This can be done in two ways:

- 1. Directly from the table without first forming a PCA object: You can then draw the Configuration or display its numbers (Gower, 1966).
- 2. Standard way: project the data table onto the PCA's eigenspace.

Multivariate statistics - Principal Component Analysis (PCA)

Gower (1966) has shown that, given $d_{a}(i, k = 1, ..., n_{n_{cl}})$ whice represents a matrix of pairwise Euclidean distances between $n_{a_{l}}$ reference objects, the co-ordinates $x_{i}, (j = 1, ..., n_{b})$ of the corresponding points P in a Euclidean subspace of x_{a} dimensions can be found by the following procedure. (1) Define $a_{a} = -iq_{a}^{2}$, $i, k = 1, ..., n_{acl}$. (2) Transform a_{a} to $b_{a} = a_{a} - \langle a_{a} \rangle_{-} \langle a_{a} \rangle_{+} \langle a_{a} \rangle_{a}$, where $\langle \cdots \rangle_{j_{a}}$, is the mean over all specified indices and $i, k = 1, ..., n_{p_{cl}}$. (3) Find all $n_{p} \le n_{ed} - 1$ positive eigenvalues $\langle \lambda_{j} \rangle$ of the matrix with elements a_{a} and all corresponding column eigenvectors $[x_{ij}, i = 1, ..., n_{p_{cl}}$. (4) Scale these $n_{p_{cl}}$ ohmun vectors which compose matrix X so that $X \setminus X - A$ where matrix $X = diag(\lambda_{1}, ..., \lambda_{q_{p_{cl}}})$.

 λ_{n_0}). The eigenvalues determined at step 3 represent the cariation along a corresponding axis. So, if eigenvalues The eigenvalues determined at step 3 represent the variation along a corresponding axis. So, if eigenvalues together with eigenvectors are sorted in decreasing order, the first two coordinates represent the best planar projec-tion. Negative eigenvalues, resulting from non-Euclidean distances in our applications, usually account for less than 20°_{20} of the variation.

Multivariate statistics – Principal Component Analysis (PCA)

Mathematical background on principal component analysis

The mathematical technique used in PCA is called eigen analysis: we solve for the eigenvalues and eigenvectors of a square symmetric matrix with sums of squares and cross products. The eigenvector associated with the largest eigenvalue has the same direction as the first principal component. The eigenvector associated with the second largest eigenvalue determines the direction of the second principal component. The sum of the eigenvalues equals the trace of the square matrix and the maximum number of eigenvectors equals the number of rows (or columns) of this matrix.









How accurate? Sequence searching: Specificity Selectivity

Sequence alignment: Column score or pairscore (correctness) Alignment score (significance)

General: Z-score

low fast? Time complexity CPU time

Statistics

Random variables described by probability density of normal distribution

$$p(x) = \frac{1}{\sqrt{2\pi}\,\sigma}\,e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$

The distribution belongs to the class of Gaussian distributions

 $p(x) = e^{-((a-b)^2/c^2)}$

In benchmarking we have two random distributions: The background distribution of all true negative scores and the target distribution of all true positive scores.









Benchmark Curve

What is a benchmark curve?

Every program that computes something can give true and false answers. In a benchmark curve one plots, for example, the sum of true answers in y-direction and the sum of false answers in the x-direction.

The terms 'true' and 'false' might be replaced by 'sensitivity' and 'specificity', or 'coverage' and 'error per query'; however, the basic idea is always the same.

An easy way to visualise a benchmark curve is a walk based on answers to a quiz: if the answer is correct go straight, if it is incorrect go right. The curve you walk is a benchmark curve: if you go mostly straight then you know a lot, if you go mostly right you know a little.

Receiver Operator Characteristics

Plot p(TP) [= sensitivity] over p(FP) [= 1 - specificity]

ROC curves are used to evaluate the results of a prediction and was first employed in the study of discriminator systems for the detection of radio signals in the presence of noise in the 1940s. In the 1960s they began to be used in psychophysics, to assess human (and occasionally animal) detection of weak signals. They also proved to be useful for the evaluation of machine learning results, such as the evaluation of Internet search engines. They are also used extensively in epidemiology and medical research.

Sometimes, the ROC is used to generate a summary statistic. Two common versions are:

the intercept of the ROC curve with the line at 90 degrees to the no-discrimination line

the area between the ROC curve and the no-discrimination line

However, any attempt to summarize the ROC curve into a single number loses information about the pattern of tradeoffs of the particular discriminator algorithm.









Important Aspects of Benchmarking

What is the 'standard of truth'?

Is the standard of truth independent of your own application?

For example, you can use a sequence alignment database derived from structural alignment as standard of truth for benchmarking a sequence alignment program, but you need to be careful when using a detabase derived from sequence alignme as standard of truth.

Is the training set random and independent of the test set?

When developing an application, you need to parametrise and test it on a standar of truth. The data set used for parametrisation is called 'training set' and the data set used for testing is called 'test set'.

 It is absolutely mandatory that training and test data sets are independent (no overlap).

2. The data for the training set have to be drawn randomly from the standard of truth. If you pre-select specific data, the application will be biased towards these data and the benchmark as well.















